DiP-GO: A <u>Di</u>ffusion <u>P</u>runer via Few-step <u>G</u>radient <u>O</u>ptimization

Anonymous Author(s) Affiliation Address email

Abstract

Diffusion models have achieved remarkable progress in the field of image genera-1 tion due to their outstanding capabilities. However, these models require substantial 2 computing resources because of the multi-step denoising process during inference. З While traditional pruning methods have been employed to optimize these models, 4 the retraining process necessitates large-scale training datasets and extensive com-5 putational costs to maintain generalization ability, making it neither convenient 6 nor efficient. Recent studies attempt to utilize the similarity of features across 7 adjacent denoising stages to reduce computational costs through simple and static 8 strategies. However, these strategies cannot fully harness the potential of the similar 9 feature patterns across adjacent timesteps. In this work, we propose a novel pruning 10 method that derives an efficient diffusion model via a more intelligent and differ-11 entiable pruner. At the core of our approach is casting the model pruning process 12 into a SubNet search process. Specifically, we first introduce a SuperNet based 13 on standard diffusion via adding some backup connections built upon the similar 14 features. We then construct a plugin pruner network and design optimization losses 15 to identify redundant computation. Finally, our method can identify an optimal 16 SubNet through few-step gradient optimization and a simple post-processing pro-17 cedure. We conduct extensive experiments on various diffusion models including 18 Stable Diffusion series and DiTs. Our DiP-GO approach achieves $4.4 \times$ speedup 19 for SD-1.5 without any loss of accuracy, significantly outperforming the previous 20 state-of-the-art methods. 21

22 1 Introduction

Diffusion models have undergone significant advancements over the past years due to the outstanding 23 capabilities of Diffusion Probabilistic Models (DPMs) [1]. DPMs typically consist of two processes: 24 the noise diffusion process and the reverse denoising process. Given their remarkable superiority 25 in content generation, diffusion models have made significant progress in various fields of general 26 27 image generation, including text-to-image generation [2, 3], layout-to-image generation [4, 5], image editing [6, 7], and image personalization [8, 9]. Furthermore, diffusion models have contributed 28 to advancements in autonomous driving, ranging from driving dataset generation [10, 11, 12] to 29 perception model enhancement [13, 14] through diffusion strategies. However, DPMs often incur 30 considerable computational overhead during both the training and inference phases. The high cost of 31 inference, due to the multi-step denoising computation during the sampling process, can significantly 32 impact their practical application. Many efforts [15, 16, 17] have been made to improve the efficiency 33 of diffusion models, which can be broadly divided into two types of optimization: inference sampling 34 optimization and model structural optimization. 35

Sampling optimization methods reduce the number of sampling steps for generation without compro-36 mising image quality. For instance, DDIM [15] reduces these steps by exploring a non-Markovian 37 process without requiring model retraining. LCM [18, 19] enable image generation in fewer steps 38 with retraining requirements. Structural optimization methods [16, 17, 20, 21] aim to reduce com-39 putational overhead through efficient model design and model pruning. These methods require 40 retraining the diffusion model, which entails significant computational overhead and large-scale 41 datasets, making them neither convenient nor efficient. DeepCache [22] proposes a novel training-free 42 paradigm based on the U-Net architecture in diffusion models, caching and retrieving features across 43 adjacent denoising stages to reduce redundant computation costs. However, DeepCache only reuses 44 the output feature of a U-Net block in a denoising step via a simple and static strategy. We believe 45 many intermediate features remain untapped, and the simple static strategy cannot fully exploit the 46 potential of similar feature patterns across adjacent timesteps during inference, as observed in recent 47 studies [15, 18, 22]. 48

To address these challenges, we introduce Diffusion Pruning via Few-step Gradient Optimization 49 (DiP-GO), a method designed to achieve efficient model pruning with enhanced dynamism and 50 intelligence. Our approach rethinks the diffusion model during inference by proposing a SuperNet 51 based on standard diffusion via adding some backup connections built upon the similar features, 52 conceptualizing the inference process as a specific SubNet derived from our proposed SuperNet. 53 We reformulate the diffusion model pruning into a SubNet search process. By addressing the out-54 of-memory issue inherent in the backward process during expanded denoising timesteps using the 55 gradient checkpoint [23] method, we introduce a plugin pruner that discovers an optimal SubNet 56 57 surpassing existing methods through carefully designed optimization losses. Extensive experiments validate the effectiveness of our approach, demonstrating a $4.4 \times$ speedup on Stable Diffusion 1.5. 58 Moreover, our method efficiently prunes the DiT model [3] without requiring retraining the diffusion 59 model, achieving significant inference speedup. Our contribution can be summarized as follows: 60 (1) We define a SuperNet based on standard diffusion and show how to obtain a SubNet. This 61 transforms the diffusion optimization problem into an efficient SubNet search process without the 62 need for retraining pretraind diffusion models. (2) We design a plugin pruner tailored specifically for 63 diffusion models. This pruner optimizes pruning constraints while maximizing synthesis capability. 64 Additionally, we introduce a post-processing method for the pruner to ensure that the SubNet meets 65 specified pruning requirements. (3) We conduct extensive pruning experiments across various 66 diffusion models, including Stable Diffusion 1.5, Stable Diffusion 2.1, Stable Diffusion XL, and DiT. 67 Extensive experiments demonstrate the superiority of our method, achieving a notable $4.4 \times$ speedup 68 during inference on Stable Diffusion 1.5 without the need for retraining the diffusion model. 69

70 2 Related Work

71 2.1 Efficient Diffusion Models

The diffusion models, celebrated for their iterative denoising process during inference, play a pivotal role in content generation but are often hindered by time-consuming operations. To mitigate this challenge, extensive research has focused on accelerating diffusion models. Acceleration efforts typically approach the problem from two primary perspectives:

Efficient Sampling Methods. Recent works focus on reducing the number of denoising steps 76 required for content generation. DDIM [15] achieves this by exploring a non-Markovian process 77 related to neural ODEs. Fast high-order solvers [24, 25] for diffusion ordinary differential equations 78 also enhance sampling speed. LCMs [18, 19] treat the reverse diffusion process as an augmented 79 probability flow ODE (PF-ODE) problem, inspired by Consistency Models (CMs) [26], enabling 80 generation in fewer steps. PNDM [27] emphasizes efficient sampling without retraining diffusion 81 model. Additionally, ADD [28] combines adversarial training and score distillation to transform 82 pretrained diffusion models into high-fidelity image generators using only single sampling steps. 83

Efficient Structural Methods. Other efforts concentrate on reducing the computational overhead
 associated with each denoising step. Previous methods [16, 17, 22] have typically conducted extensive
 empirical studies to identify and remove non-critical layers from U-Net architectures to achieve faster
 networks. BK-SDM [16] customizes three efficient U-Nets by strategically removing residual and
 attention blocks. Derived from BK-SDM, KOALA [17] develops two efficient U-Nets of varying sizes
 tailored for SD-XL applications. Diff-pruning [20] employs Taylor expansion over pruned timesteps

to pinpoint essential layer weights, optimizing model efficiency without sacrificing performance.
 DeepCache [22] enhances inference efficiency by reusing predictions from blocks in previous
 timesteps within the U-Net architecture. LAPTOP-Diff [21] tackles optimization problems with a
 one-shot pruning approach, incorporating normalized feature distillation to streamline retraining
 processes. T-GATE [29] not only reduces computation overhead but also marginally lowers FID
 scores by omitting text conditions during fidelity-improvement stages.

⁹⁶ In addition to the two primary acceleration methods, other strategies such as distillation [28, 30, 31],

early stopping [32], and quantization [33] are commonly employed to enhance performance and

98 efficiency. However, most of these strategies necessitate retraining pretrained models. Our method

falls under the category of efficient structural methods by focusing on reducing inference time at each timestep. Importantly, these efficiency gains are achieved without retraining the diffusion model.

101 2.2 Model Optimization

Network Pruning. The taxonomy of pruning methodologies typically divides into two main cate-102 gories: unstructured pruning methods [34, 35, 36] and structural pruning methods [37, 38, 39, 40]. 103 Unstructured pruning methods involve masking parameters without structural constraints by zeroing 104 them out, often requiring specialized software or hardware accelerators. In contrast, structured 105 pruning methods generally remove regular parameters or substructures from networks. Recent works 106 have been interested in accelerating transformers. Dynamic skipping blocks, which involve selectively 107 removing layers while maintaining the overall structure, have emerged as a paradigm for transformer 108 compression [41, 42, 43, 44]. However, applying structural pruning techniques to diffusion modeling 109 poses unique challenges that necessitate reevaluating conventional pruning methods. 110

111 3 Methodology

In this study, we introduce the Diffusion Pruner via Few-step Gradient Optimization (DiP-GO), 112 which utilizes a neural network to predict whether to skip or keep each computational block during 113 inference. Our primary objective is to identify the optimal subset of computational blocks that 114 facilitate denoising with minimal computational overhead. As illustrated in Figure 2, our method 115 comprises three main components: a neural network pruner, optimization losses, and a post-process 116 algorithm to derive the pruned model based on the predictions of pruner. The neural network pruner 117 is designed with learnable queries inspired by DETR [45] to predict the state of each block. Our 118 119 proposed optimization losses include sparsity and consistency constraints for generation quality, guiding the pruner to accurately assess the importance of each block. In this Section, we first revisit 120 the framework of diffusion models in Section 3.1, emphasizing their potential for exploring pruned 121 networks. In Section 3.2, we introduce a SuperNet based on diffusion models and demonstrate how 122 to derive a SubNet or pruned network from it for inference acceleration, highlighting the challenges 123 in achieving an optimal SubNet. Section 3.3 details our method, including the neural network 124 pruner, optimization losses, and post-process algorithm for obtaining a SubNet that meets pruning 125 requirements. Finally, we provide insights into the training and inference processes of our method. 126

127 3.1 Preliminary

We begin with a brief introduction to diffusion models. Diffusion models are structured to learn a series of sequential state transitions with the goal of iteratively refining random noise sampled from a known prior distribution towards a target distribution x_0 that matches the data distribution. During the forward diffusion process, the transition from x_{t-1} to x_t is initially determined by a forward transition function, which can be described as follows:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t \mathbf{I})$$
(1)

where the hyperparameter $\{\beta_t \in (0, 1)\}_{t=1}^T$ increases with each successive time step t.

To generate samples from a learned diffusion model, it involves a series of reverse state transitions from $x_T \rightarrow \cdots \rightarrow x_0$ to denoise random noise $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ into the clean data point x_0 . At each timestep, the denoised output x_{t-1} is predicted by approximating the noise prediction network, which is conditioned on the time embedding t and the previous data point x_t :

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \frac{1}{\sqrt{a_t}}(x_t - \frac{\beta_t}{\sqrt{1 - \overline{a}_t}} z_{\theta}(x_t, t)), \beta \mathbf{I})$$
(2)

where $\beta_t = 1 - \alpha_t$, $\overline{a}_t = \prod_{i=1}^T \alpha_i$, and $z_{\theta}(x_t, t)$ are the parameterized deep neural networks. With the reverse Markov chain, we can iteratively sample from the learnable transition kernel $x_{t-1} \sim p_{\theta}(x_{t-1}|x_t)$ until t = 1.

Diffusion modes typically require multi-step conditional sampling to gradually obtain the target 141 sample point x_0 . However, recent studies [15, 18, 22] have highlighted that multi-step inference pro-142 cesses involve substantial redundant feature computations, particularly in noise prediction networks 143 like UNet and Transformer. For example, in Stable Diffusion 1.4 models with 25 steps, Multiply-144 Accumulate Operations (MACs) of UNet can comprise up to 87.2% of the total computational load 145 [16]. This underscores significant potential for accelerating inference by effectively eliminating 146 these redundancies. In this work, we propose accelerating the diffusion model by integrating a 147 differentiable pruning network designed to identify and remove these redundant computations. 148

149 **3.2** SuperNet and SubNet of Diffusion Model

163

Our goal is to identify and remove unimportant blocks during inference to accelerate the process. To achieve this, we introduce a SuperNet based on the diffusion model. This SuperNet is designed to facilitate block removal while ensuring the pruned model maintains inference capability through additional connections. Our approach effectively eliminates unimportant blocks during inference, essentially deriving a SubNet from the SuperNet by skipping these unnecessary components. Thus, the pruning process can be conceptualized as a SubNet search within the SuperNet framework.

How to Construct a SuperNet. Recent studies [15, 18, 22] have observed that diffusion models often exhibit similar feature patterns across adjacent timesteps during inference. Building on this insight, we enhance the standard diffusion model's inference phase by introducing additional connections from the current timestep to the previous one. These connections serve as backups for blocks that may be removed, ensuring each block retains valid inputs even if its dependent blocks are eliminated for acceleration. Specifically, for all inputs of each block across all timesteps except the initial step during inference, we establish a backup input connection to the corresponding block in the previous timestep, as illustrated in Figure 1.



Figure 1: Overview of the SuperNet and SubNet. Standard diffusion models execute the full inference path step by step. In our framework, we propose a SuperNet based on the original flow and integrate backup connections to facilitate block removal. This allows the partial inference SubNet to efficiently eliminate redundant computational costs.

How to Obtain a SubNet. To construct the SuperNet for the standard diffusion model, we introduce additional connections that ensure a valid SubNet selects either the original input connection or the backup input connection, but not both simultaneously. This design principle mandates that if a dependent block is pruned, its original input connection is also eliminated to reflect the block's removal. Conversely, if the dependent block is retained, the backup input connection is removed to maintain efficient inference, as depicted in Figure 1.

We draw inspiration from the Lottery Ticket Hypothesis (LTH) [46], which posits the existence of a sub-network capable of achieving comparable performance to the original over-parameterized network for a given task, but with fewer unnecessary weights. Moreover, prior work [22] has explored manually removing redundant computations by caching features across adjacent steps. Thus, our



Figure 2: Overview of our diffusion pruner. a) DiP-GO employs a pruner network to learn the importance scores of blocks in the diffusion sampling process. It takes $N \times T$ queries as input and passes them through stacked self-attention (SA) and fully connected (FC) layers to capture the structural information in existing diffusion models. The network predicts the partial inference paths based on the $N \times T$ importance scores and is optimized by consistent and sparse loss. b) Once trained, the pruner network is discarded. We can infer the optimal partial inference path with expected computational costs via post-processing based on the predicted importance scores.

approach seeks to identify an optimal SubNet from the SuperNet, maximizing diffusion model acceleration while minimizing any loss in generation quality.

Hard to Obtain an Optimal SubNet. The challenge of obtaining an optimal SubNet is compounded 176 by the large number of blocks expanded during inference. In a diffusion pipeline with $N \times T$ blocks 177 (where N is the number of blocks per timestep and T is the number of timesteps), each block's 178 decision to be kept or removed results in $2^{N \times T}$ possible configurations. For instance, a 50-step 179 PLMS setup [27], considering 9 blocks in the U-Net, yields 2^{450} choices (> 10^{135}). Traditional search 180 methods like random search and genetic algorithms [47] often struggle in such vast search spaces. 181 Gradient-based optimization offers a promising approach to tackle this challenge. However, there are 182 significant hurdles to overcome. First, effectively modeling discrete block states (kept or removed) 183 with parametric methods poses difficulties. Second, training the entire model, comprising both the 184 parametric model and the expanded diffusion model with denoising timesteps, risks encountering 185 out-of-memory (OOM) issues. 186

187 3.3 Our DiP-GO Approach

In this study, we introduce a diffusion pruner network designed to predict importance scores for all 188 blocks during reverse sampling as depicted in Figure 2. To optimize the pruner network effectively, 189 we employ two key optimization losses: consistency and sparsity losses, leveraging few-step gradient 190 optimization. Addressing the OOM issue inherent in such computations, we implement gradient 191 checkpointing and half-precision floating-point representation techniques, enabling efficient search 192 processes on a single GPU. Once the pruner network trained, we extract predicted importance scores 193 for all blocks. Subsequently, we devise a post-processing algorithm to utilize these scores, generating 194 pruned SubNets of diffusion models that satisfy specific pruning criteria. 195

Pruner Network. Our pruner network comprises three main components: $N \times T$ learnable queries, a query encoder, and a prediction head. We design the learnable queries to match the number of all blocks during inference. These queries are optimized with sparsity and consistency loss constraints to learn the contextual information necessary for predicting the importance score of each block. For the query encoder, we provide two options: a simple version with several stacked linear layers, and a more complex version with several stacked self-attention layers to facilitate interaction among the learnable queries. Our experiments demonstrate that both versions can effectively obtain optimal SubNets in various diffusion models under different pruning requirements. The prediction head consists of $N \times T$ simple branches, each containing two stacked linear layers followed by a softmax operation. The final linear layer has a dimension of 2, and the softmax output represents the importance scores of a block. During training or inference, the query embeddings are transformed into output embeddings via the query encoder. These embeddings are then independently decoded into binary vectors by the multi-layer prediction head, resulting in $N \times T$ importance scores for all blocks.

Optimization Losses. The k-th predicted binary vector of importance score, denoted as s^k , represents 209 the likelihood of its corresponding block being removed or kept in the denoising process. A gate $g \in \{0,1\}^{TN}$ is derived based on s, where $g^k = 0$ or $g^k = 1$ indicate removing or keeping the k-th 210 211 computation block, respectively. Only the blocks that are kept according to q will be calculated in the 212 denoising process. However, directly converting predicted probabilities s into discrete gates g with 213 arg max is non-differentiable. To address this issue, we utilize the Straight-Through (ST) Estimator 214 [48] to approximate the real gradient $\nabla_{\theta} g$ with the gradient of the soft prediction $\nabla_{\theta} s$. To encourage 215 both high-fidelity predictions and minimal computation block usage, we design our training objective 216 function as a combination of consistent loss \mathcal{L}_c and sparse loss \mathcal{L}_s , formulated as follows: 217

$$\mathcal{L}(\boldsymbol{x}_{T};\boldsymbol{\theta},W) = \mathcal{L}_{c} + \alpha_{s}\mathcal{L}_{s} = \begin{cases} f(\boldsymbol{x}_{0}^{p},\boldsymbol{x}_{0}^{gt}) + \frac{\alpha_{s}}{NT} \sum_{k}^{NT} \gamma^{k} \boldsymbol{g}^{k} & \text{if } sparsity < \tau \\ f(\boldsymbol{x}_{0}^{p},\boldsymbol{x}_{0}^{gt}) & \text{if } sparsity \geq \tau \end{cases}$$
(3)

Here, α_s represents a hyperparameter used to balance the consistent and sparse losses. θ and W denote 218 the pruner network and pretrained diffusion model, respectively. $f(\cdot)$ denotes a distance function that evaluates the consistency between the generated clean data point \mathbf{x}_0^p from partial inference of the 219 220 pruned SubNet and the \mathbf{x}_0^{gt} from full inference. This function can be any distance measure, and in 221 this work, we utilize a negative SSIM loss [49]. The sparse loss encourages minimal computational 222 usage and is weighted by the computational flops proportion γ^k of the k-th block, thereby imposing a 223 greater penalty on heavier blocks. The calculation of γ takes into account the cascading relationships 224 between blocks. Specifically, when a block is pruned, the associated dependent blocks will also 225 pruned. Therefore, the flops reduction from pruning a block includes the block itself and its dependent 226 blocks. We denote the flops reduction ratio after pruning the k-th block as γ^k . The sparse loss is 227 only introduced when the sparsity (pruning ratio) is below a certain threshold τ . This compound loss 228 controls the trade-off between efficiency (block usage) and accuracy (generation quality). 229

230 **Post-Processing Algorithm.** After training the pruner network, our diffusion pruner is able to predict which computation blocks during inference contribute less to generation quality based on the 231 importance scores for all the blocks. As the importance scores are continuous values in inference 232 phase, they can not be utilized directly to identify which blocks should be removed to meet given 233 pruning requirements. Therefore, we present a post-process algorithm to obtain an appropriate 234 threshold for these importance score to meet the pruning requirements as shown in Algorithm 1 235 in Appendix B. Considering the required pruning sparsity, we use bisection lookup to select the 236 appropriate threshold value to identify which blocks should be removed to meet the pruning ratio. 237 Specifically, the blocks whose important scores below the threshold should be removed and the kept 238 blocks should update their input connections as mentioned in Section 3.2 to maintain the pruned 239 model inference. Thus a pruned model met the pruning ratio has been obtained. 240

Training and Inference Details. In the training phase, the prompt inputs are fed into the diffusion 241 model to obtain two kinds of outputs, one is generated by the baseline diffusion model and the other 242 is generated by the pruned model obtained via the current predictions of the pruner network. Then 243 our proposed losses are utilized to optimize the pruner network to enable distinguishing the less 244 important blocks. In the pruner's network, we initialize the weight of the last linear layer's output 245 channel to 0 and its bias to 1. This setup ensures that at the beginning of training, the consistency loss 246 is 0 and the sparsity loss is 1, facilitating smooth training. As training progresses, the sparsity loss 247 gradually decreases while the number of pruned blocks increases, causing the consistency loss to rise. 248 To maintain network fidelity after pruning, we switch to training only with the consistency loss once 249 the sparsity loss reaches 0.2, continuing until training is complete. Once the pruner is well trained, 250 we can obtain pruned models to meet the pruning requirements via our post-process algorithm. 251

252 4 Experiments

253 4.1 Experimental Setup

Pre-trained Model and Datasets. We select four official pretrained Diffusion Models (i.e., SD-1.5
[2], SD-2.1 [2], SD-XL [50] and DiT [3]) to evaluate our approach. The SD series models are
constructed on the U-Net [51] and the DiT is constructed on the transformer [52]. We utilize a subset
of the DiffusionDB [53] dataset comprising 1000 samples to train our pruner network, utilizing only
textual prompts. Following previous works [29, 22], we evaluate the DiP-GO on three public datasets,
i.e., PartiPrompts [54], MS-COCO 2017 [55] and ImageNet [56].

Evaluation Metrics. We employ the Fréchet Inception Distance (FID) [57] metrics to assess the quality of images created by the generative models. FID quantifies the dissimilarity between the Gaussian distributions of synthetic and real images. A lower FID score indicates a closer resemblance to real images in the generative model. Additionally, we utilize the CLIP Score [58] (ViT-g/14) to evaluate the relational compatibility between images and text.

Implementation Details. For Stable Diffusion models, we utilize the SGD optimizer with a cosine 265 learning schedule for 1000 steps of training. The batch size, learning rate, and weight decay are set to 266 1, 0.1, and 1×10^{-4} , respectively. The hyperparameters α_s , τ , and the query embedding dimension 267 d, along with the encoder layer number L, are set to 1, 0.2, 512, and 1, respectively. For the Diffusion 268 Transformer model, we use the same experimental configuration as for the stable diffusion model, 269 except that the learning rate set to 10. To evaluate the inference efficiency, we evaluate the Multiply 270 Accumulate Calculation (MACs), Parameters (Params), and Speedup for all models with batchsize of 271 1 in the PyTorch 2.1 environment on the AMD MI250 platform. Besides, we report MACs in those 272 tables, which refer to the totals MACs for all steps. 273

Table 1: Comparison with PLMS, BK-SDM and DeepCache on SD-1.5. We utilize prompts in PartiPrompt and COCO2017 validation set.

		PartiPrompts			COCO2017		
Method	Pruning Type	MACs↓	Speedup ↑	CLIP Score ↑	MACs↓	Speedup \uparrow	CLIP Score \uparrow
PLMS - 50 steps	Baseline	16.94T	$1.00 \times$	29.51	16.94T	$1.00 \times$	30.30
BK-SDM - Base	Structured	11.19T	$1.49 \times$	28.88	11.19T	$1.45 \times$	29.47
PLMS - 25 steps	Fast Sampler	8.47T	$2.04 \times$	29.33	8.47T	$1.91 \times$	29.99
PLMS - Skip - Interval=2	Structured	8.47T	$2.04 \times$	19.74	8.47T	$1.91 \times$	16.78
DeepCache	Structured	6.52T	$2.15 \times$	29.46	6.52T	$2.11 \times$	30.23
Ours (w/ Pruned-0.80)	Structured	3.38T	4.43 ×	29.51	3.38T	4.40 imes	30.29
BK-SDM - Small	Structured	10.88T	1.75×	27.94	10.88T	$1.68 \times$	27.96
PLMS - 15 steps	Fast Sampler	5.08T	$2.89 \times$	28.58	5.08T	$2.59 \times$	29.39
Ours (w/ Pruned-0.85)	Structured	2.54T	5.52×	29.07	2.54T	5.46 ×	29.84

Table 2: Comparison of computational complexity, inference speed, CLIP Score and FID on the MS-COCO 2017 validation set on SD-2.1.

Inference Method	MACs↓	Speedup ↑	CLIP Score ↑	FID-5K↓
SD-2.1-50 steps [2]	38.04T	$1.00 \times$	31.55	27.29
SD-2.1-20 steps [2]	15.21T	$2.49 \times$	31.53	27.83
Ours (w/ Pruned-0.7)	11.42T	$3.02 \times$	31.50	25.98
Ours (w/ Pruned-0.8)	7.61T	$3.81 \times$	30.92	27.69

274 4.2 Comparison with State-of-the-Art Methods on Different Base Models

Stable Diffusion on PartiPrompt and COCO2017. We compare our method with the state-of-the-275 art (SOTA) compression methods on Stable Diffusion 1.5 (SD-1.5), and the results are summarized in 276 Table 1. Compared to the SOTA DeepCache [22], our approach demonstrates significant performance 277 improvements, achieving nearly $2 \times$ fewer MACs while maintaining better CLIP Scores. Our method 278 can achieve $4.4 \times$ speedup compared to the baseline model. Furthermore, our method does not require 279 training the diffusion model, which preserves the pre-trained knowledge of the diffusion model. Also, 280 we apply our method on the SD-2.1 model to verify the effectiveness, as shown in Table 2, our 281 method achieves significant acceleration while maintaining generation quality, demonstrating its 282 superiority. 283

Diffusion Transformers on ImageNet. To the best of our knowledge, we are the first to apply pruning to DiT [3] model. Therefore, we have replicated a training-free acceleration method, DeepCache

Method	Pruning Type	MACs↓	FID-50K \downarrow	Speedup ↑
DiT-XL/2-250 steps	-	29.66T	2.27	$1.00 \times$
DiT-XL/2*-250 steps	Baseline	29.66T	2.97	$1.00 \times$
DiT-XL/2*-110 steps	Fast Sampler	13.05T	3.06	$2.13 \times$
DiT-XL/2*-100 steps	Fast Sampler	11.86T	3.17	$2.46 \times$
DeepCache(DiT-XL/2*)-N=2	Structured Pruning	15.88T	3.07	$1.76 \times$
Ours (DiT-XL/2* w/ Pruned-0.6)	Structured Pruning	11.86T	3.01	2.43 imes
DiT-XL/2*-70 steps	Fast Sampler	8.30T	3.35	$3.49 \times$
DeepCache(DiT-XL/2*)-N=5	Structured Pruning	6.77T	3.20	$3.44 \times$
Ours (DiT-XL/2* w/ Pruned-0.75)	Structured Pruning	7.40T	3.14	3.60 imes

Table 3: Comparison of pruning type, computational complexity, FID and inference speed on the ImageNet validation datasets on DiT. * denotes the results reproduced with diffusers [59].

with intervals = 2 and 5, on DiT for comparison. The results in Table 3 show that our method can
speed up the original DiT model by a factor of 2.4 with minimal performance loss, while DeepCache
has a lower speedup ratio when applied to the DiT model. This can be attributed to DeepCache's
overreliance on pre-defined structures, whereas our method can automatically learn the optimal
pruning strategy for the given model, thereby achieving superior performance.

291 4.3 Compatibility with Fast Sampler

We investigate the compatibility of DiP-GO with methods that prioritize reducing sampling steps 292 using faster samplers: DDIM [15], DPM-Solver [25], and LCM [18]. As shown in Table 4, it indicates 293 that our method further improves computational efficiency on existing fast samplers. Specifically, 294 we reduce MACs by a factor of 5 on the SD-1.5 with DDIM sampler and by $3.36 \times$ on the SD-2.1 295 with DPM-Solver. Our method achieves nearly unchanged performance with significant acceleration. 296 Additionally, our method benefits from information redundancy in multi-step optimization processes, 297 showing relatively limited acceleration performance on fewer-step LCM due to its low redundancy in 298 features across adjacent timesteps. 299

Table 4: Comparison with PLMS, SSIM, and LCM samplers. We evaluate the effectiveness of our methods on COCO2017 validation set.

Sampler	Bas	se Model	Ours	
	$\overline{MACs}\downarrow$	CLIP Score ↑	$\overline{MACs}\downarrow$	CLIP Score
DDIM (SD-1.5 w/ 50 steps)	16.94T	30.30	3.38T	30.29
DPM (SD-2.1 w/ 50 steps)	38.04T	31.55	11.42T	31.50
LCM (SD-XL w/ 4 steps)	11.95T	31.92	11.58T	31.30

300 4.4 Ablation Study

Compared with Different Consistent Constraints. We further compare other alternatives explored 301 for our consistent loss designs, we further scrutinize additional options, including L1, L2, SSIM, 302 and L1+SSIM losses, as depicted in Table 5. The results demonstrate that SSIM emerges as the 303 most effective choice, boasting the highest CLIP-Score. In contrast, the L1 loss function often 304 results in image blurring or distortion due to its sensitivity to pixel-level differences, the L2 loss may 305 yield overly smoothed images by penalizing squared differences between pixels. Conversely, the 306 combination of L1+SSIM loss attempts to address these limitations but can complicate the training 307 process and suffer from trade-offs. Therefore, SSIM emerges as the preferred choice in our consistent 308 loss designs, offering superior accuracy and stability while preserving image quality. 309

Table 5: Comparison with different consistent loss types. Here we conduct pruning experiments with 80% sparsity on COCO2017 validation using SD-1.5.

Loss Type	L1	L2	SSIM	L1 + SSIM
CLIP Score↑	29.94	29.71	30.29	29.77

Effect of Gradient Optimization. As the traditional search algorithm can also obtain SubNets from our proposed SuperNet. It is crucial to validate whether traditional search-based algorithms yield

positive effectiveness. We assess two search algorithms: random search and genetic algorithm-based 312 search [47] in Table 6. We have iterated the search 1000 times using the first 500 images of the test set 313 as a calibration dataset. Remarkably, we observe that the search time of traditional search algorithms 314 is significantly longer than the training time of our method due to a large number of evaluations. 315 Moreover, due to the vast search space, traditional search algorithms struggle to achieve satisfactory 316 results. Additionally, traditional search algorithms lack the "once-for-all" characteristic, requiring 317 re-execution when faced with deployment scenarios demanding different computational resources. In 318 contrast, leveraging the parametric pruner network, our method achieves superior performance with 319 reduced running time and is more adaptable to diverse development scenarios. 320

Table 6: Comparison of cost time, computational complexity and CLIP-Score between Random Search and GA search strategies on Stable Diffusion 1.5.

Method	Cost GPU Hours \downarrow	Pruning Ratio	$ $ MACs $\downarrow $	CLIP Score \uparrow
PLMS-50 steps	-	-	16.94T	30.30
Random Search	25	0.80	2.96T	28.73
GA Search	25	0.80	3.34T	29.37
Ours	2.3	0.80	3.38T	30.29
Random Search	24	0.85	2.90T	27.22
GA Search	24	0.85	2.73T	28.61
Ours	2.2	0.85	2.54T	29.84
Random Search	23	0.90	1.94T	24.07
GA Search	23	0.90	2.04T	25.14
Ours	2.2	0.90	1.69T	28.72

Qualitative Analysis of Increased Prune Ratio. In Figure 3, we visualize the generated images as we increase the pruning ratio. With the increase in pruning ratio, the model's inference speed significantly improves, allowing us to achieve up to a fourfold increase in inference speed. However, as the pruning ratio increases, some patterns in the image content deviate from those in the original images. Nevertheless, our main objects in the figures consistently adhere to the textual conditions. Subtle changes in background details typically do not compromise image quality, as quantitatively analyzed in Table 1.



Figure 3: Visualization of generated images. It shows evolving patterns as pruning ratios increase. Despite these changes, main objects in the images remain consistent with the textual conditions.

328 5 Conclusion

This work explores resolving diffusion accelerating tasks by reducing redundant feature calculations across adjacent timesteps. We present a novel diffusion pruning framework and cast the model pruning process as a SubNet search problem. Our approach introduces a plugin pruner network that identifies an optimal SubNet through few-step gradient optimization. Results on a wide range of Stable Diffusion (SD) and DiT series models verify the effectiveness of our method. We achieve a 4.4× speedup on Stable Diffusion 1.5 and effectively prune the DiT model with few step optimizations.

335 **References**

- [1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020.
- [2] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022.
- [3] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings* of the IEEE/CVF International Conference on Computer Vision, pages 4195–4205, 2023.
- [4] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, pages 3836–3847, 2023.
- [5] Chengyou Jia, Minnan Luo, Zhuohang Dang, Guang Dai, Xiaojun Chang, Mengmeng Wang, and Jingdong Wang. Ssmg: Spatial-semantic map guided diffusion model for free-form layout-to-image generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 2480–2488, 2024.
- [6] Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. Anydoor:
 Zero-shot object-level image customization. *arXiv preprint arXiv:2307.09481*, 2023.
- [7] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri,
 and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6007–6017,
 2023.
- [8] 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.
- [9] 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.
- [10] Xiaofan Li, Yifu Zhang, and Xiaoqing Ye. Drivingdiffusion: Layout-guided multi-view driving
 scene video generation with latent diffusion model. *arXiv preprint arXiv:2310.07771*, 2023.
- [11] Ruiyuan Gao, Kai Chen, Enze Xie, Lanqing Hong, Zhenguo Li, Dit-Yan Yeung, and Qiang
 Xu. Magicdrive: Street view generation with diverse 3d geometry control. *arXiv preprint arXiv:2310.02601*, 2023.
- [12] Kai Chen, Enze Xie, Zhe Chen, Lanqing Hong, Zhenguo Li, and Dit-Yan Yeung. Integrating
 geometric control into text-to-image diffusion models for high-quality detection data generation
 via text prompt. *arXiv: 2306.04607*, 2023.
- [13] Jiacheng Chen, Ruizhi Deng, and Yasutaka Furukawa. Polydiffuse: Polygonal shape reconstruction via guided set diffusion models. *ArXiv*, abs/2306.01461, 2023.
- [14] Zhendong Wang, Jianmin Bao, Wengang Zhou, Weilun Wang, Hezhen Hu, Hong Chen, and
 Houqiang Li. Dire for diffusion-generated image detection. *arXiv preprint arXiv:2303.09295*,
 2023.
- Ijiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [16] Bo-Kyeong Kim, Hyoung-Kyu Song, Thibault Castells, and Shinkook Choi. Bk-sdm: Architecturally compressed stable diffusion for efficient text-to-image generation. *ICML Workshop on Efficient Systems for Foundation Models (ES-FoMo)*, 2023.
- [17] Youngwan Lee, Kwanyong Park, Yoohrim Cho, Yong Ju Lee, and Sung Ju Hwang. Koala:
 Self-attention matters in knowledge distillation of latent diffusion models for memory-efficient
 and fast image synthesis. *arXiv preprint arXiv:2312.04005*, 2023.

- [18] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models:
 Synthesizing high-resolution images with few-step inference, 2023.
- [19] Simian Luo, Yiqin Tan, Suraj Patil, Daniel Gu, Patrick von Platen, Apolinário Passos, Longbo
 Huang, Jian Li, and Hang Zhao. Lcm-lora: A universal stable-diffusion acceleration module,
 2023.
- [20] Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. In
 Advances in Neural Information Processing Systems, 2023.
- [21] Dingkun Zhang, Sijia Li, Chen Chen, Qingsong Xie, and Haonan Lu. Laptop-diff: Layer
 pruning and normalized distillation for compressing diffusion models, 2024.
- [22] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Deepcache: Accelerating diffusion models for
 free. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [23] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear
 memory cost. *arXiv preprint arXiv:1604.06174*, 2016.
- [24] 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.
- [25] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++:
 Fast solver for guided sampling of diffusion probabilistic models, 2023.
- Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023.
- [27] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models
 on manifolds. *arXiv preprint arXiv:2202.09778*, 2022.
- 404 [28] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion
 405 distillation. *arXiv preprint arXiv:2311.17042*, 2023.
- [29] Wentian Zhang, Haozhe Liu, Jinheng Xie, Francesco Faccio, Mike Zheng Shou, and Jürgen
 Schmidhuber. Cross-attention makes inference cumbersome in text-to-image diffusion models.
 arXiv preprint arXiv:2404.02747, 2024.
- [30] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models.
 arXiv preprint arXiv:2202.00512, 2022.
- [31] David Berthelot, Arnaud Autef, Jierui Lin, Dian Ang Yap, Shuangfei Zhai, Siyuan Hu, Daniel
 Zheng, Walter Talbott, and Eric Gu. Tract: Denoising diffusion models with transitive closure
 time-distillation. *arXiv preprint arXiv:2303.04248*, 2023.
- [32] Zhaoyang Lyu, Xudong Xu, Ceyuan Yang, Dahua Lin, and Bo Dai. Accelerating diffusion
 models via early stop of the diffusion process. *arXiv preprint arXiv:2205.12524*, 2022.
- [33] Yefei He, Luping Liu, Jing Liu, Weijia Wu, Hong Zhou, and Bohan Zhuang. Ptqd: Accurate
 post-training quantization for diffusion models, 2023.
- [34] Sejun Park, Jaeho Lee, Sangwoo Mo, and Jinwoo Shin. Lookahead: A far-sighted alternative of
 magnitude-based pruning. *arXiv preprint arXiv:2002.04809*, 2020.
- [35] Xin Dong, Shangyu Chen, and Sinno Pan. Learning to prune deep neural networks via layer-wise
 optimal brain surgeon. *Advances in neural information processing systems*, 30, 2017.
- [36] Victor Sanh, Thomas Wolf, and Alexander Rush. Movement pruning: Adaptive sparsity by
 fine-tuning. Advances in neural information processing systems, 33:20378–20389, 2020.
- [37] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for
 efficient convnets, 2017.

- [38] Sara Elkerdawy, Mostafa Elhoushi, Abhineet Singh, Hong Zhang, and Nilanjan Ray. To filter
 prune, or to layer prune, that is the question. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [39] Xiaohan Ding, Guiguang Ding, Yuchen Guo, and Jungong Han. Centripetal sgd for pruning
 very deep convolutional networks with complicated structure. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4943–4953, 2019.
- [40] Liyang Liu, Shilong Zhang, Zhanghui Kuang, Aojun Zhou, Jing-Hao Xue, Xinjiang Wang,
 Yimin Chen, Wenming Yang, Qingmin Liao, and Wayne Zhang. Group fisher pruning for
 practical network compression. In *International Conference on Machine Learning*, pages
 7021–7032. PMLR, 2021.
- [41] Ji Liu, Dehua Tang, Yuanxian Huang, Li Zhang, Xiaocheng Zeng, Dong Li, Mingjie Lu,
 Jinzhang Peng, Yu Wang, Fan Jiang, Lu Tian, and Ashish Sirasao. Updp: A unified progressive
 depth pruner for cnn and vision transformer, 2024.
- [42] Minjia Zhang and Yuxiong He. Accelerating training of transformer-based language models
 with progressive layer dropping. *Advances in Neural Information Processing Systems*, 33:14011–
 14023, 2020.
- Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure
 attention loses rank doubly exponentially with depth. In *International Conference on Machine Learning*, pages 2793–2803. PMLR, 2021.
- [44] Lu Hou, Zhiqi Huang, Lifeng Shang, Xin Jiang, Xiao Chen, and Qun Liu. Dynabert: Dynamic
 bert with adaptive width and depth. *Advances in Neural Information Processing Systems*,
 33:9782–9793, 2020.
- [45] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and
 Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.
- [46] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable
 neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- ⁴⁵³ [47] John H Holland. Genetic algorithms. *Scientific american*, 267(1):66–73, 1992.
- [48] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax.
 arXiv preprint arXiv:1611.01144, 2016.
- [49] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:
 from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [50] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- [51] 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.
- [52] 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.
- [53] Zijie J. Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and
 Duen Horng Chau. DiffusionDB: A large-scale prompt gallery dataset for text-to-image
 generative models. *arXiv:2210.14896 [cs]*, 2022.

- Is4] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay
 Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive
 models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5,
 2022.
- [55] 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.
- [56] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern
 recognition, pages 248–255. Ieee, 2009.
- [57] 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. *NeurIPS*, 30, 2017.
- [58] 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.
- [59] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul,
 Mishig Davaadorj, Dhruv Nair, Sayak Paul, Steven Liu, William Berman, Yiyi Xu, and Thomas
 Wolf. Diffusers: State-of-the-art diffusion models, 2023.
- [60] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang,
 James Kwok, Ping Luo, Huchuan Lu, et al. Pixart-a: Fast training of diffusion transformer for
 photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023.
- [61] Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: Probing the reasoning skills and social
 biases of text-to-image generation models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3043–3054, 2023.
- [62] Ranjita Naik and Besmira Nushi. Social biases through the text-to-image generation lens. In
 Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society, pages 786–808,
 2023.
- [63] Candace Ross, Boris Katz, and Andrei Barbu. Measuring social biases in grounded vision and
 language embeddings. *arXiv preprint arXiv:2002.08911*, 2020.
- ⁵⁰² [64] Siwei Lyu. Deepfake detection: Current challenges and next steps. pages 1–6, 2020.

503 A Memory Optimization Details

Gradient Checkpointing. Due to the multi-step Markovian nature of sampling in diffusion models, 504 updating the entire sampling process using gradient accumulation incurs significant memory costs, 505 even with a batch size of 1. To mitigate this issue, we employ Gradient Checkpointing and half-506 precision floating-point training to reduce memory consumption. The core idea behind Gradient 507 Checkpointing is to selectively preserve a portion of activation values during forward propagation, 508 discarding the rest. During backpropagation, the gradients of the discarded activation values are 509 computed using the saved gradients of the preserved nodes, effectively reducing memory usage. 510 Additionally, we use gradient accumulation, wherein gradients computed over multiple iterations 511 are accumulated and then backpropagated in a single batch for parameter updates, thus allowing for 512 larger batch sizes under limited memory usage. 513

514 **B** Pseudo Code

⁵¹⁵ Here, we show the details of our proposed post-process algorithm via pseudo code as followings.

Algorithm 1 Diffusion Pruner

Input: A pretrained diffusion model M, importance scores S, a pruning ratio p**Output:** The pruned diffusion model M^*

```
1: left \leftarrow 0.0
 2: right \leftarrow 1.0
 3: while True do
         current \leftarrow (\text{left} + \text{right})/2
 4:
 5:
         S^* \leftarrow S
         for t in [0, 1, 2, ..., T] do
 6:
              for each block score s in S^* do
 7:
 8:
                  if s < current then
 9:
                       s \leftarrow 0
                  end if
10:
              end for
11:
12:
         end for
         update_scores_of_blocks (S^*) // remove dependent blocks to set them zeros.
13:
14:
         p^*, M^* \leftarrow \text{prune\_diffusion\_model}(S^*, M) // \text{ obtain the pruned ratio and the pruned model.}
         if abs(p^* - p) < 0.0125 then
15:
              break
16:
17:
         else if p^* < p then
              left ← current
18:
19:
         else
20:
              right \leftarrow current
21:
         end if
22:
         \mathcal{T} \leftarrow \mathcal{T}/2
23: end while
24: return M
```

516 C Additional Experiment Results



Figure 4: Visualization of DiT model generated images: samples using DDIM-250 steps (uplink) and pruned 60% MACs (downlink). The speedup ratio here is $2.4 \times$.

⁵¹⁷ We provide the original unpruned DiT model and a version pruned by 0.6 ratio to generate comparison ⁵¹⁸ images in Figure 4. It can be observed that the plots generated by the pruned model are almost

identical to those produced by the original model. Although there are slight differences in details,

such as the appearance of the dog's eyes, these do not significantly affect the overall image quality.

521 D Limitations

A limitation of our method arises from its training process of the pruner network. Our method 522 necessitates tuning an additional pruner network for the pre-trained diffusion model. This may entail 523 users investing additional time when adapting our method to specific diffusion models. For example, 524 we train DiP-GO for SD-1.5 on a single AMD Instinct MI250 GPU for ~ 2.5 hours. However, we 525 note that the introduced time is small compared to training a lightweight diffusion model. Besides, 526 same as existing work, our method struggles to maintain performance with extremely high pruning 527 ratios, presenting a challenge for deploying diffusion models in scenarios with severely limited 528 computational resources. 529

530 E Social Impact

Generative models have demonstrated promising results in content generation [50, 60, 2]. However, due to the high inference costs, current methods struggle to achieve rapid application and deployment. Our approach introduces an efficient acceleration method for diffusion models, enabling nearly lossless speedup. Moreover, our method does not require retraining of the pretrained models and is compatible with various diffusion models, making it highly generalizable. This makes it suitable for rapid deployment of generative models on mobile and edge devices.

Nevertheless, since generative models are pretrained on large-scale internet datasets, the data they generate may contain inherent social biases and stereotypes [61, 62, 63]. Additionally, there is a risk of misuse, such as in the creation of DeepFakes [64], which could pose significant social harm. While reducing the usage cost, it is crucial to prevent the low-cost generative models from being misused, leading to negative societal impacts. Therefore, it is necessary to establish relevant laws and regulations, create a well-regulated community environment, and provide guidelines to ensure responsible dissemination and use of generative models.

544 NeurIPS Paper Checklist

545 1. Claims

- 546 Question: Do the main claims made in the abstract and introduction accurately reflect the 547 paper's contributions and scope?
- 548 Answer: Yes
- Justification: The abstract provides a concise summary of the main findings and contributions of the paper, while the introduction elaborates on the problem statement and research objectives, thereby clarifying the contributions.

552 Guidelines:

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

580

581

582

583

584

585

586

587

588

589

590

591

592

593

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
 - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: In Limitation Section in Appendix D, we expound upon the limitations of the work conducted and provide a brief discussion thereof.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
 - The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
 - The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
 - If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.
- **3. Theory Assumptions and Proofs**
- Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

597	Answer: [No]
598	Justification: None.
599	Guidelines:
000	• The answer NA means that the paper does not include theoretical results
600	• The answer NA means that the paper does not include theoretical results.
601 602	• All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
603	• All assumptions should be clearly stated or referenced in the statement of any theorems.
604	• The proofs can either appear in the main paper or the supplemental material, but if
605	they appear in the supplemental material, the authors are encouraged to provide a short
000	• Inversaly, any informal proof provided in the core of the paper should be complemented
608	by formal proofs provided in appendix or supplemental material.
609	• Theorems and Lemmas that the proof relies upon should be properly referenced.
610 4.	Experimental Result Reproducibility
611	Question: Does the paper fully disclose all the information needed to reproduce the main ex-
612 613	perimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?
614	Answer: [Yes]
C1E	Institution: In Section 4.1, we introduced the details of experimental sature and model
616	training to ensure reproducibility.
617	Guidelines
017	• The ensure NA means that the paper does not include experiments
618	• The answer NA means that the paper does not include experiments.
619	• If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important regardless of
620	whether the code and data are provided or not.
622	• If the contribution is a dataset and/or model, the authors should describe the steps taken
623	to make their results reproducible or verifiable.
624	• Depending on the contribution, reproducibility can be accomplished in various ways.
625	For example, if the contribution is a novel architecture, describing the architecture fully
626	might suffice, or if the contribution is a specific model and empirical evaluation, it may
627	dataset or provide access to the model. In general, releasing code and data is often
628	one good way to accomplish this but reproducibility can also be provided via detailed
630	instructions for how to replicate the results, access to a hosted model (e.g., in the case
631	of a large language model), releasing of a model checkpoint, or other means that are
632	appropriate to the research performed.
633	• While NeurIPS does not require releasing code, the conference does require all submis-
634	sions to provide some reasonable avenue for reproducibility, which may depend on the
635	nature of the contribution. For example
636	(a) If the contribution is primarily a new algorithm, the paper should make it clear how
637	to reproduce that algorithm.
638	(b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully
640	(c) If the contribution is a new model (e.g., a large language model), then there should
640	either be a way to access this model for reproducing the results or a way to reproduce
642	the model (e.g., with an open-source dataset or instructions for how to construct
643	the dataset).
644	(d) We recognize that reproducibility may be tricky in some cases, in which case
645	authors are welcome to describe the particular way they provide for reproducibility.
646	In the case of closed-source models, it may be that access to the model is limited in
647	some way (e.g., to registered users), but it should be possible for other researchers
648	to nave some path to reproducing or verifying the results.
649 5.	Open access to data and code

650 651 652	Question: Does the paper provide open access to the data and code, with sufficient instruc- tions to faithfully reproduce the main experimental results, as described in supplemental material?
653	Answer: [No]
654	Institution: The code will be released once the submission is accepted
0.55	Guidalinas:
655	
656	• The answer NA means that paper does not include experiments requiring code.
657	• Please see the NeurIPS code and data submission guidelines (https://nips.cc/
658	• While we appeared the release of each and data we understand that this might not be
660	possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
661	including code, unless this is central to the contribution (e.g., for a new open-source
662	benchmark).
663	• The instructions should contain the exact command and environment needed to run to
664	reproduce the results. See the NeurIPS code and data submission guidelines (https:
665	//nips.cc/public/guides/CodeSubmissionPolicy) for more details.
666	• The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, atc
669	• The authors should provide scripts to reproduce all experimental results for the new
669	proposed method and baselines. If only a subset of experimental results for the new
670	should state which ones are omitted from the script and why.
671	• At submission time, to preserve anonymity, the authors should release anonymized
672	versions (if applicable).
673	• Providing as much information as possible in supplemental material (appended to the
674	paper) is recommended, but including URLs to data and code is permitted.
675	6. Experimental Setting/Details
676	Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
677	parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
678	results?
679	Answer: [Yes]
680 681	Justification: In Section 4.1, we introduced the details of experimental setup and model training and testing.
682	Guidelines:
683	• The answer NA means that the paper does not include experiments.
684	• The experimental setting should be presented in the core of the paper to a level of detail
685	that is necessary to appreciate the results and make sense of them.
686	• The full details can be provided either with the code, in appendix, or as supplemental
687	material.
688	7. Experiment Statistical Significance
689	Question: Does the paper report error bars suitably and correctly defined or other appropriate
690	information about the statistical significance of the experiments?
691	Answer: [No]
692	Justification: The experiments conducted in our paper do not involve the use of error bars or
693	statistical significance analysis, thus this aspect is not applicable to our study.
694	Guidelines:
695	• The answer NA means that the paper does not include experiments.
696	• The authors should answer "Yes" if the results are accompanied by error bars, confi-
697	dence intervals, or statistical significance tests, at least for the experiments that support
698	the main claims of the paper.
699	• The factors of variability that the error bars are capturing should be clearly stated (for
700 701	example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions)
	run man Biven experimental conditions).

702 703	• The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
704	• The assumptions made should be given (e.g., Normally distributed errors).
705	• It should be clear whether the error bar is the standard deviation or the standard error
706	of the mean.
707	• It is OK to report 1-sigma error bars, but one should state it. The authors should
708	preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
709	of Normality of errors is not verified.
710	• For asymmetric distributions, the authors should be careful not to show in tables or
711	figures symmetric error bars that would yield results that are out of range (e.g. negative
712	error rates).
713 714	• If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
715	8. Experiments Compute Resources
716	Question: For each experiment, does the paper provide sufficient information on the com-
717 718	puter resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?
719	Answer: [Yes]
720	Justification: For our experiments, we furnished detailed specifications of the GPU models
721	used along with their corresponding tasks. Furthermore, we included specific information
722	regarding the model training batch size and the number of training iterations.
723	Guidelines:
724	• The answer NA means that the paper does not include experiments.
725	• The paper should indicate the type of compute workers CPU or GPU, internal cluster,
726	or cloud provider, including relevant memory and storage.
727	• The paper should provide the amount of compute required for each of the individual
728	experimental runs as well as estimate the total compute.
729	• The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that
730	didn't make it into the paper).
732	9. Code Of Ethics
733	Ouestion: Does the research conducted in the paper conform, in every respect, with the
734	NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
735	Answer: [Yes]
736	Justification: We have carefully reviewed the NeurIPS Code of Ethics and adhere to its
737	principles.
738	Guidelines:
739	• The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
740	• If the authors answer No, they should explain the special circumstances that require a
741	deviation from the Code of Ethics.
742	• The authors should make sure to preserve anonymity (e.g., if there is a special consid-
743	a b b b b b b b b b b b b b b b b b b b
744 I	0. Broader Impacts
745 746	Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?
747	Answer: [Yes]
748	Justification: We discuss both potential positive societal impacts and negative societal
749	impacts of the work performed in Appendix.
750	Guidelines:
751	• The answer NA means that there is no societal impact of the work performed.

752 753	• If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
754	• Examples of negative societal impacts include potential malicious or unintended uses
755	(e.g., disinformation, generating fake profiles, surveillance), fairness considerations
756	(e.g., deployment of technologies that could make decisions that unfairly impact specific
757	groups), privacy considerations, and security considerations.
758	• The conference expects that many papers will be foundational research and not tied
759	to particular applications, let alone deployments. However, if there is a direct path to
760	any negative applications, the authors should point it out. For example, it is legitimate
761	to point out that an improvement in the quality of generative models could be used to
762	generate deeplakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train
763	models that generate Deepfakes faster
765	• The authors should consider possible harms that could arise when the technology is
766	being used as intended and functioning correctly, harms that could arise when the
767	technology is being used as intended but gives incorrect results, and harms following
768	from (intentional or unintentional) misuse of the technology.
769	• If there are negative societal impacts, the authors could also discuss possible mitigation
770	strategies (e.g., gated release of models, providing defenses in addition to attacks,
771	mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
772	feedback over time, improving the efficiency and accessibility of ML).
773 1	1. Safeguards
774	Question: Does the paper describe safeguards that have been put in place for responsible
775	image generators, or scraped datasets)?
777	Answer: [NA]
778	Justification: Our paper poses no such risks.
770	Guidelines:
//9	
780	• The answer NA means that the paper poses no such risks.
781	• Released models that have a high risk for misuse or dual-use should be released with
782	the users adhere to usage guidelines or restrictions to access the model or implementing
783	safety filters
785	• Datasets that have been scraped from the Internet could nose safety risks. The authors
786	should describe how they avoided releasing unsafe images.
787	• We recognize that providing effective safeguards is challenging and many papers do
788	not require this, but we encourage authors to take this into account and make a best
789	faith effort.
790 1	2. Licenses for existing assets
791	Question: Are the creators or original owners of assets (e.g., code, data, models), used in
792	the paper, properly credited and are the license and terms of use explicitly mentioned and
793	properly respected?
794	Answer: [Yes]
795	Justification: The creators or original owners of assets, such as code, data, or models, used
796	in the paper, are properly credited. Additionally, the license and terms of use associated
797	with these assets are explicitly mentioned and respected in accordance with ethical and legal
798	standards.
799	Guidelines:
800	
	• The answer INA means that the paper does not use existing assets.
801	 The answer NA means that the paper does not use existing assets. The authors should cite the original paper that produced the code package or dataset.
801 802	 The answer NA means that the paper does not use existing assets. The authors should cite the original paper that produced the code package or dataset. The authors should state which version of the asset is used and, if possible, include a UPI
801 802 803	 The answer NA means that the paper does not use existing assets. The authors should cite the original paper that produced the code package or dataset. The authors should state which version of the asset is used and, if possible, include a URL. The neuro of the ligner (a p. CC BN 4.0) has the held formula bet formula.

805 806		• For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
807		• If assets are released, the license, copyright information, and terms of use in the
808		package should be provided. For popular datasets, paperswithcode.com/datasets
809		has curated licenses for some datasets. Their licensing guide can help determine the
810		license of a dataset.
811 812		• For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
813 814		• If this information is not available online, the authors are encouraged to reach out to the asset's creators.
815	13.	New Assets
816 817		Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?
818		Answer: [NA]
819		Justification: Our paper does not release new assets.
820		Guidelines:
821		• The answer NA means that the paper does not release new assets.
822		• Researchers should communicate the details of the dataset/code/model as part of their
823		submissions via structured templates. This includes details about training, license,
824		limitations, etc.
825		• The paper should discuss whether and how consent was obtained from people whose
826		asset is used.
827		• At submission time, remember to anonymize your assets (if applicable). You can either
828		create an anonymized URL or include an anonymized zip file.
829	14.	Crowdsourcing and Research with Human Subjects
830		Question: For crowdsourcing experiments and research with human subjects, does the paper
831		include the full text of instructions given to participants and screenshots, if applicable, as
832		well as details about compensation (if any)?
833		Answer: [NA]
834		Justification: Our paper does not involve crowdsourcing nor research with human subjects.
835		Guidelines:
836		• The answer NA means that the paper does not involve crowdsourcing nor research with
837		human subjects.
838		• Including this information in the supplemental material is fine, but if the main contribu-
839		tion of the paper involves human subjects, then as much detail as possible should be
840		included in the main paper.
841		• According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
842		or other labor should be paid at least the minimum wage in the country of the data
843		collector.
844	15.	Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
845		Subjects
846		Question: Does the paper describe potential risks incurred by study participants, whether
847		such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
848		approvals (or an equivalent approval/review based on the requirements of your country or
849		institution) were obtained?
850		Answer: [NA]
851		Justification: Our paper does not involve crowdsourcing nor research with human subjects.
852		Guidelines:
853		• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects
854		numan subjects.

855	• Depending on the country in which research is conducted, IRB approval (or equivalent)
856	may be required for any human subjects research. If you obtained IRB approval, you
857	should clearly state this in the paper.
858	• We recognize that the procedures for this may vary significantly between institutions
859	and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
860	guidelines for their institution.
861	• For initial submissions, do not include any information that would break anonymity (if
862	applicable), such as the institution conducting the review.