

Make a Cheap Scaling: A Self-Cascade Diffusion Model for Higher-Resolution Adaptation

Lanqing Guo^{1,2†}, Yingqing He^{2,3†}, Haoxin Chen², Menghan Xia²,
Xiaodong Cun², Yufei Wang¹, Siyu Huang⁴, Yong Zhang^{2*}, Xintao
Wang², Qifeng Chen³, Ying Shan², and Bihan Wen^{1*}

¹ Nanyang Technological University

² Tencent AI Lab

³ The Hong Kong University of Science and Technology

⁴ Clemson University

Project page: <https://guolanqing.github.io/Self-Cascade/>

Abstract. Diffusion models have proven to be highly effective in image and video generation; however, they encounter challenges in the correct composition of objects when generating images of varying sizes due to single-scale training data. Adapting large pre-trained diffusion models to higher resolution demands substantial computational and optimization resources, yet achieving generation capabilities comparable to low-resolution models remains challenging. This paper proposes a novel self-cascade diffusion model that leverages the knowledge gained from a well-trained low-resolution image/video generation model, enabling rapid adaptation to higher-resolution generation. Building on this, we employ the pivot replacement strategy to facilitate a tuning-free version by progressively leveraging reliable semantic guidance derived from the low-resolution model. We further propose to integrate a sequence of learnable multi-scale upsampler modules for a tuning version capable of efficiently learning structural details at a new scale from a small amount of newly acquired high-resolution training data. Compared to full fine-tuning, our approach achieves a $5\times$ training speed-up and requires only 0.002M tuning parameters. Extensive experiments demonstrate that our approach can quickly adapt to higher-resolution image and video synthesis by fine-tuning for just 10k steps, with virtually no additional inference time.

Keywords: Image Synthesis · Video Synthesis · Diffusion Model · Higher-Resolution Adaptation

1 Introduction

Over the past two years, stable diffusion (SD) [7, 21] has sparked great interest in generative models, gathering attention from both academic and industry. It has demonstrated impressive outcomes across diverse generative applications, *e.g.*,

[†] Equal Contributions

^{*} Corresponding Authors

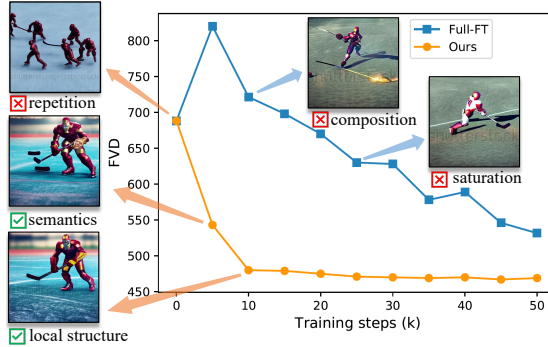


Fig. 1: The FVD \downarrow score averages for both the full fine-tuning (Full-FT) and our proposed fast adaptation method (Ours) are assessed every 5k iterations on the *Webvid-10M* [1] benchmark. We observe that full fine-tuning necessitates a large number of training steps and suffers from poor composition ability and desaturation issues. In contrast, our method enables rapid adaptation to the higher-resolution domain while preserving reliable semantic and local structure generation capabilities.

text-to-image generation [6, 10, 21, 23], image-to-image translation [25, 31], and text-to-video generation [2, 13, 28, 36, 38]. To scale up SD models to high-resolution content generation, a commonly employed approach is progressive training [7, 21], *i.e.*, training the SD model with lower-resolution images before fine-tuning with higher-resolution images. This warm-up approach enhances the model’s semantic composition ability, leading to the generation of high-quality images. However, even a well-trained diffusion model for low-resolution images demands extensive fine-tuning and computational resources when transferring to a high-resolution domain due to its large size of model parameters. For instance, SD 2.1 [7] requires 550k steps of training at resolution 256^2 before fine-tuning with $>1000k$ steps at resolution 512^2 to enable 512^2 image synthesis. Insufficient tuning steps may severely degrade the model’s composition ability, resulting in issues such as pattern repetition, desaturation, and unreasonable object structures as in Fig. 1.

Several tuning-free methods, such as those proposed in [19] and ScaleCrafter [12], attempted to seamlessly adapt the SD to higher-resolution image generation with reduced efforts. In [19], the authors explored SD adaptation for variable-sized image generation using attention entropy, while ScaleCrafter [12] utilized dilated convolution to enlarge the receptive field of convolutional layers and adapt to new resolution generation. However, these tuning-free solutions require careful adjustment of factors such as the dilated stride and injected step, potentially failing to account for the varied scales of object generation. More recent methods, such as those proposed in [40], have attempted to utilize LORA [18] as additional parameters for fine-tuning. However, this approach is not specifically designed for scale adaptation and still requires a substantial number of tuning steps. Other works [15, 35, 39] proposed to cascade the super-resolution mechanisms based on diffusion models for scale enhancement. However, the use of extra super-resolution models necessitates a doubling of training parameters and limits the scale extension ability for a higher resolution.

In this paper, we present a novel self-cascade diffusion model that harnesses the rich knowledge gained from a well-trained low-resolution generation model to enable rapid adaptation to higher-resolution generation. Our approach begins with the introduction of a tuning-free version, which utilizes a pivot replacement strategy to enforce the synthesis of detailed structures at a new scale by injecting reliable semantic guidance derived from the low-resolution model. Building on this baseline, we further propose time-aware feature upsampling modules as plugins to a base low-resolution model to conduct a tuning version. To enhance the robustness of scale adaptation while preserving the model’s original composition and generation capabilities, we fine-tune the plug-and-play and lightweight upsampling modules at different feature levels, using a small amount of acquired high-quality data with a few tuning steps.

The proposed upsampler modules can be flexibly plugged into any pre-trained SD-based models, including both image and video generation models. Compared to full fine-tuning, our approach offers a training speed-up of more than 5 times and requires only 0.002M trainable parameters. Extensive experiments demonstrated that our proposed method can rapidly adapt to higher-resolution image and video synthesis with just 10k fine-tuning steps and virtually no additional inference time.

Our main contributions are summarized as follows:

- We propose a novel self-cascade diffusion model for fast-scale adaptation to higher resolution generation, by cyclically re-utilizing the low-resolution diffusion model. Based on that, we employ a pivot replacement strategy to enable a tuning-free version as the baseline.
- We further construct a series of plug-and-play, learnable time-aware feature upsampler modules to incorporate knowledge from a few high-quality images for fine-tuning. This approach achieves a $5\times$ training speed-up compared to full fine-tuning and requires only 0.002M learnable parameters.
- Comprehensive experimental results on image and video synthesis demonstrate that the proposed method attains state-of-the-art performance in both tuning-free and tuning settings across various scale adaptations.

2 Related Work

Stable diffusion. Building upon the highly effective and efficient foundations established by the Latent Diffusion Model (LDM) [24], diffusion models [14, 30] have recently demonstrated remarkable performance in various practical applications, *e.g.*, text-to-image generation [6, 10, 21, 23], image-to-image translation [25, 31], and text-to-video generation [2, 13, 28, 36, 38]. In this field, stable diffusion (SD) [21, 24] has emerged as a prominent model for generating photo-realistic images from text. However, despite its impressive synthesis capabilities at specific resolutions (*e.g.*, 512^2 for SD 2.1 and 1024^2 for SD XL), it often produces extremely unnatural outputs for unseen image sizes. This limitation mainly arises from the fact that current SD models are trained exclusively on fixed-size images, leading to a lack of varied resolution generalizability. In this paper, we aim to explore the fast

adaptation ability of the original diffusion model with limited image size to a higher resolution.

High-resolution synthesis and adaptation. Although existing stable diffusion-based synthesis methods have achieved impressive results, high-resolution image generation remains challenging and demands substantial computational resources, primarily due to the complexity of learning from higher-dimensional data. Additionally, the practical difficulty of collecting large-scale, high-quality image and video training datasets further constrains synthesis performance. To address these challenges, prior work can be broadly categorized into three main approaches:

1. Training from scratch. This type of work can be further divided into two categories: cascaded models [9, 15, 16, 32] and end-to-end models [3, 5, 17, 21]. Cascade diffusion models employ an initial diffusion model to generate lower-resolution data, followed by a series of super-resolution diffusion models to successively upsample it. End-to-end methods learn a diffusion model and directly generate high-resolution images in one stage. However, they all necessitate sequential, separate training and a significant amount of training data at high resolutions.
2. Fine-tuning. Parameter-efficient tuning is an intuitive solution for higher-resolution adaptation. DiffFit [37] utilized a customized partial parameter tuning approach for general domain adaptation. Zheng *et al.* [40] adopted the LORA [18] as the additional parameters for fine-tuning, which is still not specifically designed for the scale adaptation problem and still requires huge of tuning steps.
3. Tuning-free. Several methods [8, 11, 12, 19] have explored expanding low-resolution diffusion models to higher resolutions without tuning. Recently, Jin *et al.* [19] explored a tuning-free approach for variable sizes but did not address high-resolution generation. ScaleCrafter [12] employed dilated convolution to expand the receptive field of convolutional layers for adapting to new resolutions. Besides, DemoFusion [8] used low-resolution semantic guidance as well as dilated sampling strategy to achieve a high-resolution generation. However, these approaches require careful adjustments, such as dilated stride and injected step, which lack semantic constraints and result in artifacts for various generation scales.

3 Preliminary

Our proposed method is based on the recent text-to-image diffusion model (*i.e.*, stable diffusion (SD) [21, 24]), which formulates the diffusion and denoising process in a learned low-dimensional latent space. An autoencoder first conducts perceptual compression to significantly reduce the computational cost, where the encoder E converts image $x_0 \in \mathbb{R}^{3 \times H \times W}$ to its latent code $z_0 \in \mathbb{R}^{4 \times H' \times W'}$ and the decoder D reconstructs the image x_0 from z_0 as follows,

$$z_0 = E(x_0) , \quad \hat{x}_0 = D(z_0) \approx x_0 . \quad (1)$$

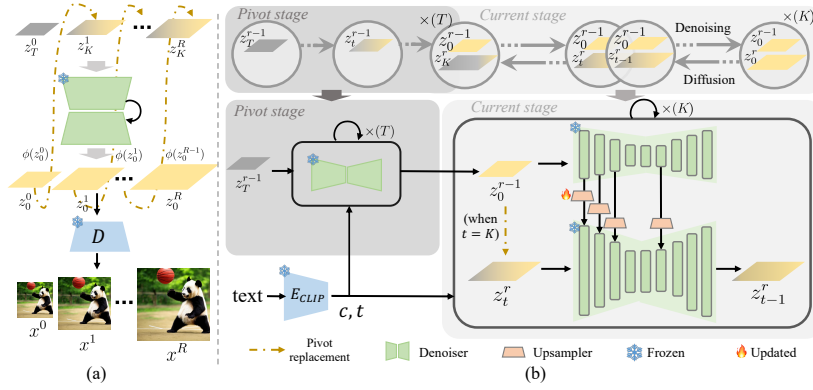


Fig. 2: Illustration of the proposed self-cascade diffusion model, which is implemented in both tuning-free and tuning versions. (a) For the tuning-free version, we cyclically re-utilize the low-resolution model to progressively adapt it to the higher-resolution generation; (b) For the tuning version, we additionally plug feature upsamplers (Φ) into the base low-resolution generation model: the denoising process of image z_t^r in step t will be guided by the pivot guidance z_0^{r-1} from the pivot stage (last stage) with a series of plugged-in tuneable upsampler modules.

Then, the diffusion model formulates a fixed forward diffusion process to gradually add noise to the latent code $z_0 \sim p(z_0)$:

$$q(z_t|z_0) = \mathcal{N}(z_t; \sqrt{\alpha_t}z_0, (1 - \alpha_t)\mathbf{I}). \quad (2)$$

In the inference stage, we sample latent features from the conditional distribution $p(z_0|c)$ with the conditional information c (e.g., the text embedding extracted by a CLIP encoder [22] E_{CLIP}):

$$p_\theta(z_{0:T}|c) = p(z_T) \prod_{t=1}^T p_\theta(z_{t-1}|z_t, c). \quad (3)$$

The U-Net denoiser ϵ_θ consists of a sequential transformer and convolution blocks to perform denoising in the latent space. The corresponding optimization process can be defined as the following formulation:

$$\mathcal{L} = \mathbb{E}_{z_t, c, \epsilon, t} (\|\epsilon - \epsilon_\theta(z_t, t, c)\|^2), \quad (4)$$

where $z_t = \sqrt{\alpha_t}z_0 + \sqrt{1 - \alpha_t}\epsilon$ represents the noised feature map at step t .

4 Methodology

In this section, we first introduce the overall framework of the proposed self-cascade diffusion model (Sec. 4.1). Based on that, we propose a tuning-free version using a pivot replacement strategy as the baseline, as well as an improved tuning version by plugging tunable feature upsamplers (Sec. 4.2). We then provide an analysis and discussion on our self-cascade diffusion model (Sec. 4.3).

4.1 Self-Cascade Diffusion Model

Given a pre-trained stable diffusion (SD) model with the denoiser $\epsilon_\theta(\cdot)$ for synthesizing low-resolution images (latent code) $z \in \mathbb{R}^d$, our goal is to generate higher-resolution images $z^R \in \mathbb{R}^{d_R}$ in a time/resource and parameter-efficient manner with an adapted model $\tilde{\epsilon}_\theta(\cdot)$. To achieve such a goal, we aim to reuse the rich knowledge from the well-trained low-resolution model and only learn the low-level details at a new scale. Thus, we propose a self-cascade diffusion model to cyclically re-utilize the low-resolution image synthesis model. We intuitively define a **scale decomposition** to decompose the whole scale adaptation $\mathbb{R}^d \rightarrow \mathbb{R}^{d_R}$ into multiple progressive adaptation processes such that $d = d_0 < d_1 \dots < d_R$ where $R = \lceil \log_4 d_R/d \rceil$. We first progressively synthesize a low-resolution image (latent code) z^{r-1} and then utilize it as the pivot guidance to synthesize the higher resolution result z^r in the next stage, where the reverse process of the self-cascade diffusion model can be extended by Eq. (3) for each z^r , $r = 1, \dots, R$ as follows:

$$p_\theta(z_{0:T}^r | c, z_0^{r-1}) = p(z_T^r) \prod_{t=1}^T p_\theta(z_{t-1}^r | z_t^r, c, z_0^{r-1}), \quad (5)$$

where the reverse transition $p_\theta(z_{t-1}^r | z_t^r, c, z_0^{r-1})$ not only conditions on denoising step t and text embedding c , but also on lower-resolution latent code z_0^{r-1} generated in the last stage. Different from previous works, *e.g.*, [16], LAVIE [35], and SHOW-1 [39] where p_θ in Eq. 9 is implemented by a new super-resolution model, we cyclically re-utilize the base low-resolution synthesis model to inherit the prior knowledge of the base model thus improve the efficiency.

Pivot replacement. According to the *scale decomposition*, the whole scale adaptation process will be decoupled into multiple moderate adaptation stages, *e.g.*, $4\times$ more pixels than the previous stage. The information capacity gap between z^r and z^{r-1} is not significant, especially in the presence of noise (intermediate step of diffusion). Consequently, we assume that $p(z_K^r | z_0^{r-1})$ can be considered as the proxy for $p(z_K^r | z_0^r)$ to manually set the initial diffusion state for current adaptation stage $\mathbb{R}^{d_{r-1}} \rightarrow \mathbb{R}^{d_r}$, where $K < T$ is an intermediate step. Specifically, let ϕ denote a deterministic resize interpolation function (*i.e.*, bilinear interpolation) to upsample from scale d_{r-1} to d_r . We upsample the generated lower-resolution image z_0^{r-1} from last stage into $\phi(z_0^{r-1})$ to maintain dimensionality. Then we can diffuse it by K steps and use it to replace z_K^r as follows:

$$z_K^r \sim \mathcal{N}(\sqrt{\bar{\alpha}_K} \phi(z_0^{r-1}), \sqrt{1 - \bar{\alpha}_K} \mathbf{I}). \quad (6)$$

Regarding z_K^r as the initial state for the current stage and conduct denoising with $K \rightarrow 0$ steps as Eq. (3) to generate the z_0^r , which is the generated higher-resolution image in the current stage.

We can employ such a pivot replacement strategy at all decoupled scale adaptation stages. Hence, the whole synthesis process for a higher-resolution image with resolution d_R using pivot replacement strategy can be illustrated as Fig. 2(a). So far, we have devised a **tuning-free version** of self-cascade

diffusion model (denoted as Ours-TF in experiments) to progressively expand the model capacity for higher-resolution adaptation with cyclically re-utilizing the totally frozen low-resolution model. Although the tuning-free self-cascade diffusion model built upon pivot replacement strategy (Sec. 4.1) can achieve a feasible and scale-free higher-resolution adaptation, it has limitations on synthesis performance especially the detailed low-level structures due to the unseen higher-resolution ground-truth images. To achieve a more practical and robust scale adaptation performance, we further introduce an improved **tuning version** of the self-cascade diffusion model (denoted as Ours-T in experiments) in Sec. 4.2.

4.2 Feature Upsampler Tuning

In this section, we propose a tuning version of the self-cascade diffusion model that enables a cheap scaling, by inserting very lightweight time-aware feature upsamplers as illustrated in Fig. 2(b). The proposed upsamplers can be plugged into any diffusion-based synthesis methods. The detailed tuning and inference processes of our tuning version self-cascade diffusion model are in Algorithm 1 and 2, respectively. Note that by omitting the tuning process and solely executing the inference step in Algorithm 2, it turns into our tuning-free version.

Specifically, given an intermediate z_t^r in step t and the pivot guidance z_0^{r-1} from the last stage, we can achieve corresponding intermediate multi-scale feature groups h_t^r and h_0^{r-1} via the pre-trained UNet denoiser ϵ_θ , respectively, as follows:

$$\begin{aligned} h_0^{r-1} &= \{h_{1,0}^{r-1}, h_{2,0}^{r-1}, \dots, h_{N,0}^{r-1}\} \\ h_t^r &= \{h_{1,t}^r, h_{2,t}^r, \dots, h_{N,t}^r\} \end{aligned} \quad (7)$$

where N represents the number of features within each feature group and the details are included in the **supplementary**. In short, inspired by the recent work [27] that investigated the impact of various components in the UNet architecture on synthesis performance, we choose to use skip features as a feature group. These features have a negligible effect on the quality of the generated images while still providing semantic guidance. We define a series of time-aware feature upsamplers $\Phi = \{\phi_1, \phi_2, \dots, \phi_N\}$ to upsample and transform pivot features at each corresponding scale. During the diffusion generation process, the focus shifts from high-level semantics to low-level detailed structures as the signal-to-noise ratio progressively increases as noise is gradually removed. Consequently, we propose that the learned upsampler transformation should be adaptive to different time steps. The upsampled features $\phi_n(h_{n,0}^{r-1}, t)$ is then added to original features $h_{n,t}^r$ at each scale:

$$\hat{h}_{n,t}^r = h_{n,t}^r + \phi_n(h_{n,0}^{r-1}, t), \quad n \in \{1, \dots, N\}. \quad (8)$$

Optimization details. For each training iteration for scale adaptation $\mathbb{R}^{d_{r-1}} \rightarrow \mathbb{R}^{d_r}$, we first randomly sample a step index $t \in (0, K]$. The corresponding optimization process can be defined as the following formulation:

$$\mathcal{L} = \mathbb{E}_{z_t^r, z_0^{r-1}, t, c, \epsilon, t} (\|\epsilon - \tilde{\epsilon}_{\theta, \theta_\Phi}(z_t^r, t, c, z_0^{r-1})\|^2), \quad (9)$$

Algorithm 1 Feature upsampler tuning process.

```

1: while not converged do
2:    $(x_0, c) \sim p(x_0, c)$ 
3:    $z_0^r = E(x_0)$ 
4:    $z_0^{r-1} = E(\text{Downsample}(x_0))$ 
5:    $t \sim \text{Uniform}\{1, \dots, K\}$ 
6:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
7:    $z_t^r = \sqrt{\bar{\alpha}_t} z_0^r + \sqrt{1 - \bar{\alpha}_t} \epsilon$ 
8:    $\theta_\phi \leftarrow \theta_\phi - \eta \nabla_{\theta_\phi} \|\tilde{\epsilon}_{\theta, \theta_\phi}(z_t^r, t, c, z_0^{r-1}) - \epsilon\|^2$ 
9: end while
10: return  $\theta_\phi$ 

```

Algorithm 2 Inference process for $\mathbb{R}^{d_{r-1}} \rightarrow \mathbb{R}^{d_r}$.

```

Input: text embedding  $c$ 
1: if  $r = 1$  then
2:    $z_T^r \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
3:   for  $t = T, \dots, 1$  do
4:      $z_{t-1}^r \sim p_\theta(z_{t-1}^r | z_t^r, c)$ 
5:   end for
6: else
7:    $z_K^r \sim q(z_K^r | z_0^{r-1})$ 
8:   for  $t = K, \dots, 1$  do
9:      $z_{t-1}^r \sim p_\theta(z_{t-1}^r | z_t^r, c, z_0^{r-1})$ 
10:  end for
11: end if
12: return  $z_0^r$ 

```

where θ_ϕ denotes the trainable parameters of the plugged-in upsamplers and θ denotes the frozen parameters of pre-trained diffusion denoiser. Each upsampler is simple and lightweight, consisting of one bilinear upsampling operation and two residual blocks. In all experiments, we set $N = 4$, resulting in a total of 0.002M trainable parameters. Therefore, the proposed tuning self-cascade diffusion model requires only a few tuning steps (*e.g.*, $10k$) and the collection of a small amount of higher-resolution new data.

4.3 Analysis and Discussion

Drawing inspiration from previous explorations on scale adaptation [12], we found that directly applying the SD 2.1 model trained with 512^2 images to generate 1024^2 images led to issues such as object repetition and diminished composition capacity (see Fig. 1). We observed that the local structural details of the generated images appeared reasonable and abundant without smoothness when the adapted scale was not large (*e.g.*, $4\times$ more pixels). In summary, the bottleneck for adapting to higher resolutions lies in the semantic component and composition capacity. Fortunately, the original pre-trained low-resolution diffusion model can generate a reliable low-resolution pivot, naturally providing proper semantic guidance by injecting the pivot semantic features during the higher-resolution diffusive sampling process. Simultaneously, the local structures can be completed based on the rich texture prior learned by the diffusion model itself, under strong semantic constraints.

Compared to existing cascaded diffusion frameworks for high-fidelity image and video generation [16], our work is the first to conduct self-cascade by cyclically re-utilizing pre-trained diffusion model on low-resolution with the following major advantages:

- **Lightweight upsampler module.** Conventional cascade diffusion models comprise a pipeline of multiple diffusion models that generate images of

Table 1: Quantitative results of different methods on the dataset of *Laion-5B* with $4\times$ adaptation on 1024^2 resolution. The best results are highlighted in **bold**. Note that Ours-TF and Ours-T denote the tuning-free version and the upsampler tuning version, respectively. “#Param” denotes the number of trainable parameters and “Infer Time” denotes the inference time of different methods *v.s.* Direct Inference. pFID_r/pKID_r denote the patch- We put ‘-’ since FID_b/KID_b are unavailable for SD+SR.

Methods	#Param	Training Step	Infer Time	FID _r ↓	KID _r ↓	pFID _r ↓	pKID _r ↓	FID _b ↓	KID _b ↓
Direct Inference	0	-	1×	29.89	0.010	20.88	0.0070	24.21	0.007
Attn-SF [19]	0	-	1×	29.95	0.010	21.07	0.0072	22.75	0.007
ScaleCrafter [12]	0	-	1×	20.88	0.008	21.00	0.0071	16.67	0.005
Ours-TF (Tuning-Free)	0	-	1.04×	12.25	0.004	19.59	0.0071	6.09	0.001
Full Fine-tuning	860M	18k	1×	21.88	0.007	19.33	0.0077	17.14	0.005
LORA-R32	15M	18k	1.22×	17.02	0.005	18.65	0.0076	11.33	0.003
LORA-R4	1.9M	18k	1.20×	14.74	0.005	18.06	0.0074	9.47	0.002
SD+SR	184M	1.25M	5×	12.59	0.005	17.21	0.0053	-	-
Ours-T (Tuning)	0.002M	4k	1.06×	12.40	0.004	15.35	0.0058	3.15	0.0005

Table 2: Quantitative results of different methods on the dataset of *Laion-5B* with $16\times$ image scale adaptation to 2048^2 resolution. The best results are highlighted in **bold**. 10k and 20k denote the training steps used for tuning.

Methods	FID _r ↓	KID _r ↓	FID _b ↓	KID _b ↓
Direct Inference	104.70	0.043	104.10	0.040
Attn-SF [19]	104.34	0.043	103.61	0.041
ScaleCrafter [12]	59.40	0.021	57.26	0.018
Ours-TF (Tuning-Free)	38.99	0.015	34.73	0.013
Full Fine-tuning (20k)	43.55	0.014	41.58	0.012
LORA-R4 (20k)	50.72	0.020	51.99	0.019
Ours-T (Tuning) (10k)	18.46	0.005	8.99	0.001

increasing resolution, which results in a multiplicative increase in the number of model parameters. Our model is built upon the shared diffusion model at each stage with only very lightweight upsampler modules (*i.e.*, 0.002M parameters) to be tuned.

- **Less fine-tuning data.** Previous cascaded model chains necessitate sequential, separate training, with each model being trained from scratch, thereby imposing a significant training burden. Our model is designed to quickly adapt the low-resolution synthesis model to higher resolutions using a small amount of high-quality data for fine-tuning.

5 Experiments

5.1 Implementation Details

The proposed method is implemented using PyTorch and trained on two NVIDIA A100 GPUs. The original base diffusion model’s parameters are frozen, with the only trainable component being the integrated upsampling modules. The initial learning rate is 5×10^{-5} . We used 1000 diffusion steps T for training, and 50 steps for DDIM [29] inference. We set $N = 4$ and $K = 700$ for all experiments. We conduct evaluation experiments on text-to-image models, specifically Stable Diffusion (SD), focusing on two widely-used versions: SD 2.1 [7] and SD XL 1.0 [21], as they adapt to two unseen higher-resolution domains. For the original SD 2.1, which is trained with 512^2 images, the inference resolutions are 1024^2 and 2048^2 , corresponding to $4\times$ and $16\times$ more pixels than the training, respectively. We also conduct evaluation experiments on text-to-video models, where we select the LVDM [13] as the base model which is trained with 16×256^2 videos (16 frames), the inference resolutions are 16×512^2 , $4\times$ more pixels than the base resolution. We left the experiments for SD XL 1.0 in the **supplementary**.

5.2 Evaluation on Image Generation

Dataset and evaluation metrics. We select the Laion-5B [26] as the benchmark dataset which contains 5 billion image-caption pairs. We fine-tune all tuning-based competing methods by applying online filtering on Laion-5B for high-resolution images larger than the target resolution. We randomly sample $30k$ images with text prompts from the dataset and evaluate the generated image quality and diversity using the Inception Distance (FID) and Kernel Inception Distance (KID) metrics, which are measured between the generated images and real images, denoted as FID_r and KID_r . Following previous work [12], we sample $10k$ images when the inference resolution is higher than 1024^2 . Besides, to address the issue of squeezed resolutions in standard FID_r/KID_r , we randomly cropped local patches to calculate these metrics instead of resizing, denoted as $pFID_r/pKID_r$ [4]. To ensure consistency in image pre-processing steps, we use the clean-fid implementation [20]. Since pre-trained models can combine different concepts that are not present in the training set, we also measure the FID and KID metrics between the generated samples under the base training resolution and inference resolution, denoted as FID_b and KID_b . This evaluation assesses how well our method preserves the model’s original ability when adapting to a new higher resolution.

Comparison with state-of-the-art. We conduct the comparison experiments on two settings, *i.e.*, tuning-free and fine-tuning. For the tuning-free setting,

We follow the same comparison settings of ScaleCrafter [12]. Since FID_b/KID_b are evaluated on the original low-resolution by down-sampling, the down-sampling results of SD+SR will be roughly the same as the reference real image set which denotes “zero distance”.

we compare our tuning-free version, denoted as Ours-TF, with the vanilla text-to-image diffusion model (Direct Inference) that directly samples the higher resolution images via the original checkpoint, as well as two tuning-free methods, *i.e.*, Attn-SF [19] and ScaleCrafter [12]. Besides, we also compare our fine-tuning version, denoted as Ours-T, with the full fine-tuning model, and LORA tuning (consisting of two different ranks, *i.e.*, 4 and 32, denoted as LORA-R4 and LORA-R32). Tab. 1 and Tab. 2 show the quantitative results on Laion-5B [26] over $4\times$ and $16\times$ more pixels compared to base model. Our methods outperform existing methods in both tuning-free and fine-tuning settings, especially when adapting to a higher-resolution domain, *e.g.*, $16\times$ scale adaptation. Besides, with the merits of injecting newly acquired higher-resolution data for tuning, our tuning version can achieve a better and more robust generation performance, especially for $16\times$ scale adaptation. We show random samples from our method (Ours-T) on adapted various resolutions in Fig. 3, comparing to results of Full fine-tuning (Full-FT) and LORA-R4. Our method can achieve not only high-quality image results with rich structural details and accurate object composition, *e.g.*, the relationship between the bear and motorbike as shown in Fig. 3. Visual comparisons with the competing methods are included in the **supplementary**.

5.3 Evaluation on Video Generation

Dataset and evaluation metrics. We select the Webvid-10M [1] as the benchmark dataset which contains 10M high-resolution collected videos. We randomly sample 2048 videos with text prompts from the dataset and evaluate the generated video quality and diversity using video counterpart Frechet Video Distance (FVD) [33] and Kernel Video Distance (KVD) [34], denoted as FVD_r and KVD_r .

Comparison with state-of-the-art. To comprehensively verify the effectiveness of our proposed method, we also conduct comparison experiments on a video generation base model [13]. Thus, we compare our method with a full fine-tuning model and LORA tuning with different ranks, as well as the previous tuning-free method, *i.e.*, ScaleCrafter. Tab. 3 shows the quantitative results on *Webvid-10M* [1] and visual comparisons are shown in Fig. 4. Our method achieves better FVD and KVD results in approximately 20% of the training steps compared to the competing approaches. With the merits of the reuse of reliable semantic guidance from a well-trained low-resolution diffusion model, our method can achieve better object composition ability (*e.g.*, the reaction between cat and yarn ball and the times square as shown in the second and fourth examples of Fig. 4, respectively) and rich local structures compared to the competing methods (*e.g.*, the details of the teddy bear as shown in the third example of Fig. 4). In contrast, for full fine-tuning models, the issue of low saturation and over-smoothness requires many training steps to resolve and it is difficult to achieve results as good as those obtained with low-resolution models. Besides, the generated results of both full fine-tuning and LORA tuning methods will have motion shift or motion inconsistency issues as shown in the bag of the astronaut in the first example of Fig. 4, while our method can better maintain the original



Fig. 3: Visual examples of Ours-T (Tuning) on the higher-resolution adaptation to various higher resolutions, *e.g.*, 1024^2 , 3072×1536 , 1536×3072 , and 2048^2 , with the pre-trained SD 2.1 trained with 512^2 images, comparing to 1024^2 results of Full fine-tuning (Full-FT) and LORA-R4 (right down corner: red dashed box). Please zoom in for more details.

model’s temporal consistency capabilities, generating more coherent videos (video examples refer to **supplementary**).

5.4 Network Analysis

Efficiency comparison. To demonstrate the training and sampling efficiency of our method, we compare our approach with selected competing methods in Tab. 1 for generating 1024^2 resolution images on the *Laion-5B* dataset. Our model has only 0.002M trainable parameters, utilizing approximately the parameters compared to LORA-R4 (with a rank of 4). Although our proposed method requires a cascaded generation process, *i.e.*, starting with low-resolution generation followed by progressively pivot-guided higher-resolution generation, the inference time of our method is similar to that of direct inference (with a factor of $1.04\times$ for the tuning-free version and $1.06\times$ for the tuning version), resulting in virtually no additional sampling time. Besides, we present the FID and FVD scores for several methods every $5k$ iteration on image (Laion-5B) and video (Webvid-10M) datasets as shown in Fig. 5. Our observations demonstrate that our method can

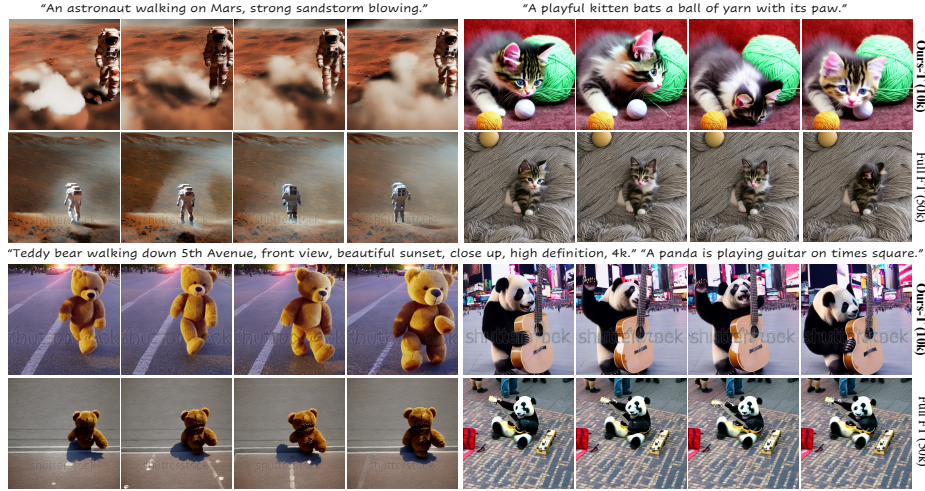


Fig. 4: Visual quality comparisons between full fine-tuning (50k) and Ours-T (10k) on higher-resolution video synthesis of 16×512^2 .

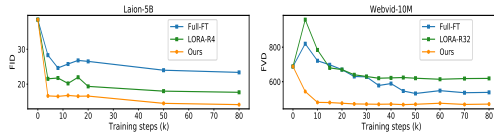


Fig. 5: Average FID and FVD scores of three methods every 5k iterations on image (Laion-5B) and video (Webvid-10M) datasets. Our observations indicate that our method can rapidly adapt to the higher-resolution domain while maintaining a robust performance among both image and video generation.

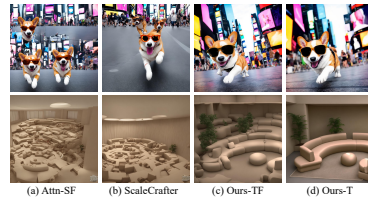


Fig. 6: Visual quality comparisons between the tuning-free methods and ours on higher-resolution adaptation to 1024^2 resolutions. Please zoom-in to see more details.

rapidly adapt to the desired higher resolution. By cyclically reusing the frozen diffusion base model and incorporating only lightweight upsampler modules, our approach maximally retains the generation capacity of the pretrained base model, resulting in improved fine-tuned performance.

Tuning-free or fine-tuning? Although our tuning-free self-cascade diffusion model can inject the correct semantic information to higher-resolution adaptation, some extreme examples still make it difficult to completely suppress repetition issues and composition capabilities, such as repetitive legs and sofas as shown in Fig. 6. Such failure case is particularly evident in the repetition of very fine-grain objects or texture, which is a common occurrence among all tuning-free competing methods, like Attn-SF [19] and ScaleCrafter [12]. By tuning plug-and-play and lightweight upsampler modules with a small amount of higher-resolution data, the diffusion model can learn the low-level details at a new scale.

Table 3: Quantitative results of different methods on the dataset of *Webvid-10M* [1] with $4\times$ video scale adaptation on 16×512^2 resolution (16 frames). The best results are highlighted in **bold**. 10k and 50k denote the training steps used for tuning.

Methods	FVD $_r\downarrow$	KVD $_r\downarrow$
Direct Inference	688.07	67.17
ScaleCrafter [12]	562.00	44.52
Ours-TF	553.85	33.83
Full Fine-tuning (10k)	721.32	94.57
Full Fine-tuning (50k)	531.57	33.61
LORA-R4 (10k)	1221.46	263.62
LORA-R32 (10k)	959.68	113.07
LORA-R4 (50k)	623.72	74.13
LORA-R32 (50k)	615.75	76.99
Ours-T (10k)	494.19	31.55



Fig. 7: Visual examples of video generation of the (a) low-resolution pivot samples generated by the pre-trained base model, (b) super-resolution result by SD+SR, and (c) high-resolution final output of our tuning approach.

Relation to the super-resolution methods. We also compare our approach to using a pre-trained Stable Diffusion super-resolution (SD 2.1-upscaler- $4\times$) as post-processing, denoted as SD+SR, for the higher-resolution generation as shown in Tab. 1. Our approach achieves better performance and reduced inference time, even in a tuning-free manner (Ours-TF). In contrast, SD+SR still requires a large amount of high-resolution data for training a new diffusion model with around 184M extra parameters to be trained. Furthermore, our method not only increases the resolutions of pivot samples like SD+SR, but also *explores the potential of the pre-trained diffusion model for fine-grained details generation and inheriting the composition capacity*. We illustrate one example of video generation in Fig. 7, demonstrating two key advantages of our method over SD+SR: (1) while the low-resolution pivot sample from the base model predicts an “object shift” result across temporal frames, our method effectively corrects such inconsistencies, which is not achievable by simply applying SD+SR; (2) our approach excels in synthesizing finer details and textures compared to using SR solely as post-processing, as shown in the enhanced zoom-in on the tiger region.

Limitations. Our method adapts well to higher-resolution domains but still has limitations. Since the number of parameters in the upsampler modules we insert is very small, there is an upper bound to the performance of our method when there is sufficient training data, especially when the scale gap is too large, *e.g.*, higher resolution than $4k$ resolution data. We will further explore the trade-off between adaptation efficiency and generalization ability in future work.

6 Conclusion

In this work, we present a novel self-cascade diffusion model for rapid higher-resolution adaptation. Our approach first introduces a pivot-guided noise re-

schedule strategy in a tuning-free manner, cyclically re-utilizing the well-trained low-resolution model. We then propose an efficient tuning version that incorporates a series of plug-and-play, learnable time-aware feature upsampler modules to interpolate knowledge from a small amount of newly acquired high-quality data. Our method achieves over 5x training speed-up with only 0.002M tuning parameters and negligible extra inference time. Experimental results demonstrate the effectiveness and efficiency of our approach plugged into various image and video synthesis base models over different scale adaptation settings.

Acknowledgements

This research was carried out at the Rapid-Rich Object Search (ROSE) Lab at Nanyang Technological University in Singapore. This work was supported in part by the National Research Foundation Singapore Competitive Research Program (award number CRP29-2022-0003), and in part by the National Key R&D Program of China under grant number 2022ZD0161501.

References

1. Bain, M., Nagrani, A., Varol, G., Zisserman, A.: Frozen in time: A joint video and image encoder for end-to-end retrieval. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1728–1738 (2021) 2, 11, 14
2. Blattmann, A., Rombach, R., Ling, H., Dockhorn, T., Kim, S.W., Fidler, S., Kreis, K.: Align your latents: High-resolution video synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 22563–22575 (2023) 2, 3
3. Bond-Taylor, S., Willcocks, C.G.: ∞ -diff: Infinite resolution diffusion with subsampled mollified states. arXiv preprint arXiv:2303.18242 (2023) 4
4. Chai, L., Gharbi, M., Shechtman, E., Isola, P., Zhang, R.: Any-resolution training for high-resolution image synthesis. In: European Conference on Computer Vision. pp. 170–188. Springer (2022) 10
5. Chen, T.: On the importance of noise scheduling for diffusion models. arXiv preprint arXiv:2301.10972 (2023) 4
6. Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. Advances in Neural Information Processing Systems 34, 8780–8794 (2021) 2, 3
7. Diffusion, S.: Stable diffusion 2-1 base. https://huggingface.co/stabilityai/stable-diffusion-2-1-base/blob/main/v2-1_512-ema-pruned.ckpt (2022) 1, 2, 10
8. Du, R., Chang, D., Hospedales, T., Song, Y.Z., Ma, Z.: Demofusion: Democratising high-resolution image generation with no \$\$\$\$. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6159–6168 (2024) 4
9. Gu, J., Zhai, S., Zhang, Y., Susskind, J., Jaitly, N.: Matryoshka diffusion models. arXiv preprint arXiv:2310.15111 (2023) 4
10. Gu, S., Chen, D., Bao, J., Wen, F., Zhang, B., Chen, D., Yuan, L., Guo, B.: Vector quantized diffusion model for text-to-image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10696–10706 (2022) 2, 3

11. Haji-Ali, M., Balakrishnan, G., Ordonez, V.: Elasticdiffusion: Training-free arbitrary size image generation through global-local content separation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6603–6612 (2024) [4](#)
12. He, Y., Yang, S., Chen, H., Cun, X., Xia, M., Zhang, Y., Wang, X., He, R., Chen, Q., Shan, Y.: Scalecrafter: Tuning-free higher-resolution visual generation with diffusion models. arXiv preprint arXiv:2310.07702 (2023) [2](#), [4](#), [8](#), [9](#), [10](#), [11](#), [13](#), [14](#)
13. He, Y., Yang, T., Zhang, Y., Shan, Y., Chen, Q.: Latent video diffusion models for high-fidelity video generation with arbitrary lengths. arXiv preprint arXiv:2211.13221 (2022) [2](#), [3](#), [10](#), [11](#)
14. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. *Advances in neural information processing systems* **33**, 6840–6851 (2020) [3](#)
15. Ho, J., Saharia, C., Chan, W., Fleet, D.J., Norouzi, M., Salimans, T.: Cascaded diffusion models for high fidelity image generation. *The Journal of Machine Learning Research* **23**(1), 2249–2281 (2022) [2](#), [4](#)
16. Ho, J., Saharia, C., Chan, W., Fleet, D.J., Norouzi, M., Salimans, T.: Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.* **23**, 47–1 (2022) [4](#), [6](#), [8](#)
17. Hoogeboom, E., Heek, J., Salimans, T.: simple diffusion: End-to-end diffusion for high resolution images. arXiv preprint arXiv:2301.11093 (2023) [4](#)
18. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021) [2](#), [4](#)
19. Jin, Z., Shen, X., Li, B., Xue, X.: Training-free diffusion model adaptation for variable-sized text-to-image synthesis. arXiv preprint arXiv:2306.08645 (2023) [2](#), [4](#), [9](#), [11](#), [13](#)
20. Parmar, G., Zhang, R., Zhu, J.Y.: On aliased resizing and surprising subtleties in gan evaluation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11410–11420 (2022) [10](#)
21. Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., Rombach, R.: Sdxl: improving latent diffusion models for high-resolution image synthesis. arXiv preprint arXiv:2307.01952 (2023) [1](#), [2](#), [3](#), [4](#), [10](#)
22. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021) [5](#)
23. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10684–10695 (2022) [2](#), [3](#)
24. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022) [3](#), [4](#)
25. Saharia, C., Chan, W., Chang, H., Lee, C., Ho, J., Salimans, T., Fleet, D., Norouzi, M.: Palette: Image-to-image diffusion models. In: ACM SIGGRAPH 2022 Conference Proceedings. pp. 1–10 (2022) [2](#), [3](#)
26. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., Schramowski, P., Kundurthy, S., Crowson, K., Schmidt, L., Kaczmarczyk, R., Jitsev, J.: Laion-5b: An open large-scale dataset for training next generation image-text models (2022) [10](#), [11](#)

27. Si, C., Huang, Z., Jiang, Y., Liu, Z.: Freeu: Free lunch in diffusion u-net. arXiv preprint arXiv:2309.11497 (2023) [7](#)
28. Singer, U., Polyak, A., Hayes, T., Yin, X., An, J., Zhang, S., Hu, Q., Yang, H., Ashual, O., Gafni, O., et al.: Make-a-video: Text-to-video generation without text-video data. arXiv preprint arXiv:2209.14792 (2022) [2](#), [3](#)
29. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 (2020) [10](#)
30. Song, Y., Sohl-Dickstein, J., Kingma, D.P., Kumar, A., Ermon, S., Poole, B.: Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456 (2020) [3](#)
31. Su, X., Song, J., Meng, C., Ermon, S.: Dual diffusion implicit bridges for image-to-image translation. arXiv preprint arXiv:2203.08382 (2022) [2](#), [3](#)
32. Teng, J., Zheng, W., Ding, M., Hong, W., Wangni, J., Yang, Z., Tang, J.: Relay diffusion: Unifying diffusion process across resolutions for image synthesis. arXiv preprint arXiv:2309.03350 (2023) [4](#)
33. Unterthiner, T., van Steenkiste, S., Kurach, K., Marinier, R., Michalski, M., Gelly, S.: Towards accurate generative models of video: A new metric & challenges. arXiv preprint arXiv:1812.01717 (2018) [11](#)
34. Unterthiner, T., van Steenkiste, S., Kurach, K., Marinier, R., Michalski, M., Gelly, S.: Towards accurate generative models of video: A new metric & challenges. ICLR (2019) [11](#)
35. Wang, Y., Chen, X., Ma, X., Zhou, S., Huang, Z., Wang, Y., Yang, C., He, Y., Yu, J., Yang, P., et al.: Lavie: High-quality video generation with cascaded latent diffusion models. arXiv preprint arXiv:2309.15103 (2023) [2](#), [6](#)
36. Wu, J.Z., Ge, Y., Wang, X., Lei, S.W., Gu, Y., Shi, Y., Hsu, W., Shan, Y., Qie, X., Shou, M.Z.: Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7623–7633 (2023) [2](#), [3](#)
37. Xie, E., Yao, L., Shi, H., Liu, Z., Zhou, D., Liu, Z., Li, J., Li, Z.: Diffit: Unlocking transferability of large diffusion models via simple parameter-efficient fine-tuning. arXiv preprint arXiv:2304.06648 (2023) [4](#)
38. Yu, S., Sohn, K., Kim, S., Shin, J.: Video probabilistic diffusion models in projected latent space. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18456–18466 (2023) [2](#), [3](#)
39. Zhang, D.J., Wu, J.Z., Liu, J.W., Zhao, R., Ran, L., Gu, Y., Gao, D., Shou, M.Z.: Show-1: Marrying pixel and latent diffusion models for text-to-video generation. arXiv preprint arXiv:2309.15818 (2023) [2](#), [6](#)
40. Zheng, Q., Guo, Y., Deng, J., Han, J., Li, Y., Xu, S., Xu, H.: Any-size-diffusion: Toward efficient text-driven synthesis for any-size hd images. arXiv preprint arXiv:2308.16582 (2023) [2](#), [4](#)