# FiTv2: Scalable and Improved Flexible Vision Transformer for Diffusion Model

ZiDong Wang, Zeyu Lu, Di Huang, Cai Zhou, Wanli Ouyang, and Lei Bai

*Abstract*—*Nature is infinitely resolution-free*. In the context of this reality, existing diffusion models, such as Diffusion Transformers, often face challenges when processing image resolutions outside of their trained domain. To address this limitation, we conceptualize images as sequences of tokens with dynamic sizes, rather than traditional methods that perceive images as fixed-resolution grids. This perspective enables a flexible training strategy that seamlessly accommodates various aspect ratios during both training and inference, thus promoting resolution generalization and eliminating biases introduced by image cropping. On this basis, we present the Flexible Vision Transformer (FiT), a transformer architecture specifically designed for generating images with *unrestricted resolutions and aspect ratios*. We further upgrade the FiT to FiTv2 with several innovative designs, includingthe Query-Key vector normalization, the AdaLN-LoRA module, a rectified flow scheduler, and a Logit-Normal sampler. Enhanced by a meticulously adjusted network structure, FiTv2 exhibits $2\times$ convergence speed of FiT. When incorporating advanced training-free extrapolation techniques, FiTv2 demonstrates remarkable adaptability in both resolution extrapolation and diverse resolution generation. Additionally, our exploration of the scalability of the FiTv2 model reveals that larger models exhibit better computational efficiency. Furthermore, we introduce an efficient post-training strategy to adapt a pre-trained model for the high-resolution generation. Comprehensive experiments demonstrate the exceptional performance of FiTv2 across a broad range of resolutions. We have released all the codes and models at https://github.com/whlzy/FiT to promote the exploration of diffusion transformer models for arbitrary-resolution image generation.

*Index Terms*—Vision Transformer, Diffusion Model.

## I. INTRODUCTION

Natural images inherently possess various resolutions, as illustrated in Fig. 2, the images in *ImageNet* [1] showcase diverse resolutions and aspect ratios. However, current image generation models struggle with generalizing across arbitrary resolutions. The Diffusion Transformer (DiT) [2] family, while excelling within certain resolution ranges, falls short when dealing with images of varying resolutions. This arises from the inability of DiT to incorporate images with dynamic resolutions during its training process, impeding its capability to adapt to diverse token lengths or resolutions effectively.

ZiDong Wang and Zeyu Lu contribute equally to this project.

Zidong Wang, Zeyu Lu, Wanli Ouyang, and Lei Bai are with the Shanghai AI Laboratory, Shanghai, 200000, China.

Zidong Wang and Wanli Ouyang are with the Chinese University of Hong Kong, Shatin, 999077, Hong Kong.

Zeyu Lu is with the Shanghai Jiao Tong University, Shanghai, 200000, China.

Di Huang is with the University of Sydney, Camperdown NSW 2050, Australia.

Cai Zhou is with the Tsinghua University, Beijing, 100084, China.

Corresponding author is Lei Bai. Email: baisanshi@gmail.com.

To bridge this gap, **Flexible Vision Transformer** (FiT) [3] proposes a novel architecture adept at generating images at *unrestricted resolutions and aspect ratios*. The core motivation lies in a fundamental shift in image data conceptualization: FiT conceptualizes images as sequences of variable-length tokens, departing from the traditional perspective of static grids with fixed dimensions. This paradigm shift enables dynamic adjustment of sequence length, facilitating image generation at arbitrary resolutions unconstrained by predefined grids. By efficiently managing and padding these variable-length token sequences to a specified maximum, FiT achieves resolution-independent image synthesis.

FiT represents this paradigm shift through three significant advancements: the flexible training pipeline, network architecture, and inference process. FiT introduces a flexible training pipeline that preserves original image aspect ratios by treating images as token sequences, accommodating varied resolutions within a predefined maximum token limit. This approach, unique among transformer-based generation models, enables adaptive resizing without cropping or disproportionate scaling. Building upon the DiT [2] architecture, FiT incorporates 2-D Rotary Positional Embedding (2-D RoPE) [4], Swish-Gated Linear Units (SwiGLU) [5], and Masked Multi-Head Self-Attention to effectively handle diverse image sizes. For inference, FiT adapts length extrapolation techniques from LLMs, tailoring them for 2-D RoPE to enhance performance across a wide range of resolutions and aspect ratios.

Despite its innovations, FiT exhibits several limitations. It underperforms on the standard ImageNet 256×256 benchmark, and its architecture results in increased parameter count and computational costs compared to DiT. Furthermore, there exists instability issues in the training of FiT, presenting additional challenges for practical implementation.

To achieve better performance, we adopt several advanced enhancements to build FiTv2, an improved version of FiT. Extensive experiments on both class-guided image generation and text-to-image generation tasks demonstrate that our FiTv2 outperforms or achieves competitive performance, compared to other state-of-the-art CNN models [6], [7] and transformer models [2], [8]. Specifically, our FiTv2-3B/2 model, after training only $1000K$ steps on *ImageNet* [1] dataset, achieves competitive performance on standard *ImageNet*-$256 \times 256$ benchmark while outperforming all SOTA models by a significant margin across resolutions of $160 \times 320$, $128 \times 384$, $320 \times 320$, $224 \times 448$, and $160 \times 480$. With merely $200K$ extra post-training steps, our FiTv2-3B/2 model exceeds all SOTA models by a great margin across $512 \times 512$, $320 \times 640$, and $256 \times 768$ resolutions. Moreover, FiTv2-XL/2 model

Fig. 1: **Selected samples from FiTv2-3B/2 models at resolutions of** $256 \times 256$, $512 \times 512$, $768 \times 768$, $256 \times 768$ **and** $768 \times 256$**.** All the images are sampeld with CFG=4.0. FiT is capable of generating images at unrestricted resolutions and aspect ratios. FiTv2 pushes the image generation ability of FiT to a new level, capable of generating better and higher-resolution images.

holistically surpass the DiT [2]-XL/2 and SiT [8]-XL/2 with the same parameters and $28.6\%$ of the training costs. Further, with the same training steps, our FiTv2-XL/2 model surpasses the SiT [8]-XL/2 model greatly on text-to-image tasks.

A preliminary version (i.e., FiT) of this work was published in [3]. In this paper, we extend FiT in the following aspects:

- We propose an improved version of FiT by incorporating Query-Key Vector Normalization (QK-Norm) into the attention layer for stability, as well as decreasing the hidden size of Swish-Gated Linear Unit (SwiGLU) [5] and adopting the Adaptive Layer Normalization with Low-Rank Adaptation [9] (AdaLN-LoRA) for efficiency. These improvements lead to FiTv2, a more efficient and

scalable version of FiT, which achieves state-of-the-art performance in many image generation tasks.

- We improve the training strategy by switching the noise scheduler from denoising diffusion probabilistic model (DDPM) [10] to rectified flow [11] and adopting the Logit-Normal sampling for timesteps, which results in faster convergence. Furthermore, we analyze the limitations of the original FiT and propose a novel mixed data preprocessing strategy that benefits image synthesis across various resolutions. Combining the above architectural and training strategy improvement enables FiTv2 to achieve $2\times$ the convergence speed of the original FiT.
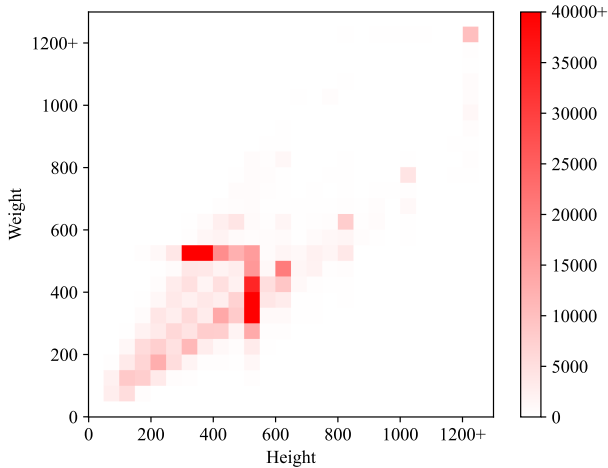- We provide comprehensive analytical experiments and

Fig. 2: The Height/Width distribution of the original *ImageNet* [1] dataset.

visualization to evaluate the effectiveness of FiTv2. We comprehensively analyze the effect of each modification from FiT to FiTv2, as detailed in Sec. V-B. We explore different training-free resolution extrapolation methods for arbitrary resolution generation in FiTv2. Moreover, we benchmark the generalization and extrapolation performance of FiTv2 and other state-of-the-art methods. We also scale our FiTv2 model to 3 billion parameters to study the scalability. Furthermore, we conduct an efficient post-training experiment to investigate the transfer from low resolution to high resolution. To validate the effectiveness of FiTv2 beyond the class-guided image generation, we extend it to text-to-image generation tasks, demonstrating its superiority over the previous state-of-the-art SiT [8] model.

## II. RELATED WORKS

### A. Diffusions and Flows

Denoising diffusion probabilistic models (DDPMs) [10], [12]–[15] and score-based generative models [16], [17] have exhibited remarkable progress in the context of image generation tasks [6], [12], [18]–[21]. The Denoising Diffusion Implicit Model (DDIM) [22], offers an accelerated sampling procedure. Latent Diffusion Models (LDMs) [6] establishes a new benchmark of training deep generative models to reverse a noise process in the latent space, through the use of VAE [23]. Normalizing Flows [24], [25] are a distinct category of generative models which represent data as intricate and complex distributions. Recent flow models [11], [26]–[28] present an alternative approach by learning a neural ordinary differential equation (ODE) that transports between two distributions. By solving a nonlinear least squares optimization problem, rectified flow model [11] learns to map the points drawn from two distributions following the straight paths, which are the shortest paths between two points and hence yield computational efficiency. We follow the rectified flow implementation for image synthesis with fewer sampling steps.

### B. Transformer for Image Generation

The Transformer models [29] have been also explored in the DDPMs [10] and rectified flows [11] to synthesize images. DiT [2] is the seminal work that utilizes a vision transformer as the backbone of LDMs and can serve as a strong baseline. Based on DiT architecture, MDT [30] introduces a masked latent modeling approach, which requires two forward runs in training and inference. U-ViT [31] treats all inputs as tokens and incorporates U-Net architectures into the ViT backbone of LDMs. DiffiT [32] introduces a time-dependent self-attention module into the DiT backbone to adapt to different stages of the diffusion process. Furthermore, SiT [8] utilizes the same architecture as DiT and explores different rectified flow configurations. Efficient-DiT [33] incorporates dynamic mediator tokens into the transformer of SiT and decreases the generation computation. Flag-DiT [34] and SD3 [35] scale up the rectified transformers and achieve better performance. We follow the LDM paradigm of the above methods and further propose a novel flexible image synthesis pipeline.

### C. Long Context Extrapolation

Rotary Position Embedding (RoPE) [4] is a pivotal advancement in positional embedding techniques for large language models (LLM) [36]–[40]. Although RoPE enjoys valuable properties, such as the flexibility of sequence length, its performance drops when the input sequence surpasses the training length. Many training-free approaches have been proposed to solve this issue. Position Interpolation (PI) [41] linearly down-scales the input position indices to match the original context window size, while NTK-aware Scaled RoPE Interoplation [42] changes the rotary base of RoPE based on the Neural Tangent Kernel (NTK) theory. YaRN (Yet another RoPE extensioN) [43] is an improved method to efficiently extend the context window. While these methods scale the the positional embedding to accommodate longer contexts during inference, another paradigm [44], [45] directly scales the attention logits to aggregate information based on entropy theory. Our work provides a comprehensive benchmark for diverse methods on image resolution extrapolation and generalization.

## III. PRELIMINARIES

### A. Rectified Flow

DDPM [10] and score-based generative model [17], both formulated through stochastic differential equations (SDE) [17], produce high-quality samples but suffer from slow inference due to iterative denoising. DDIM [22], an implicit probabilistic model based on ordinary differential equations (ODE), accelerates sampling with fewer steps but at the cost of lower generation quality compared to SDE methods.

To tackle the aforementioned problem, [11] proposes rectified flow, an ODE-based model that transports two empirical distributions $\pi_0$ to $\pi_1$ by following straight line paths as much as possible. The straight paths are both theoretically desired since they are the shortest paths between two endpoints, and computationally efficient because they can be simulated exactly without time discretization, allowing for few-step and even one-step sampling.

Given two target distributions $\pi_0, \pi_1$ and empirical observations $X_0 \sim \pi_0, X_1 \sim \pi_1$, the rectified flow induced from $(X_0, X_1)$ is an ODE model on time $t \in [0, 1]$,

$$\mathrm{d}Z_t = v(Z_t, t)\mathrm{d}t \tag{1}$$

which converts $Z_0$ from $\pi_0$ to $Z_1$ from $\pi_1$. Here the drift force $v : \mathbb{R}^d \to \mathbb{R}^d$ aims to drive the flow to follow the straight direction $(X_1 - X_0)$ as much as possible. To learn this force following the linear path pointing from $X_0$ to $X_1$, a simple least square regression problem needs to be solved:

$$\min_v \int_0^1 \mathbb{E}\Big[||(X_1 - X_0) - v(X_t, t)||^2\Big]\mathrm{d}t \tag{2}$$

where $X_t := tX_1 + (1-t)X_0$ is the linear interpolation of $X_0$ and $X_1$. In practice, $v$ is parameterized with neural networks.

Rectified flow yields several desired properties. First, the flows avoid crossing different paths, which is the condition that the ODE is well-defined, i.e., its solution exists and is unique. While $Z_t$ causalizes, Markovianizes, and derandomizes $X_t$, it preserves the marginal distributions all the time because the continuity equation always holds. Theoretically, rectified flow provably reduces the convex transport costs: $\mathbb{E}[c(Z_1 - Z_0)] \leq \mathbb{E}[c(X_1 - X_0)]$ for any convex $c : \mathbb{R}^d \to \mathbb{R}$. An intuitive explanation is that the paths of the flow $Z_t$ is a rewiring of the straight paths connecting $(X_0, X_1)$, thus the convex transport costs are guaranteed to decrease. Furthermore, on the practical computational efficiency side, the flow becomes nearly straight with just one step of reflow, hence a very few number of Euler discretization steps or even a single Euler step is needed to simulate the ODE. This not only reduces discretization error but also largely improves the sample efficiency.

### B. Rotary Positional Embedding

**1-D RoPE (Rotary Positional Embedding).** 1-D RoPE [4] is a a dominant positional embedding technique for large language models (LLM). By applying a rotary transformation to the embeddings, it incorporates relative position information into absolute positiaonal embedding. Given the $m$-th key and $n$-th query vector as $\mathbf{q}_m, \mathbf{k}_n \in \mathbb{R}^{|D|}$, 1-D RoPE multiplies the bias to the key and query vector in the complex vector space:

$$f_q(\mathbf{q}_m, m) = e^{im\Theta}\mathbf{q}_m, \quad f_k(\mathbf{k}_n, n) = e^{in\Theta}\mathbf{k}_n \tag{3}$$

where $\Theta = \mathrm{Diag}(\theta_1, \cdots, \theta_{|D|/2})$ is rotary frequency matrix with $\theta_d = b^{-2d/|D|}$ and rotary base $b = 10000$. In the real space, given $l = |D|/2$, the rotary matrix $e^{im\Theta}$ equals to:

$$\begin{bmatrix} \cos m\theta_1 & -\sin m\theta_1 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & \cos m\theta_l & -\sin m\theta_l \\ 0 & 0 & \cdots & \sin m\theta_l & \cos m\theta_l \end{bmatrix} \tag{4}$$

The attention score with 1-D RoPE is calculated as:

$$A_n = \mathrm{Re}\langle f_q(\mathbf{q}_m, m), f_k(\mathbf{k}_n, n)\rangle \tag{5}$$

**2-D RoPE.** 2-D RoPE is introduced by our previous work FiT [3] to enhance the resolution generalization in diffusion transformer. Given 2-D coordinates of width and height as $\{(w, h) | 1 \leqslant w \leqslant W, 1 \leqslant h \leqslant H\}$, the 2-D RoPE is:

$$\begin{aligned} f_q(\mathbf{q}_m, h_m, w_m) &= [e^{ih_m\Theta}\mathbf{q}_m \parallel e^{iw_m\Theta}\mathbf{q}_m], \\ f_k(\mathbf{k}_n, h_n, w_n) &= [e^{ih_n\Theta}\mathbf{k}_n \parallel e^{iw_n\Theta}\mathbf{k}_n], \end{aligned} \tag{6}$$

where $\Theta = \mathrm{Diag}(\theta_1, \cdots, \theta_{|D|/4})$, and $\parallel$ denotes concatenating two vectors in the last dimension. Note that we divide the $|D|$-dimension space into $|D|/4$-dimension subspace to ensure the consistency of dimension, which differs from $|D|/2$-dimension subspace in 1-D RoPE. Analogously, the attention score with 2-D RoPE is:

$$A_n = \mathrm{Re}\langle f_q(\mathbf{q}_m, h_m, w_m), f_k(\mathbf{k}_n, h_n, w_n)\rangle. \tag{7}$$

### C. Flexible Vision Transformer Architecture

A prvious version [3] of our work proposes FiT, a transformer architecture that can stably train across various resolutions and generate images with arbitrary resolutions and aspect ratios. Built upon DiT [2], FiT has made some substantial improvements to support flexible training and inference.

Motivated by some significant architectural advances in LLMs [37], [38], [46], FiT replaces the absolute positional embedding with 2-D RoPE and replaces the MLP in Feed-forward Neural Network (FFN) with SwiGLU, further improving the extrapolation capability. Furthermore, FiT uses Masked Multi-Head Self-Attention (MHSA) to replace the standard MHSA in DiT to maintain training integrity with dynamic sequences. Such design enables interaction between noised tokens while isolating padding tokens during the transformer's forward pass. Given sequence mask $M$ where noised tokens are assigned the value of 0, and padding tokens are assigned the value of negative infinity (-inf), masked attention is defined as follows:

$$\mathrm{MaskedAttn}_i = \mathrm{Softmax}\left(\frac{Q_i K_i^T}{\sqrt{d_k}} + M\right)V_i, \tag{8}$$

where $Q_i$, $K_i$, $V_i$ are the query, key, and value matrices for the $i$-th head.

## IV. ENHANCED FLEXIBLE VISION TRANSFORMER

### A. Overview

The overview of FiTv2 is illustrated in Fig. 3. *In the training phase*, FiTv2 encodes preprocessed images into image latents using a pre-trained VAE encoder. These latents are then patchified into sequences of varying lengths $L$. To batch these sequences, the latent tokens are padded to a maximum length $L_{\max} = 256$ with padding tokens. Positional embeddings are similarly padded with zero. The loss function is computed only for the denoised output tokens, ignoring padding tokens.

*In the inference phase*, a position map is defined for the generated image, and noised tokens are sampled from a Gaussian distribution as input. After $K$ iterations of denoising, the tokens are reshaped and unpatchified according to the position map to produce the final image.
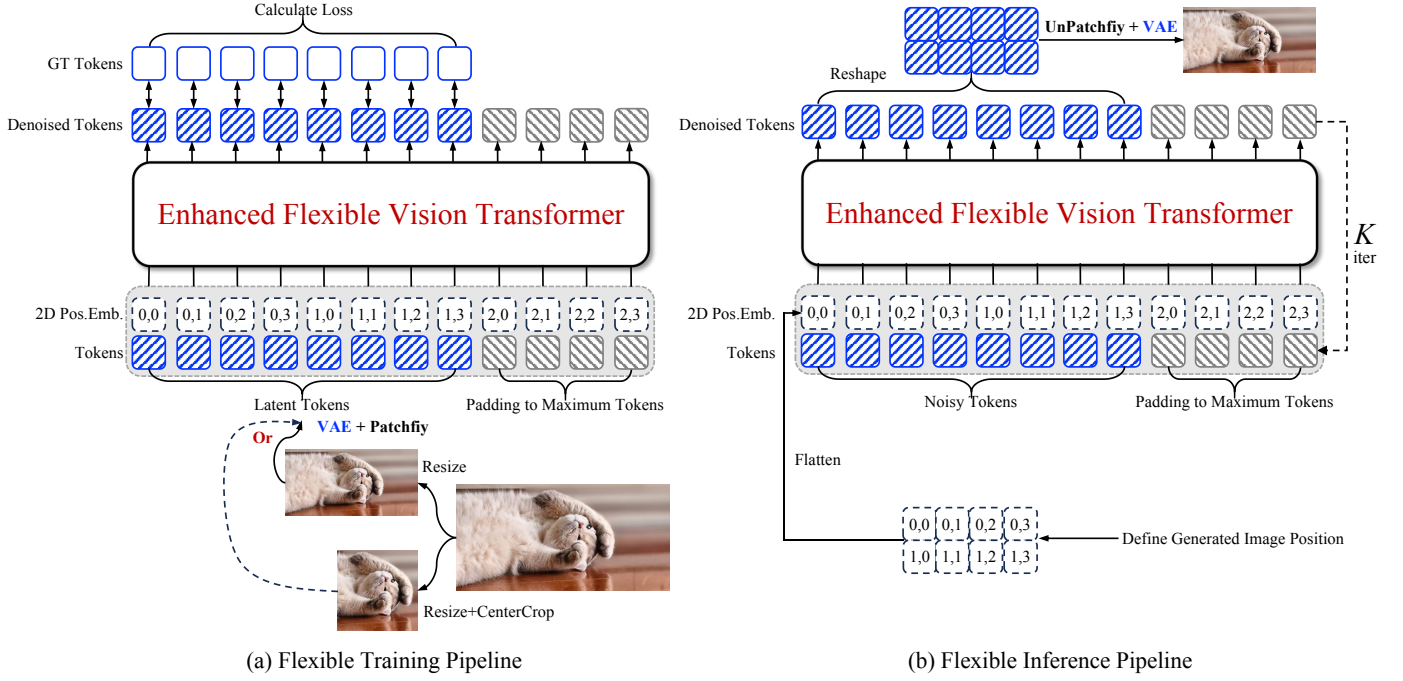
Fig. 3: **Overview of (a) flexible training pipeline, and (b) flexible inference pipeline.** We conceptualize images as dynamic sequences of tokens, allowing for flexible image generation across different resolutions and aspect ratios.
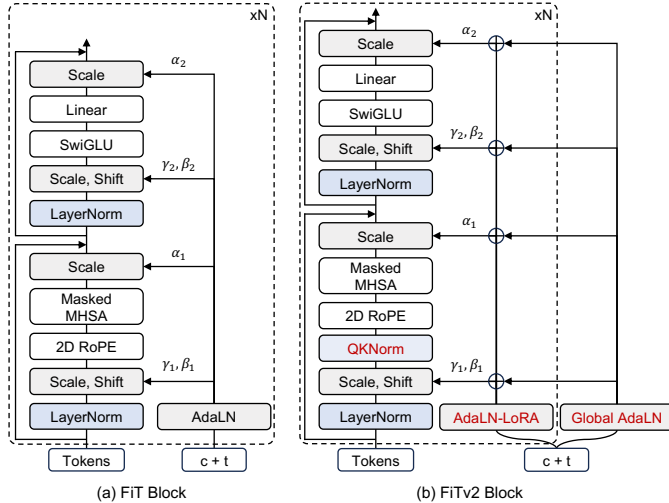


Fig. 4: **Block comparison between (a) FiT and (b) FiTv2.** New modules, QKNorm, AdaLN-LoRA and Global AdaLN, are marked by red color.

### B. Enhancing FiT to FiTv2

We conduct extensive experiments to further improve the design of FiT blocks that enable more stable and efficient training, as detailed in Sec. V-B. The architecture changes from FiT to FiTv2 block are illustrated in Fig. 4.

**Adding QK-Norm to stabilize training.** We observe a vanishing loss problem when scaling up the training steps of the original FiT under mixed-precision training, as in Tab. II. Inspired by the ViT-22B [47], we apply LayerNorm (LN) to the Query (Q) and Key (K) vectors before the attention calculation. Formally, the attention weights in Eq. (8) is

modified to:

$$\text{Softmax}(\frac{1}{\sqrt{d_k}}\text{LN}(Q_i)\text{LN}(K_i)^T + M). \quad (9)$$

By applying this technique, we can effectively eliminate excessively large values in attention logits, which stabilizes the training process, particularly during mixed-precision training. **Reassigning model parameters.** We find that directly using SwiGLU with the same hidden size as the original MLP in DiT [19] will incur more parameters and computational cost, as detailed in Tab. I. To align the parameters and FLOPs with the baseline (the MLP in DiT), the hidden size of SwiGLU in FiTv2 is set to $\frac{2}{3} \times$ of that in the original FiT.

Given the hidden size as $d$, the main parameters of a FiT block are composed of:

$$N = N_{\text{attn}} + N_{\text{swiglu}} + N_{\text{AdaLN}} = 4 \cdot d^2 + 8 \cdot d^2 + 6 \cdot d^2 \quad (10)$$

The parameter ratio of Attention, SwiGLU, and AdaLN module is $2 : 4 : 3$. We argue that too many parameters are occupied by the AdaLN module, which reduces the capacity available for self-attention blocks and potentially affects the scalability of the model. Inspired by W.A.L.T. [48], we adopt AdaLN-LoRA in our FiTv2 block. Additionally, a global AdaLN module is utilized to capture overlapping condition information and reduce the redudancy of condition information of each block. This global AdaLN module is shared by $N$ blocks, as shown in Fig. 4.

Let $S^i = [\beta_1^i, \beta_2^i, \gamma_1^i, \gamma_2^i, \alpha_1^i, \alpha_2^i] \in \mathbb{R}^{6 \times d}$ denote the tuple of all scale and shift parameters, $\mathbf{c} \in \mathbb{R}^d$ and $\mathbf{t} \in \mathbb{R}^d$ represent the embedding for class and time step respectively. For the $i$-th FiTv2 block, we compute the scale and shift parameters:

$$\begin{aligned} S^i &= \text{AdaLN}_{\text{global}}(\mathbf{c} + \mathbf{t}) + \text{AdaLN}_{\text{LoRA}}(\mathbf{c} + \mathbf{t}) \\ &= W^g(\mathbf{c} + \mathbf{t}) + W_2^i W_1^i(\mathbf{c} + \mathbf{t}), \end{aligned} \quad (11)$$
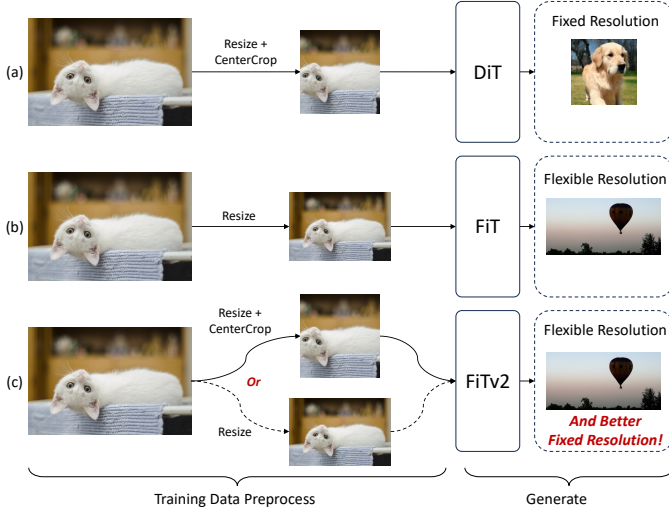
Fig. 5: **Pipeline comparison between (a) DiT, (b) FiT, and (c) FiTv2.** In FiTv2, we incorporate both fixed-resolution images and the flexible-resolution images into training process.

where $W^{\mathrm{g}} \in \mathbb{R}^{(6 \times d) \times d}, W_2^i \in \mathbb{R}^{(6 \times d) \times r}, W_1^i \in \mathbb{R}^{r \times d}$, and the bias parameters are omitted for simplicity. We can adjust the LoRA rank $r$ to change the parameter ratio in FiTv2 blocks. This flexibility allows us to reduce $r$ while simultaneously increasing the number of attention layers $N$, leading to enhanced model performance. In practice, we set $r = \frac{1}{4}d$ (see ablation studies of $r$ in appendices) and the final parameters of a FiTv2 block are composed as:

$$
\begin{aligned}
N &= N_{\mathrm{attn}} + N_{\mathrm{swiglu}} + N_{\mathrm{AdaLN\text{-}LoRA}} \\
&= 4 \cdot d^2 + 8 \cdot d^2 + 1.75 \cdot d^2.
\end{aligned}
\tag{12}
$$

Compared with Eq. (10), we decrease the model parameters occupied by the AdaLN module, enabling us to increase $N$ accordingly while maintaining the model parameters in line with the baseline, as shown in Tab. I.

### C. Improved Training Strategy

**Switching from DDPM to rectified flow.** DDPM [10] is a widely used framework for diffusion models, however, it often exhibits limitations in sampling efficiency. Recently, the rectified flow [11] framework proposes a more flexible manner than DDPM which constructs a transport between two distributions through ordinary differential equations. Unlike DDPM relying on discretized time steps, rectified flow follows straight paths, enabling faster simulation. This elimination of time discretization not only enhances sampling efficiency but also simplifies the overall process. Such inherent advantages have enabled the development of advanced generative models, such as SiT [8] and SD3 [35]. We follow the rectified flow implementation in SiT, linearly connecting the noise and data distributions, and predicting the velocity fields.

**Mixed data preprocessing.** Although the original FiT achieves state-of-the-art performance across unrestricted resolutions and aspect ratios, it underperforms on the standard *ImageNet*-$256 \times 256$ benchmark. We posit that this discrepancy arises from the methodological difference in dataset preparation. As illustrated in Fig. 5 (b), the initial reliance on image

---

**Algorithm 1:** Mixed Data Preprocessing

---
**Input** : Image $I \in \mathbb{R}^{C,H,W}$, target resolution size $S$.
if H > S and W > S:
    if random.random() > 0.5:
        return CenterCrop(Resize($I$))
    else:
        return Resize($I$)
else:
    return Resize($I$)

---

resizing alone of FiT, is opposed to the standard resizing and cropping operations employed in the ADM [7] ImageNet reference dataset used for our FID [49] evaluation.

To bridge this gap, we propose a mixed-data preprocessing strategy, as shown in Fig. 5 (c). To mitigate the blurriness from upscaling low-resolution images, we only crop images whose width and height are both larger than the target resolution size. Exactly, in preprocessing, for images whose sizes are both larger than the target resolution size, we randomly select between resizing only or resizing and cropping with a probability of $\frac{1}{2}$, as in Algorithm 1. For images that do not meet these criteria, we simply resize them to satisfy the sequence length limitation.

As we incorporate resized and cropped images into our training process, we can align the generation distribution of our model with the distribution of the ADM ImageNet reference dataset used for our FID evaluation. Unlike the universal application of resizeing and croping of DiT, as in Fig. 5 (a), our method imposes strict limitations on cropping operations. This strategy, combined with the integration of flexible-resolution images, mitigates the blurring and information loss issues prevalent in previous methods. As a result, this modification enables our FiTv2 to achieve competitive performance on the standard *ImageNet*-$256 \times 256$ benchmark and *ImageNet*-$512 \times 512$ benchmark, while still maintaining the ability to generate images across arbitrary resolutions and aspect ratios.

**Improved sampling strategy.** Typically, the rectified flow scheduler samples timesteps uniformly from the $[0,1]$ interval. Recent studies conducted by SD3 [35] have investigated the choice of timestep sampling strategies and found that the Logit-Normal sampler outperforms the original uniform sampler as well as other variants. Formally, the Logit-Normal sampler is defined as:

$$
u \sim \mathcal{N}(\mathbf{0}, \mathbf{1}), \quad t = \log(\frac{u}{1-u}),
\tag{13}
$$

where $\mathcal{N}(\mathbf{0}, \mathbf{1})$ denotes the standard normal distribution with the mean of $0$ and the standard deviation of $1$. Statistically, this transformation via the logit function ensures that the tails of normal distribution map to the extremes of the $[0,1]$ interval in a way that naturally gives more weight to the central part of the diffusion process. Therefore, the logit-normal sampler facilitates the challenge of learning velocity in the middle of the schedule, as highlighted by EDM [50], and significantly accelerates the model convergence.
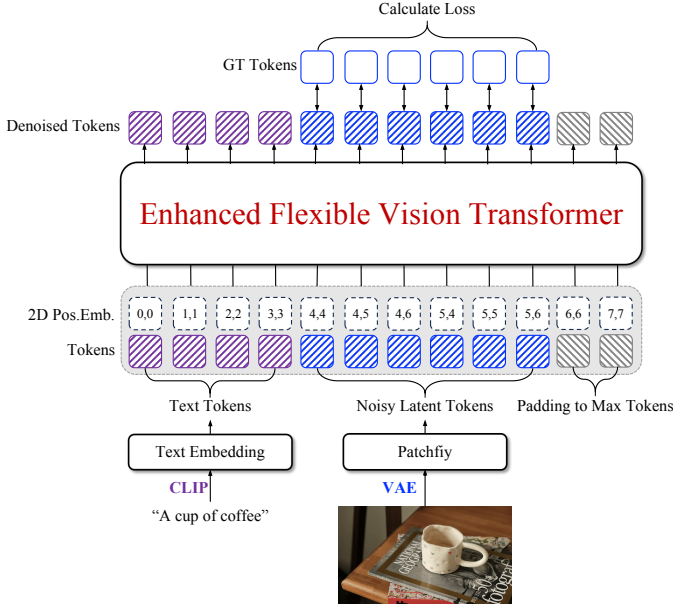
Fig. 6: **Overview of our text-to-image generation model flexible training pipeline.** We utilize CLIP-L to encode text prompts and SD-XL VAE to encode image latents.

### D. High-resolution Post-training

Previous state-of-the-art methods typically train high-resolution models from scratch, thus incurring substantial computational costs. We hypothesize that models trained on low-resolution images have already learned the essential semantic information from the *ImageNet* dataset, but have not been adapted to high-resolution. Therefore, we freeze the majority of parameters of the model and adapt this model through parameter-efficient fine-tuning on the high-resolution data.

Inspired by BitFit [51], our post-training keeps most parameters of the model frozen, only unfreezing specific parameters related to bias and normalization. Considering the increased image resolution, we also unfreeze the parameters of the image patch embedder and the final output layer, leading to only $14.15\%$ of the overall parameters to be trained. Additionally, we apply the NTK Interpolation to the 2-D RoPE embedding to facilitate the transition to higher resolutions.

### E. Text-to-Image Generation

We further evaluate the effectiveness of our FiTv2 model on the text-to-image [35], [52], [53] (T2I) generation task. As illustrated in Fig. 6, we encode an image into image latents with a pre-trained SDXL-VAE [53] encoder and patchify the image latents to latent tokens. We use CLIP-L [54] text encoder to encode the image caption into text features and embed them into text tokens with an MLP. The FiTv2-T2I model processes the concatenated text tokens and noised latent image latents to predict the denoised latent image tokens. The output text tokens and padding tokens are discarded when calculating loss. To accommodate the 1D text tokens with our 2-D RoPE, we convert each single text positional index to a 2D index tuple. Formally, given text tokens $\mathcal{T} \in \mathbb{R}^{M \times D}$

and latent image tokens $\mathcal{I} \in \mathbb{R}^{(H \times W) \times D}$, the text positional indices and image positional indices are defined as follows:

$$P_{\mathcal{T}} = [(0,0), (1,1) \cdots, (M-1, M-1)],$$

$$P_{\mathcal{I}} = \begin{bmatrix} (M, M) & \cdots & (M, M+W-1) \\ \vdots & \ddots & \vdots \\ (M+H-1, M) & \cdots & (M+H-1, M+W-1) \end{bmatrix}.$$
(14)

We also leverage the modulation mechanism of AdaLN module for text conditioning. Specifically, we average-pool the text tokens $\mathcal{T}$ into a semantic text embedding $c_{\mathcal{T}} \in \mathbb{R}^D$, replacing the original class embedding. This pooled text embedding, along with the time embedding, is then used as input for the global AdaLN and the AdaLN-LoRA modules.

### F. Training Free Resolution Extrapolation

**Vision Positional Interpolation.** To achieve resolution extrapolation, we employ training-free positional interpolation techniques, including two widely recognized methods in LLMs: NTK-aware Scaled RoPE Interpolation [42] and YaRN (Yet another RoPE extensioN) Interpolation [43]. Furthermore, we leverage the advanced VisionNTK and VisionYaRN methodologies [3]. These vision-specific adaptations of the original interpolation techniques are tailored to address the unique challenges posed by two-dimensional image data, which are especially effective in generating images with arbitrary aspect ratios. The detailed formulations of these techniques are comprehensively documented in the appendices.

**Attention Scale for Longer Context** In the context of resolution-extrapolation, another approach beyond positional embedding interpolation is scaling the attention logits to aggregate information effectively. Previous studies [44], [45] have theoretically demonstrated that longer contexts result in higher attention entropy of models trained on shorter contexts, leading to widespread aggregation for each token. For higher-resolution image generation, this can cause redundancy in spatial information and disordered object presentations, thereby destroying aesthetics and fidelity. Therefore, a scale factor, which is defined as $s = max(1.0, \sqrt{\log \frac{H_{\text{test}} \times W_{\text{test}}}{H_{\text{train}} \times W_{\text{train}}}})$, is introduced to mitigate the entorpy fluctuations. The formulation of scaled attention is as follows:

$$\text{Softmax}(\frac{1}{\sqrt{d_k}} \text{LN}(Q_i)\text{LN}(K_i)^T \cdot s + M). \quad (15)$$

## V. EXPERIMENTS

### A. FiTv2 Implementation

We present the implementation details of FiTv2, including model architecture, training details, and evaluation metrics.
**Model architecture.** The detailed model architecture is shown in Tab. I. For FiT, we follow SiT-B and SiT-XL to set the same layers, hidden size, and attention heads for base model FiT-B and x-large model FiT-XL. For FiTv2, as described in Sec. IV-B, we reassign parameters to increase the model layers, thereby aligning the parameters with those of DiT [2] and SiT [8]. As DiT and SiT reveal stronger synthesis performance when using a smaller patch size, we use a patch size p=2,

| Model | Layers $N$ | Hidden size $d$ | Heads | Params | GFLOPs |
|---|---|---|---|---|---|
| SiT-B | 12 | 768 | 12 | 131M | 21.8 |
| FiT-B | 12 | 768 | 12 | 159M | 29.1 |
| FiTv2-B | 15 | 768 | 12 | 128M | 27.3 |
| SiT-XL | 28 | 1152 | 16 | 675M | 114 |
| FiT-XL | 28 | 1152 | 16 | 824M | 153 |
| FiTv2-XL | 36 | 1152 | 16 | 671M | 147 |
| FiTv2-3B | 40 | 2304 | 24 | 3B | 653 |

TABLE I: **Details of FiTv2 model architecture.** We follow our original FiT to set the base model and XL model for FiTv2. We also scale up our FiTv2 to 3 billion parameters as our largest model.

denoted by FiT-B/2 and FiTv2-B/2. We adopt the same off-the-shelf pre-trained VAE [6] as SiT, which is provided by the Stable Diffusion [6] to encode/decode the image/latent tokens. The VAE encoder has a downsampling ratio of $1/8$ and a feature channel dimension of 4. An image of size $160 \times 320 \times 3$ is encoded into latent codes of size $20 \times 40 \times 4$. The latent codes of size $20 \times 40 \times 4$ are patchified into latent tokens of length $L = 10 \times 20 = 200$.

**Training details.** We train class-conditional latent FiTv2 models under predetermined maximum resolution limitation, i.e., $H \cdot W \leqslant 256^2$ (equivalent to token length $L \leqslant 256$) for pre-training and $H \cdot W \leqslant 512^2$ (equivalent to token length $L \leqslant 1024$) for post-training, on the *ImageNet* [1] dataset. We down-resize the high-resolution images to meet the $HW \leqslant 256^2$ limitation while maintaining the aspect ratio. We follow SiT to use Horizontal Flip Augmentation. For the pre-training process, we employ a linear learning rate warm-up over the first 5000 steps for stability. Subsequently, we use a constant learning rate of $1 \times 10^{-4}$ using AdamW [55], no weight decay, and a batch size of 256, consistent with SiT. To reduce the training costs, all the experiments are conducted using mixed-precision training. Following common practice in the generative models, we adopt an exponential moving average (EMA) of model weights over training with a decay of 0.9999. All results are reported using the EMA model. We retain the same rectified flow hyper-parameters as SiT.

**Evaluation details and metrics.** We evaluate models with some commonly used metrics, *i.e.* Fre'chet Inception Distance (FID) [56], sFID [49], Inception Score (IS) [57], improved Precision and Recall [58]. For fair comparisons, we follow DiT to use the TensorFlow evaluation from ADM [7] to report FID-50K and other results. Images of FiT and DiT are sampled with 250 DDPM sampling steps, while FiTv2 and SiT both use the adaptive-step ODE sampler (i.e., dopri5) to generate images. FID is used as the major metric as it measures both diversity and fidelity, while IS, sFID, Precision, and Recall are reported as secondary metrics. We report the exact CFG scale if used. The ablation evaluation results on the CFG scale of our FiTv2 model are shown in Fig. 7.

**Evaluation resolution.** Following FiT, we conduct evaluation on different aspect ratios, which are $1 : 1$, $1 : 2$, and $1 : 3$. Besides, we divide the assessment into resolution within the training distribution and resolution out of the training distribution. For the resolution in distribution, we mainly use $256 \times 256$ (1:1), $160 \times 320$ (1:2), and $128 \times 384$ (1:3) for



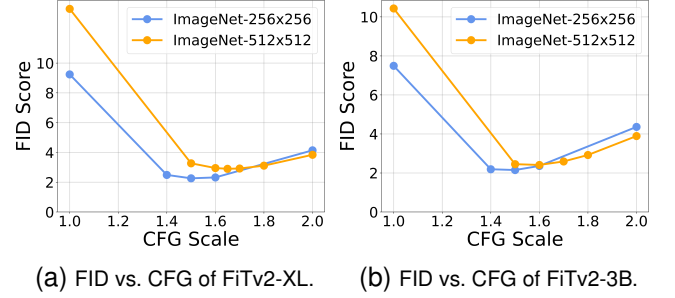(a) FID vs. CFG of FiTv2-XL.  (b) FID vs. CFG of FiTv2-3B.

Fig. 7: **Effect of classifier-free guidance scale on FID score for *ImageNet*-$256 \times 256$ and *ImageNet*-$512 \times 512$ experiments with (a) FiTv2-XL/2 and (b) FiTv2-3B/2 models.** (a) For FiTv2-XL/2 model, the optimal performance is achieved with CFG=1.5 for $256 \times 256$ resolution and CFG=1.65 for $512 \times 512$ resolution. (b) For FiTv2-3B/2 model, the optimal performance is observed with CFG=1.5 for $256 \times 256$ resolution and CFG=1.6 for $512 \times 512$ resolution.

evaluation, with 256, 200, 192 latent tokens respectively. All token lengths are smaller than or equal to 256, leading to respective resolutions within the pre-training distribution. For the resolution out of distribution, we mainly use $320 \times 320$ (1:1), $224 \times 448$ (1:2), and $160 \times 480$ (1:3) for evaluation, with 400, 392, 300 latent tokens respectively. All token lengths are larger than 256, resulting in the resolutions out of pre-training distribution. Through such division, we holistically evaluate the image synthesis and resolution extrapolation ability of FiTv2 at various resolutions and aspect ratios.

### B. From FiT to FiTv2

In this section, we conduct an ablation study to validate the architecture design in FiTv2. We report the results of various variants of FiTv2-B/2, utilizing FID at $256 \times 256$ resolution, and compare these with the DiT-B/2, and SiT-B/2. We train all the models to $2000K$ steps to access the training stability. **Rectified Flow vs. DDPM.** *Rectified Flow scheduler significantly improves the performance and training stability in our FiT model.* Specifically, *Config A* replaces the DDPM scheduler in the original FiT-B/2 with the rectified flow scheduler, leading to substantial performance improvement, both with and without classifier-free guidance (CFG), as in Tab. II. Notably, *Config A* successfully trains to $2000K$ steps, while the training of FiT-B/2 fails after $1500K$ steps, highlighting the stability benefits of the rectified flow scheduler. **QK-Norm vs. No Norm.** *QK-Norm contributes to stabilizing the training process and provides a slight performance enhancement.* We implement LayerNorm for the query and key vectors of attention (*Config B* and *Config E* compared to *Config A* and *Config D*, respectively). As in Tab. II, *Config B* generally achieves better FID scores than *Config A*. Remarkably, we observe that *Config E* maintains a stable training process up to $2000K$ steps, while *Config D* fails to reach this training step. Furthermore, *Config E* outperforms *Config D* at $1500K$ steps in terms of FID score. **Reassigned parameters vs Original parameters.** *Parameter reassignment enhances the efficiency and effectiveness of our*

| Method | Scheduler | QK-Norm | Parameters | Data | Sampling | 256×256 (400k) cfg=1.0 | cfg=1.5 | 256×256 (1000k) cfg=1.0 | cfg=1.5 | 256×256 (1500k) cfg=1.0 | cfg=1.5 | 256×256 (2000k) cfg=1.0 | cfg=1.5 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DiT-B/2 | DDPM | - | - | - | - | 45.33 | 22.21 | 33.27 | 12.59 | ✗ | ✗ | ✗ | ✗ |
| SiT-B/2 | Rectified Flow | - | - | - | - | 36.7 | 16.31 | 27.13 | 9.3 | ✗ | ✗ | ✗ | ✗ |
| FiT-B/2 | DDPM | No | Original | Flexible | Uniform | 36.36 | 18.86 | 29.14 | 11.06 | 26.08 | 9.23 | ✗ | ✗ |
| *Config A* | Rectified Flow | No | Original | Flexible | Uniform | 30.74 | 13.14 | 23.48 | 8.67 | 22.32 | 8.25 | 21.23 | 7.61 |
| *Config B* | Rectified Flow | LayerNorm | Original | Flexible | Uniform | 30.83 | 13.21 | 23.64 | 8.57 | 21.64 | 7.70 | 20.73 | 7.10 |
| *Config C* | Rectified Flow | LayerNorm | Reassigned | Flexible | Uniform | 28.59 | 12.74 | 21.16 | 8.05 | 19.56 | 7.16 | 18.42 | 6.60 |
| *Config D* | Rectified Flow | No | Original | Mixed | Uniform | 34.15 | 13.99 | 25.54 | 8.27 | 23.63 | 7.24 | ✗ | ✗ |
| *Config E* | Rectified Flow | LayerNorm | Original | Mixed | Uniform | 34.55 | 14.19 | 25.94 | 8.37 | 23.45 | 6.99 | 22.04 | 6.31 |
| *Config F* | Rectified Flow | LayerNorm | Original | Mixed | Logit-Normal | 28.49 | 9.98 | 21.93 | 6.16 | 20.09 | 5.23 | 19.21 | 4.84 |
| FiTv2-B/2 | Rectified Flow | LayerNorm | Reassigned | Mixed | Logit-Normal | **26.03** | **9.45** | **19.02** | **5.51** | **17.70** | **4.73** | **16.52** | **4.30** |

TABLE II: **Ablation results from FiT-B/2 to FiTv2-B/2 without using classifier-free guidance.** We train the models to $2000k$ steps to assess stability. A ✗ indicates that the training process breaks down before reaching this evaluation point.

| Method | 320×320 (1:1) FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | 224×448 (1:2) FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | 160×480 (1:3) FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SiT-XL/2 | 19.72 | 54.91 | 144.06 | 0.63 | 0.47 | 46.17 | 67.89 | 73.32 | 0.43 | 0.43 | 104.57 | 91.47 | 23.43 | 0.16 | 0.41 |
| SiT-XL/2 + EI | 8.93 | 19.68 | 212.99 | 0.72 | 0.5 | 78.87 | 48.97 | 43.57 | 0.27 | 0.45 | 131.04 | 71.18 | 17.63 | 0.11 | 0.43 |
| SiT-XL/2 + PI | 8.55 | 20.74 | 217.74 | 0.73 | 0.49 | 82.51 | 50.83 | 41.67 | 0.26 | 0.44 | 133.47 | 72.81 | 17.57 | 0.11 | 0.43 |
| FiTv2-XL/2 | 5.79 | 13.7 | 233.03 | 0.75 | 0.55 | 10.46 | 17.24 | 184.06 | 0.68 | 0.54 | 16.4 | 19.55 | 127.72 | 0.59 | 0.51 |
| FiTv2-XL/2 + PI | 11.47 | 21.131 | 197.04 | 0.67 | 0.51 | 154.59 | 77.21 | 13.18 | 0.10 | 0.14 | 169.4 | 9.81 | 78.31 | 0.06 | 0.06 |
| FiTv2-XL/2 + YaRN | 5.87 | 15.38 | 250.66 | 0.77 | 0.52 | 21.41 | 34.70 | 146.31 | 0.56 | 0.38 | 36.73 | 35.81 | 78.55 | 0.42 | 0.26 |
| FiTv2-XL/2 + NTK | 6.04 | 14.35 | 232.91 | 0.75 | 0.55 | 10.82 | 17.84 | 184.68 | 0.66 | 0.53 | 16.3 | 20.13 | 131.8 | 0.58 | 0.50 |
| FiTv2-XL/2 + VisionYaRN | 5.87 | 15.38 | 250.66 | 0.77 | 0.52 | 6.62 | 18.22 | 245.47 | 0.76 | 0.48 | 16.17 | 27.35 | **151.99** | 0.62 | 0.39 |
| FiTv2-XL/2 + VisionNTK | 6.04 | 14.35 | 232.91 | 0.75 | **0.55** | 10.11 | 17.08 | 188.4 | 0.68 | **0.53** | 15.44 | 19.48 | 135.57 | 0.60 | 0.50 |
| FiTv2-XL/2 + VisionNTK + Attn-Scale | **3.55** | **9.60** | **274.48** | **0.82** | 0.52 | **5.54** | **14.53** | 233.11 | **0.77** | 0.51 | **13.55** | 19.47 | 144.62 | **0.63** | **0.50** |

TABLE III: **Benchmarking class-conditional image generation with out-of-distribution resolution on ImageNet.** The official SiT-XL/2 at $7000k$ training steps and our FiTv2-XL/2 at $2000k$ training steps are adopted in this experiment. Metrics are calculated using classifier-free guidance (cfg=1.5). YaRN and NTK mean the vanilla implementation of such two methods. Our FiTv2-XL/2 demonstrates stable extrapolation performance, which can be significantly improved combined with VisionNTK and attention scale methods.

| Method | Images | Params | 256×256 (1:1) FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | 160×320 (1:2) FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | 128×384 (1:3) FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BigGAN-deep | - | - | 6.95 | 7.36 | 171.4 | 0.87 | 0.28 | - | - | - | - | - | - | - | - | - | - |
| StyleGAN-XL | - | - | 2.30 | 4.02 | 265.12 | 0.78 | 0.53 | - | - | - | - | - | - | - | - | - | - |
| MaskGIT | 355M | - | 6.18 | - | 182.1 | 0.80 | 0.51 | - | - | - | - | - | - | - | - | - | - |
| CDM | - | - | 4.88 | - | 158.71 | - | - | - | - | - | - | - | - | - | - | - | - |
| Large-DiT-7B | 256M | 7.3B | 6.09 | 5.59 | 153.32 | 0.70 | 0.59 | - | - | - | - | - | - | - | - | - | - |
| Efficient-DiT-G (cfg=1.5) | - | 675M | 2.01 | *4.49* | 271.04 | 0.82 | 0.60 | - | - | - | - | - | - | - | - | - | - |
| MaskDiT-G | 2048M | - | 2.28 | 5.67 | 276.56 | 0.80 | **0.61** | - | - | - | - | - | - | - | - | - | - |
| SimpleDiffusion-G (cfg=1.1) | 1024M | 2B | 2.44 | - | 256.3 | - | - | - | - | - | - | - | - | - | - | - | - |
| Flag-DiT-3B-G* | 256M | 4.23B | 1.96 | **4.43** | 284.8 | 0.82 | **0.61** | - | - | - | - | - | - | - | - | - | - |
| Large-DiT-3B-G* | 435M | 4.23B | 2.10 | 4.52 | **304.36** | 0.82 | 0.60 | 118.98 | 62.00 | 12.24 | 0.14 | 0.28 | 142.76 | 80.62 | 10.74 | 0.075 | 0.26 |
| U-ViT-H/2-G (cfg=1.4) | 512M | 501M | 2.35 | 5.68 | 265.02 | 0.82 | 0.57 | 6.93 | 12.64 | 175.08 | 0.67 | 0.63 | 196.84 | 95.90 | 7.54 | 0.06 | 0.27 |
| ADM-G,U | 507M | 673M | 3.94 | 6.14 | 215.84 | 0.83 | 0.53 | 10.26 | 12.28 | 126.99 | 0.67 | 0.59 | 56.52 | 43.21 | 32.19 | 0.30 | 0.50 |
| LDM-4-G (cfg=1.5) | 214M | 395M | 3.60 | 5.12 | 247.67 | **0.87** | 0.48 | 10.04 | 11.47 | 119.56 | 0.65 | 0.61 | 29.67 | 26.33 | 57.71 | 0.44 | **0.61** |
| MDT-G† (cfg=3.8,s=4) | 1664M | 676M | **1.79** | 4.57 | 283.01 | 0.81 | 0.61 | 135.6 | 73.08 | 9.35 | 0.15 | 0.20 | 124.9 | 70.69 | 13.38 | 0.13 | 0.42 |
| DiT-XL/2-G (cfg=1.5) | 1792M | 675M | 2.27 | 4.60 | 278.24 | 0.83 | 0.57 | 20.14 | 30.50 | 97.28 | 0.49 | **0.67** | 107.2 | 68.89 | 15.48 | 0.12 | 0.52 |
| SiT-XL/2-G (cfg=1.5) | 1792M | 675M | 2.15 | 4.50 | 258.09 | 0.81 | 0.60 | 17.38 | 28.59 | 110.32 | 0.52 | 0.65 | 87.40 | 57.41 | 23.45 | 0.16 | 0.56 |
| FiT-XL/2-G (cfg=1.5) | 512M | 824M | 4.21 | 10.01 | 254.87 | **0.84** | 0.51 | 5.48 | 9.95 | 192.93 | **0.74** | 0.56 | 16.59 | **20.81** | 111.59 | 0.57 | 0.52 |
| FiTv2-XL/2-G (cfg=1.5) | 512M | 671M | 2.26 | 4.53 | 260.95 | 0.81 | 0.59 | 5.50 | 11.42 | 211.26 | 0.74 | 0.55 | 14.46 | 23.20 | 135.31 | 0.60 | 0.47 |
| FiTv2-3B/2-G (cfg=1.5) | 256M | 3B | 2.15 | 4.49 | 276.32 | 0.82 | 0.59 | 6.72 | 13.13 | **233.31** | **0.76** | 0.50 | **13.73** | 23.26 | **145.38** | **0.61** | 0.48 |

TABLE IV: **Benchmarking class-conditional image generation with in-distribution resolution on *ImageNet* dataset.** "-G" denotes the results with classifier-free guidance. *: Flag-DiT-3B and Large-DiT-3B actually have 4.23 billion parameters, where 3B means the parameters of all transformer blocks. †: MDT-G adopts an improved classifier-free guidance strategy: $w_t = (1 - \cos \pi (\frac{t}{t_{max}})^s)w/2$, where $w = 3.8$ is the maximum guidance scale and $s = 4$ is the controlling factor.

*FiTv2 model.* As detailed in Sec. IV-B, we reassign the parameters in FiTv2 to optimize the architecture, comparing the reassigned parameters (FiTv2-B/2) with the original parameters (FiT-B/2) in Tab. II. *Config C*, which adopts the reassigned parameters, shows consistent FID improvements across all evaluation points compared with *Config B*.

**Mixed training vs. Flexible training.** *Mixed training improves the model performance when using CFG.* As shown in Tab. II, *Config E* employs a mixed training strategy and exhibits FID performance gains at 1000k, 1500k, and 2000k steps compared to *Config B*.

**Logit-Normal sampling vs. Uniform sampling.** *Logit-*

*Normal sampling significantly accelerates the convergence speed, compared with uniform sampling.* As demonstrated in Tab. II, *Config F* obtains better results than *Config E* at all evaluation points, both with and without CFG.

**From FiT to FiTv2.** *FiTv2 demonstrates significant superiority compared with the original FiT, as well as DiT and SiT.* As reported in Tab. II, experiments on DiT and SiT both break down after $1000K$ training steps, revealing the instability of their architecture. In contrast, FiTv2 exhibits superior training stability, as well as achieves an approximately $2\times$ faster convergence speed compared with FiT, DiT, and SiT.

| Method | Images | Params | 320×320 (1:1) | | | | | 224×448 (1:2) | | | | | 160×480 (1:3) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ |
| U-ViT-H/2-G (cfg=1.4) | 512M | 501M | 7.65 | 16.30 | 208.01 | 0.72 | 0.54 | 67.10 | 42.92 | 45.54 | 0.30 | 0.49 | 95.56 | 44.45 | 24.01 | 0.19 | 0.47 |
| ADM-G,U | 507M | 774M | 9.39 | 9.01 | 161.95 | 0.74 | 0.50 | 11.34 | 14.50 | 146.00 | 0.71 | 0.49 | 23.92 | 25.55 | 80.73 | 0.57 | 0.51 |
| LDM-4-G (cfg=1.5) | 214M | 395M | 6.24 | 13.21 | 220.03 | 0.83 | 0.44 | 8.55 | 17.62 | 186.25 | 0.78 | 0.44 | 19.24 | 20.25 | 99.34 | 0.59 | 0.50 |
| DiT-XL/2-G (cfg=1.5) | 1792M | 675M | 9.98 | 23.57 | 225.72 | 0.73 | 0.48 | 94.94 | 56.06 | 35.75 | 0.23 | 0.46 | 140.2 | 79.50 | 14.70 | 0.09 | 0.45 |
| SiT-XL/2-G (cfg=1.5) | 1792M | 675M | 8.55 | 20.74 | 217.74 | 0.73 | 0.49 | 82.51 | 50.83 | 41.67 | 0.26 | 0.44 | 133.5 | 72.81 | 17.57 | 0.11 | 0.43 |
| FiT-XL/2-G (cfg=1.5) | 512M | 824M | 5.11 | 13.32 | 256.15 | 0.81 | 0.47 | 7.60 | 17.15 | 218.74 | 0.74 | 0.47 | 15.20 | 20.96 | 135.17 | 0.62 | 0.48 |
| FiTv2-XL/2-G* (cfg=1.5) | 512M | 671M | 3.55 | 9.60 | 274.48 | 0.82 | 0.55 | 5.54 | 14.53 | 233.11 | 0.77 | 0.51 | 13.55 | 19.47 | 144.62 | 0.63 | 0.50 |
| FiTv2-3B/2-G* (cfg=1.5) | 256M | 3B | 3.22 | 9.96 | 291.13 | 0.83 | 0.53 | 4.87 | 14.47 | 263.27 | 0.80 | 0.49 | 12.15 | 19.47 | 162.24 | 0.65 | 0.48 |

TABLE V: **Benchmarking class-conditional image generation with out-of-distribution resolution on *ImageNet* dataset.** *: FiTv2 adopts VisionNTK and attention scale for resolution extrapolation. Our FiTv2 model achieves state-of-the-art performance across all the resolutions and aspect ratios, demonstrating a strong extrapolation capability.

| Method | Images | Params | 512×512 (1:1) | | | | | 320×640 (1:2) | | | | | 256×768 (1:3) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | FID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ |
| DiM-Huge-G (cfg=1.7) | +26M | 860M | 3.78 | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| DiffusionSSM-XL-G | 302M | 660M | 3.41 | 5.84 | 255.06 | 0.85 | 0.49 | - | - | - | - | - | - | - | - | - | - |
| MaskGiT | 384M | 227M | 7.32 | - | 156.0 | 0.78 | 0.50 | - | - | - | - | - | - | - | - | - | - |
| SimpleDiffusion-G (cfg=1.1) | 1024M | 2B | 3.02 | - | 248.7 | - | - | - | - | - | - | - | - | - | - | - | - |
| DiffiT-G (cfg=1.49) | - | 561M | 2.67 | - | 252.12 | 0.83 | 0.55 | - | - | - | - | - | - | - | - | - | - |
| MaskDiT-G | 1024M | - | 2.50 | 5.10 | 256.27 | 0.83 | 0.56 | - | - | - | - | - | - | - | - | - | - |
| Large-DiT-3B-G (cfg=1.5) | 471M | 4.23B | 2.52 | 5.01 | 303.70 | 0.82 | 0.57 | - | - | - | - | - | - | - | - | - | - |
| U-ViT-H/2-G (cfg=1.4) | 512M | 501M | 4.05 | 6.44 | 263.79 | 0.84 | 0.48 | 9.79 | 14.64 | 188.8 | 0.76 | 0.49 | 146.58 | 78.69 | 12.47 | 0.21 | 0.36 |
| ADM-G,U | 1385M | 774M | 3.85 | 5.86 | 221.72 | 0.84 | 0.53 | 13.31 | 10.67 | 113.69 | 0.73 | 0.64 | 33.35 | 25.04 | 59.23 | 0.61 | 0.62 |
| DiT-XL/2-G (cfg=1.5) | 768M | 675M | 3.04 | 5.02 | 240.82 | 0.84 | 0.54 | 41.25 | 66.83 | 54.84 | 0.54 | 0.59 | 148.25 | 154.39 | 6.64 | 0.13 | 0.36 |
| FiTv2-XL/2-G (cfg=1.65) | +102M | 671M | 2.90 | 5.73 | 263.11 | 0.83 | 0.53 | 4.87 | 10.75 | 228.09 | 0.80 | 0.53 | 18.55 | 21.69 | 126.55 | 0.69 | 0.53 |
| FiTv2-3B/2-G (cfg=1.6) | +51M | 3B | 2.41 | 5.34 | 284.49 | 0.82 | 0.58 | 4.54 | 11.04 | 240.30 | 0.80 | 0.56 | 16.08 | 19.75 | 140.10 | 0.72 | 0.52 |

TABLE VI: **Benchmarking class-conditional image generation with high-resolution image generation on *ImageNet* dataset.** Our FiTv2 can directly generates images with different aspect ratios with stable and state-of-the-art performance.

## C. Resolution Extrapolation Design

In this part, we adopt the official SiT-XL/2 model at $7000K$ training steps and our FiTv2-XL/2 model at $2000K$ training steps to evaluate the extrapolation performance on three out-of-distribution resolutions: $320 \times 320$, $224 \times 448$ and $160 \times 480$. Direct extrapolation does not perform well on larger resolutions outside of training distribution. So we conduct a comprehensive benchmarking analysis focused on higher resolution extrapolation.

**PI and EI.** PI (Position Interpolation) and EI (Embedding Interpolation) are two baseline positional embedding interpolation methods. PI linearly down-scales the inference position coordinates to match the original coordinates. EI resizes the positional embedding with bilinear interpolation. Following ViT [59], EI is used for absolute positional embedding.

**NTK, YaRN, VisionNTK and VisionYaRN.** The implementation of these interpolation techniques strictly follows the implementation in FiT [3].

**Attention Scale.** The attention scale is defined in Eq. (15), we apply this technique combined with the VisonNTK.

**Analysis.** We present in Tab. III that our FiTv2-XL/2 shows stable performance when directly extrapolating to larger resolutions. When combined with PI, the extrapolation performance of FiTv2-XL/2 at all three resolutions decreases. When directly combined with YaRN, the FID score on $320 \times 320$ changes slightly, but the performance on $224 \times 448$ and $168 \times 480$ descends. Our VisionYaRN solves this dilemma and reduces the FID score by **3.84** on $224 \times 448$ compared with YaRN. NTK interpolation method demonstrates stable extrapolation performance but increases the FID score slightly at $320 \times 320$ and $224 \times 448$ resolutions. Our VisionNTK method slightly exceeds the performance of direct extrapolation on $224 \times 448$ and $160 \times 480$ resolutions. When combining

VisionNTK and attention scale, the performance significantly surpasses all the other extrapolation methods, with FID improvement **2.24** on $320 \times 320$, **4.92** on $224 \times 448$ and **2.89** on $160 \times 480$ compared with direct extrapolation. In conclusion, our FiTv2-XL/2 model demonstrates robust extrapolation capabilities. Additionally, VisionYaRN and VisonNTK can enhance the generation performance on varied aspect ratios. Furthermore, the combination of VisionNTK with attention scale greatly improves high-resolution extrapolation ability.

However, the official SiT-XL/2 model demonstrates poor extrapolation ability, in Tab. III. When combined with PI, the FID score achieves **19.72** at $320 \times 320$ resolution, which still falls behind our FiTv2-XL/2. At $224 \times 448$ and $160 \times 480$ resolutions, PI and EI interpolation methods cannot improve the extrapolation performance.

## D. Pre-trained Model Results

*1) In-Distribution Resolution Results:* In this part, we compare our FiTv2 model with other baselines. Our FiTv2-XL model is trained with $2000K$ steps, consuming only $28.6\%$ of the cost of SiT but with better performance. Furthermore, we scale our FiTv2 up to 3B parameters, which is trained with $1000K$ steps. We conduct experiments to evaluate the performance of FiTv2 at three different in-distribution resolutions: $256 \times 256$, $160 \times 320$, and $128 \times 384$. We show samples from the FiTv2 in Fig 1, and we compare against some state-of-the-art class-conditional generative models: BigGAN [60], StyleGAN-XL [61], MaskGIT [62], CDM [63], Large-DiT [34], MaskDiT [64], Efficient-DiT [33], SimpleDiffusion [65], Flag-DiT [34], U-ViT [31], ADM [7], LDM [6], MDT [30], DiT [2], SiT [8] and our original FiT. When generating images of $160 \times 320$ and $128 \times 384$ resolution, we adopt PI on the positional embedding of the DiT and SiT
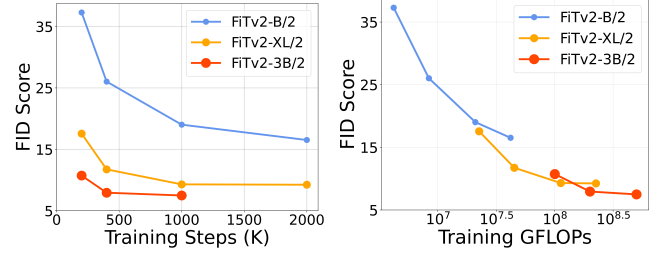
model, as stated in Sec. V-C. EI is employed in the positional embedding of U-ViT and MDT models, as they use learnable positional embedding. ADM and LDM can directly synthesize images with resolutions different from the training resolution. For Large-DiT, we directly generate images of different resolutions as it uses 1D-RoPE as position embedding. For the FiT and FiTv2 models, we directly generate images with different aspect ratios without any extrapolation techniques.

As shown in Tab. IV, FiTv2-XL/2 and FiTv2-3B/2 outperform all prior diffusion models, demonstrating exceptional performance on both the standard $256 \times 256$ benchmark and varied resolutions. FiTv2-XL/2 reduces the FID by 1.95 compared to the original FiT-XL/2 with the same training steps and a smaller model size. Our FiTv2-XL/2 and FiTv2-3B/2 can be competitive with any other SOTA methods on $256 \times 256$ resolution. FiT-XL/2 and FiTv2-XL/2 achieve superior performance on $160 \times 320$ resolution, decreasing the previous best FID of **6.93** achieved by U-ViT-H/2-G to **5.48** and **5.50** respectively. On $128 \times 384$ resolution, FiTv2-XL/2 and FiTv2-3B/2 show significant superiority, decreasing the previous SOTA FID-50K of **29.67** achieved by LDM-4/G to **14.46** and **13.73** respectively. In conclusion, these results suggest that our FiTv2 model has improved performance on standard benchmarks while maintaining the enhanced capability to generate images with arbitrary aspect ratios.

*2) Out-of-Distribution Resolution Results:* We evaluate our FiTv2-XL/2 on three different out-of-distribution resolutions: $320 \times 320$, $224 \times 448$, and $160 \times 480$ and compare against some SOTA class-conditional generative models: U-ViT, ADM, LDM-4, MDT, DiT, SiT, and the original FiT. PI is employed in DiT and SiT, while EI is adopted in U-ViT, as in Sec. V-D1. U-Net-based methods, such as ADM and LDM-4 can directly generate images with resolution out of distribution. VisionNTK is adopted in FiT, and we combine VisionNTK and attention scale to our FiTv2 model. Note that we do not evaluate the MDT and Large-DiT, as they fall short of generating images whose resolution differs from the training resolution in Tab. IV.

As shown in Tab. V, FiTv2-XL/2 and FiTv2-3B/2 achieve the best FID-50K, IS, and Precision, on all three resolutions, indicating their outstanding extrapolation ability. In terms of other metrics, such as sFID and Recall, the FiTv2 model demonstrates competitive performance. FiTv2-XL/2 surpasses FiT-XL/2 on all three resolutions with fewer parameters and FLOPs. Compared with the previous SOTA LDM-4, FiTv2-3B/2 gains FID improvement by **3.02**, **3.68** and **7.09** on $320 \times 320$, $224 \times 448$ and $160 \times 480$ resolutions, respectively.

*3) Analysis of the Pretraining Results:* **Scalibility analysis.** In Fig. 8a, we demonstrate how model performance changes as training steps increase. In Fig. 8b, we present the relation of model performance with the training GFLOPs, which is calculated as GFLOPs $\times$ batch size $\times$ training steps $\times$ 3, following DiT [2]. As the GFLOPs increase, whether by increasing training steps or enlarging the model size, the FID score and aesthetic quality consistently improve. Additionally, we observe that with the same training GFLOPs, the larger FiTv2 model always shows better qualitative and quantitative results. In contrast, the smaller FiTv2 model, even when trained for



(a) FID vs. Training Steps.  (b) FID vs. Training GFLOPs.

Fig. 8: **Effect of scaling FiTv2 model.** All the images are sampled without using CFG. We demonstrate FID over training iterations (a) and training GFLOPs (b) of our FiTv2 model of three sizes. Scaling our FiTv2 model yieds better quantitative and qualitative performance.

more steps, fails to reach the performance of larger FiTv2 models trained for fewer steps. We conclude that scaling model size is a more efficient approach to managing to compute costs, consistent with the findings from DiT.

**Flexibility analysis.** LDMs with transformer backbones are known to have difficulty in generating images out of training resolution, such as DiT, U-ViT, MDT, SiT, and Large-DiT. More seriously, MDT almost has no ability to generate images beyond the training resolution. We speculate this is because both learnable absolute PE and learnable relative PE are used in MDT. Large-Dit also encounters difficulty in generating images with varied resolutions, as the usage of 1D-RoPE makes it hard to encode spatial structure in images. DiT, U-ViT, and SiT show a certain degree of extrapolation ability and achieve FID scores of **9.98**, **7.65** and **8.55** respectively at $320 \times 320$ resolution. However, when the aspect ratio is not equal to one, their generation performance drops significantly, as $128 \times 384$, $224 \times 448$, and $160 \times 480$ resolutions. Benefiting from the advantage of the local receptive field of the CNN, ADM and LDM show stable performance on resolution extrapolation and generalization ability to various aspect ratios. Our FiTv2 model solves the problem of insufficient extrapolation and generalization capabilities of the transformer in image synthesis. At $160 \times 320$, $128 \times 384$, $320 \times 320$, $224 \times 448$, and $160 \times 480$ resolutions, FiTv2-XL/2 exceeds the previous SOTA CNN methods, like ADM and LDM.

### E. High-resolution Post-trained Model Results

We extend the context length to $1024$ (equivalent to $H \cdot W \leqslant 512^2$) to conduct high-resolution post-training. As detailed in Sec. IV-D, we utilize the model pre-trained with the context length $L \leqslant 256$, keeping the major parameters frozen. We only update the parameters associated with bias, normalization, image patch embedder, and the final layer, leading to merely $14.15\%$ of the overall parameters. Training is conducted using a constant learning rate of $1 \times 10^{-4}$ using AdamW, no weight decay, and a batch size of 256, same with the DiT and SiT training setting. Specifically, we train the FiTv2-XL/2 model for $200K$ steps and the FiTv2-3B/2 model for $100K$ steps.

| The picture on the wall is framed by two lamps | An ornate clock with four sides is shown in this photograph | A black cat with a tiny knitted hat on its head | A wedding cake decorated with a tree design | A bowl with rice and steamed broccoli next to a bottle of sauce |
| A view of a train and a city street while it's snowing | A freshly baked pizza that has not been eaten yet | A tray with two sandwich halves, a salad, a cup of coffee and a fork | A yellow train in an outside train station | Two grey teddy bears wearing bright knit hats and sweaters |

Fig. 9: **Selected samples from FiTv2-XL/2 models at resolutions of** $256 \times 256$ **on text-to-image generation tasks.** All the images are sampled with CFG=4.0. With only $400K$ training steps, our model is capable of generating releastic images according to text descriptions.



Fig. 10: Comparision of **FID and CLIP-L score** across different CFG scales for two text-to-image models: FiTv2-XL/2 and SiT-XL/2. FiTv2-XL/2 significantly outperforms SiT-XL/2 in terms of FID score and CLIP-L score.

The model performance is evaluated on three resolutions: $512 \times 512$ (1:1), $320 \times 320$ (1:2), and $256 \times 768$ (1:3), offering a comprehensive assessment of the image synthesis capability. Our FiTv2 is compared with several state-of-the-art baselines, including DiM [66], DiffusionSSm [67], MaskGiT [62], SimpleDiffusion [65], DiffiT [32], MaskDiT [64], Large-DiT [34], U-ViT [31], ADM [7], and DiT [2]. The open-source baseline models are evaluated on $320 \times 640$ and $256 \times 768$ resolutions. Consistent with Sec. V-D1, PI is adopted in DiT while EI is employed in U-ViT. For ADM and our FiTv2, images with different resolutions are directly generated.

As demonstrated in Tab. VI, FiTv2-XL/2 beats DiT-XL/2 on all three resolutions, with comparable parameters and significantly lower training costs. Remarkably, our FiTv2-XL/2 surpasses DiT-XL/2 on the FID score by **36.38** at $320 \times 640$ resolution and by **129.7** at $256 \times 768$ resolution. Furthermore, our FiTv2-3B/2 consistently outperforms all other baseline models on all three resolutions. FiTv2-3B/2 surpasses the previous SOTA Large-DiT-3B and MaskDiT at $512 \times 512$ resolution. At $320 \times 640$ and $256 \times 768$ resolutions, FiTv2-3B/2 shows great superiority, exceeding the previous SOTA U-ViT by **5.25** at $320 \times 640$ resolution on FID and surpassing previous SOTA ADM by **17.27** at $256 \times 768$ resolution.

*F. Text-to-Image Results*

We conduct text-to-image (T2I) generation experiments to further evaluate the effectiveness of our FiTv2 architecture. We use the filtered and recaptioned CC12M [68] subset from PixelProse [69] for training, which comprises 8.6 million high-quality images with descriptive captions. The CLIP-L [54] text encoder is employed to extract text features, resulting in 77 text tokens, each with 768 dimensions. We use the penultimate hidden representation from the CLIP-L text encoder as the text features following Imagen [70]. We use the SDXL-VAE [53] to extract image latents and the training pipeline follows the class-guided image generation methodology. The procedure aligns with the training recipe outlined in Sec. V-E, with our FiTv2-XL/2 model trained for $400K$ steps. Additionally, a baseline SiT-XL/2 model is trained for the same $400K$ steps for comparative analysis. To ensure a fair comparison, SiT-XL/2 processes the text features and image latents in the same manner as our FiTv2, detailed in Sec. IV-E.

We evaluate our FiTv2-XL/2 and SiT-XL/2 for T2I generation on the standard MS-COCO benchmark at $256 \times 256$ resolution. Consistent with previous literature, we randomly sample $30K$ prompts from the MS-COCO validation set and generate images according to those prompts to compute the FID score and CLIP-L score. The Pareto curve is shown in Fig. 10 with classifier-free guidance factor of $[1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 9.0]$. With the same training steps, our FiTv2 achieves stronger results both on FID and CLIP scores, attaining an optimal FID of **27.88** and an optimal CLIP score of **0.2535** at cfg=4.0. In comparison, the SiT model reaches an optimal FID of **40.8** and an optimal CLIP score of **0.2278** at CFG=4.0. Combined with the qualitative results in Fig. 9, it is evident that our FiTv2 model beats the SiT model on T2I architecture.

## VI. CONCLUSION

In this work, we aim to contribute to the ongoing research on flexible generating arbitrary resolutions and aspect ratios. We propose an Enhanced Flexible Vision Transformer (FiTv2) for the diffusion model, a refined transformer architecture with a flexible training pipeline specifically designed for generating images with arbitrary resolutions and aspect ratios. FiTv2 surpass all previous models, whether transformer-based or CNN-based, across various resolutions. With our resolution extrapolation method, VisionNTK, and attention scale, the performance of FiTv2 has been significantly enhanced further. We also scale the FiTv2 to 3 billion to investigate the scalability of our model. Extensive experiments on class-guided image generation, flexible image generation, high-resolution image generation, and text-to-image generation demonstrate the effectiveness of our FiTv2. We hope our work can inspire insights towards designing more powerful diffusion transformer models.

## ACKNOWLEDGMENT

## REFERENCES

[1] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *EEE/CVF Conference on Computer Vision and Pattern Recognition*, 2009. 1, 3, 8

[2] W. Peebles and S. Xie, "Scalable diffusion models with transformers," in *IEEE/CVF International Conference on Computer Vision*, 2023. 1, 2, 3, 4, 7, 10, 11, 12

[3] Z. Lu, Z. Wang, D. Huang, C. Wu, X. Liu, W. Ouyang, and L. Bai, "Fit: Flexible vision transformer for diffusion model," in *International Conference on Machine Learning*, 2024. 1, 2, 4, 7, 10

[4] J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu, "Roformer: Enhanced transformer with rotary position embedding," *Neurocomputing*, 2024. 1, 3, 4

[5] N. Shazeer, "Glu variants improve transformer," *arXiv preprint arXiv:2002.05202*, 2020. 1, 2

[6] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *EEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 1, 3, 8, 10

[7] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in Neural Information Processing Systems*, 2021. 1, 6, 8, 10, 12

[8] N. Ma, M. Goldstein, M. S. Albergo, N. M. Boffi, E. Vanden-Eijnden, and S. Xie, "Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers," *arXiv preprint arXiv:2401.08740*, 2024. 1, 2, 3, 6, 7, 10

[9] E. J. Hu, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen *et al.*, "Lora: Low-rank adaptation of large language models," in *International Conference on Learning Representations*, 2022. 2

[10] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, 2020. 2, 3, 6

[11] X. Liu, C. Gong, and qiang liu, "Flow straight and fast: Learning to generate and transfer data with rectified flow," in *The Eleventh International Conference on Learning Representations*, 2023. [Online]. Available: https://openreview.net/forum?id=XVjTT1nw5z 2, 3, 6

[12] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Advances in Neural Information Processing Systems*, 2022. 3

[13] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*, 2021. 3

[14] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 9, pp. 10 850–10 869, 2023. 3

[15] S. Bond-Taylor, A. Leach, Y. Long, and C. G. Willcocks, "Deep generative modelling: A comparative review of vaes, gans, normalizing flows, energy-based and autoregressive models," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 11, pp. 7327–7347, 2021. 3

[16] A. Hyvärinen and P. Dayan, "Estimation of non-normalized statistical models by score matching." *Journal of Machine Learning Research*, 2005. 3

[17] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *arXiv preprint arXiv:2011.13456*, 2020. 3

[18] F. Ling, Z. Lu, J.-J. Luo, L. Bai, S. K. Behera, D. Jin, B. Pan, H. Jiang, and T. Yamagata, "Diffusion model-based probabilistic downscaling for 180-year east asian climate reconstruction," *npj Climate and Atmospheric Science*, 2024. 3

[19] C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon, "Sdedit: Guided image synthesis and editing with stochastic differential equations," *arXiv preprint arXiv:2108.01073*, 2021. 3, 5

[20] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman, "Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation," in *EEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 3

[21] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall, "Dreamfusion: Text-to-3d using 2d diffusion," *arXiv preprint arXiv:2209.14988*, 2022. 3

[22] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020. 3

[23] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013. 3

[24] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud, "Neural ordinary differential equations," *Advances in neural information processing systems*, vol. 31, 2018. 3

[25] D. Rezende and S. Mohamed, "Variational inference with normalizing flows," in *International conference on machine learning*. PMLR, 2015, pp. 1530–1538. 3

[26] M. S. Albergo and E. Vanden-Eijnden, "Building normalizing flows with stochastic interpolants," *arXiv preprint arXiv:2209.15571*, 2022. 3

[27] Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le, "Flow matching for generative modeling," *arXiv preprint arXiv:2210.02747*, 2022. 3

[28] M. S. Albergo, N. M. Boffi, and E. Vanden-Eijnden, "Stochastic interpolants: A unifying framework for flows and diffusions," *arXiv preprint arXiv:2303.08797*, 2023. 3

[29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems*, 2017. 3

[30] S. Gao, P. Zhou, M.-M. Cheng, and S. Yan, "Masked diffusion transformer is a strong image synthesizer," *arXiv preprint arXiv:2303.14389*, 2023. 3, 10

[31] F. Bao, S. Nie, K. Xue, Y. Cao, C. Li, H. Su, and J. Zhu, "All are worth words: A vit backbone for diffusion models," in *EEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 3, 10, 12

[32] A. Hatamizadeh, J. Song, G. Liu, J. Kautz, and A. Vahdat, "Diffit: Diffusion vision transformers for image generation," *arXiv preprint arXiv:2312.02139*, 2023. 3, 12

[33] Y. Pu, Z. Xia, J. Guo, D. Han, Q. Li, D. Li, Y. Yuan, J. Li, Y. Han, S. Song *et al.*, "Efficient diffusion transformer with step-wise dynamic attention mediators," *arXiv preprint arXiv:2408.05710*, 2024. 3, 10

[34] P. Gao, L. Zhuo, Z. Lin, C. Liu, J. Chen, R. Du, E. Xie, X. Luo, L. Qiu, Y. Zhang *et al.*, "Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers," *arXiv preprint arXiv:2405.05945*, 2024. 3, 10, 12

[35] P. Esser, S. Kulal, A. Blattmann, R. Entezari, J. Müller, H. Saini, Y. Levi, D. Lorenz, A. Sauer, F. Boesel *et al.*, "Scaling rectified flow transformers for high-resolution image synthesis," in *Forty-first International Conference on Machine Learning*, 2024. 3, 6, 7

[36] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, and P. B. et al, "Palm: Scaling language modeling with pathways," *Journal of Machine Learning Research*, 2023. 3

[37] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, and B. R. et al, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023a. 3, 4

[38] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, and N. B. et al, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023b. 3, 4

[39] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang *et al.*, "Qwen technical report," *arXiv preprint arXiv:2309.16609*, 2023. 3

[40] A. Yang, B. Yang, B. Hui, B. Zheng, B. Yu, C. Zhou, C. Li, C. Li, D. Liu, F. Huang *et al.*, "Qwen2 technical report," *arXiv preprint arXiv:2407.10671*, 2024. 3

[41] S. Chen, S. Wong, L. Chen, and Y. Tian, "Extending context window of large language models via positional interpolation," *arXiv preprint arXiv:2306.15595*, 2023. 3

[42] LocalLLaMA, "Ntk-aware scaled rope allows llama models to have extended (8k+) context size without any fine-tuning and minimal perplexity degradation," https://www.reddit.com/r/LocalLLaMA/comments/14lz7j5/ntkaware_scaled_rope_allows_llama_models_to_have/, accessed: 2024-2-1. 3, 7

[43] B. Peng, J. Quesnelle, H. Fan, and E. Shippole, "Yarn: Efficient context window extension of large language models," *arXiv preprint arXiv:2309.00071*, 2023. 3, 7

[44] Z. Jin, X. Shen, B. Li, and X. Xue, "Training-free diffusion model adaptation for variable-sized text-to-image synthesis," *Advances in Neural Information Processing Systems*, vol. 36, 2024. 3, 7

[45] J. Su, "Revisiting attention scale operation from the invariance of entropy," https://kexue.fm/archives/8823. 3, 7

[46] C. Wu, Y. Gan, Y. Ge, Z. Lu, J. Wang, Y. Feng, P. Luo, and Y. Shan, "Llama pro: Progressive llama with block expansion," *Association for Computational Linguistics*, 2024. 4

[47] M. Dehghani, J. Djolonga, B. Mustafa, P. Padlewski, J. Heek, J. Gilmer, A. P. Steiner, M. Caron, R. Geirhos, I. Alabdulmohsin *et al.*, "Scaling vision transformers to 22 billion parameters," in *International Conference on Machine Learning*. PMLR, 2023, pp. 7480–7512. 5

[48] A. Gupta, L. Yu, K. Sohn, X. Gu, M. Hahn, L. Fei-Fei, I. Essa, L. Jiang, and J. Lezama, "Photorealistic video generation with diffusion models," *arXiv preprint arXiv:2312.06662*, 2023. 5

[49] C. Nash, J. Menick, S. Dieleman, and P. W. Battaglia, "Generating images with sparse representations," *arXiv preprint arXiv:2103.03841*, 2021. 6, 8

[50] T. Karras, M. Aittala, T. Aila, and S. Laine, "Elucidating the design space of diffusion-based generative models," *Advances in neural information processing systems*, vol. 35, pp. 26565–26577, 2022. 6

[51] E. B. Zaken, S. Ravfogel, and Y. Goldberg, "Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models," *arXiv preprint arXiv:2106.10199*, 2021. 7

[52] G. Sun, W. Liang, J. Dong, J. Li, Z. Ding, and Y. Cong, "Create your world: Lifelong text-to-image diffusion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 7

[53] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, "Sdxl: Improving latent diffusion models for high-resolution image synthesis," *arXiv preprint arXiv:2307.01952*, 2023. 7, 12

[54] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763. 7, 12

[55] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017. 8

[56] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," *Advances in Neural Information Processing Systems*, 2017. 8

[57] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," *Advances in Neural Information Processing Systems*, 2016. 8

[58] T. Kynkäänniemi, T. Karras, S. Laine, and T. Lehtinen, J.and Aila, "Improved precision and recall metric for assessing generative models," *Advances in Neural Information Processing Systems*, 2019. 8

[59] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020. 10

[60] A. Brock, J. Donahue, and K. Simonyan, "Large scale gan training for high fidelity natural image synthesis," *arXiv preprint arXiv:1809.11096*, 2018. 10

[61] A. Sauer, K. Schwarz, and A. Geiger, "Stylegan-xl: Scaling stylegan to large diverse datasets," in *ACM SIGGRAPH 2022 conference proceedings*, 2022. 10

[62] H. Chang, H. Zhang, L. Jiang, C. Liu, and W. T. Freeman, "Maskgit: Masked generative image transformer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 10, 12

[63] J. Ho, C. Saharia, W. Chan, D. J. Fleet, M. Norouzi, and T. Salimans, "Cascaded diffusion models for high fidelity image generation," *Journal of Machine Learning Research*, 2022. 10

[64] H. Zheng, W. Nie, A. Vahdat, and A. Anandkumar, "Fast training of diffusion models with masked transformers," *Transactions on Machine Learning Research*, 2023. 10, 12

[65] E. Hoogeboom, J. Heek, and T. Salimans, "simple diffusion: End-to-end diffusion for high resolution images," in *International Conference on Machine Learning*. PMLR, 2023, pp. 13213–13232. 10, 12

[66] Y. Teng, Y. Wu, H. Shi, X. Ning, G. Dai, Y. Wang, Z. Li, and X. Liu, "Dim: Diffusion mamba for efficient high-resolution image synthesis," *arXiv preprint arXiv:2405.14224*, 2024. 12

[67] J. N. Yan, J. Gu, and A. M. Rush, "Diffusion models without attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 8239–8249. 12

[68] S. Changpinyo, P. Sharma, N. Ding, and R. Soricut, "Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 3558–3568. 12

[69] V. Singla, K. Yue, S. Paul, R. Shirkavand, M. Jayawardhana, A. Ganjdanesh, H. Huang, A. Bhatele, G. Somepalli, and T. Goldstein, "From pixels to prose: A large dataset of dense image captions," *arXiv preprint arXiv:2406.10328*, 2024. 12

[70] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Advances in neural information processing systems*, vol. 35, pp. 36479–36494, 2022. 12

[71] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European Conference on Computer Vision*, 2014.

[72] T. Brooks, B. Peebles, C. Holmes, W. DePue, Y. Guo, L. Jing, D. Schnurr, J. Taylor, T. Luhman, E. Luhman, C. Ng, R. Wang, and A. Ramesh, "Video generation models as world simulators," https://openai.com/sora, 2024, accessed: 2024-5-1.

[73] Z. Wang, Z. Lu, D. Huang, T. He, X. Liu, W. Ouyang, and L. Bai, "Predbench: Benchmarking spatio-temporal prediction across diverse disciplines," *arXiv preprint arXiv:2407.08418*, 2024.

[74] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Generative image inpainting with contextual attention," in *EEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.

[75] Z. Lu, J. Jiang, J. Huang, G. Wu, and X. Liu, "Glama: Joint spatial and frequency loss for general image inpainting," in *EEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.