# Scalable High-Resolution Pixel-Space Image Synthesis with Hourglass Diffusion Transformers

**Katherine Crowson** [* 1]   **Stefan Andreas Baumann** [* 2]   **Alex Birch** [* 3]   **Tanishq Mathew Abraham** [1]
**Daniel Z. Kaplan** [4]   **Enrico Shippole** [5]

Figure 1: Samples generated directly in RGB pixel space using our ⧗ HDiT models trained on FFHQ-$1024^2$ and ImageNet-$256^2$.
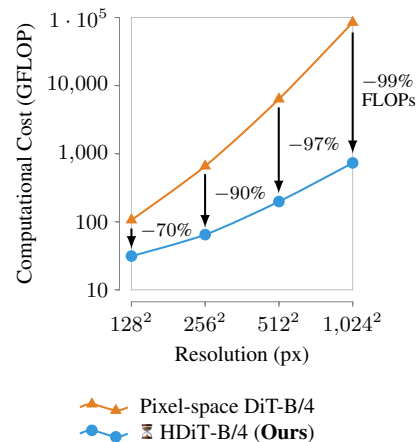


Figure 2: Scaling of computational cost w.r.t. target resolution of our ⧗ HDiT-B/4 model vs. DiT-B/4 (Peebles & Xie, 2023a), both in pixel space. At megapixel resolutions, our model incurs less than 1% of the computational cost compared to the standard diffusion transformer DiT at a comparable size.

## Abstract

We present the Hourglass Diffusion Transformer (HDiT), an image generative model that exhibits linear scaling with pixel count, supporting training at high-resolution (e.g. $1024 \times 1024$) directly in pixel-space. Building on the Transformer architecture, which is known to scale to billions of parameters, it bridges the gap between the efficiency of convolutional U-Nets and the scalability of Transformers.

HDiT trains successfully without typical high-resolution training techniques such as multiscale architectures, latent autoencoders or self-conditioning. We demonstrate that HDiT performs competitively with existing models on ImageNet $256^2$, and sets a new state-of-the-art for diffusion models on FFHQ-$1024^2$.

Code and additional results are available on the project page: crowsonkb.github.io/hourglass-diffusion-transformers.

---
[*]Equal contribution  [1]Stability AI [2]LMU Munich [3]Birchlabs [4]Independent Researcher [5]Independent Researcher. Correspondence to: Katherine Crowson <crowsonkb@gmail.com>, Stefan Baumann <stefan.baumann@lmu.de>, Alex Birch <alex@birchlabs.co.uk>.

1

# 1. Introduction

Diffusion models have emerged as the pre-eminent method for image generation, as evidenced by state-of-the-art approaches like Stable Diffusion (Rombach et al., 2022), Imagen (Saharia et al., 2022), eDiff-I (Balaji et al., 2023), or Dall-E 2 (Ramesh et al., 2022). Their success extends beyond static images to various modalities like video and audio (Blattmann et al., 2023; Kong et al., 2021), showcasing the versatility of diffusion architectures. This recent success can be attributed to their scalability, stability in training, and the diversity of generated samples.

Within the space of diffusion models, there is a large amount of variation in the backbone architectures used, spanning CNN-based (Ho et al., 2020), transformer-based (Peebles & Xie, 2023a; Bao et al., 2023a), CNN-transformer-hybrid (Hoogeboom et al., 2023), or even state-space models (Yan et al., 2023). There is likewise variation in the approaches used to scale these models to support high-resolution image synthesis. Current approaches add complexity to training, necessitate additional models, or sacrifice quality.

Latent diffusion (Rombach et al., 2022) reigns as the dominant method for achieving high-resolution image synthesis. In practice, it fails to represent fine detail (Dai et al., 2023), impacting sample quality and limiting its utility in applications such as image editing. Other approaches to high-resolution synthesis include cascaded super-resolution (Saharia et al., 2022), multi-scale losses (Hoogeboom et al., 2023), the addition of inputs and outputs at multiple resolutions (Gu et al., 2023), or the utilization of self-conditioning and the adaptation of fundamentally new architecture schemes (Jabri et al., 2023).

Our work tackles high-resolution synthesis via backbone improvements. We introduce a pure transformer architecture inspired by the hierarchical structure introduced in (Nawrot et al., 2022), which we call the Hourglass Diffusion Transformer (⧗ HDiT). By introducing a range of architectural improvements, we obtain a backbone that is capable of high-quality image generation at megapixel scale in standard diffusion setups. This architecture, even at low spatial resolutions such as $128 \times 128$ is substantially more efficient than common diffusion transformer backbones such as DiT (Peebles & Xie, 2023a) (see Figure 2) while being competitive in generation quality. Using our method for adapting the model architecture to different target resolutions, we obtain $\mathcal{O}(n)$ computational complexity scaling with the target number of image tokens $n$ in place of the $\mathcal{O}(n^2)$ scaling of normal diffusion transformer architectures, making this the first transformer-based diffusion backbone architecture that is competitive in computational complexity with convolutional U-Nets for pixel-space high-resolution image synthesis.

Our main contributions are as follows:

- We investigate how to adapt transformer-based diffusion backbones for efficient, high-quality pixel-space image generation

- We introduce the Hourglass Diffusion Transformer (⧗ HDiT) architecture for high-resolution pixel-space image generation with subquadratic scaling of compute cost with resolution

- We demonstrate that this architecture scales to high-quality direct pixel-space generation at resolutions of $1024 \times 1024$ without requiring high-resolution-specific training tricks such as progressive growing or multi-scale losses while still being competitive with previous transformer-based architectures at lower resolutions

# 2. Related Work

## 2.1. Transformers

Transformers (Vaswani et al., 2017) reign as the state-of-the-art architectures in various domains (OpenAI, 2023; Zong et al., 2022; Zhang et al., 2022b; Yu et al., 2022; Piergiovanni et al., 2023). Notably, they offer great scalability, up to tens of billions of parameters in the vision space, (Dehghani et al., 2023) and beyond that in other domains such as natural language processing (Chowdhery et al., 2023; Fedus et al., 2022). Transformers consider interactions between all elements in the sequence via the attention mechanism. This enables them to learn long-range interactions efficiently but has the downside of causing their computational complexity to scale quadratically with the length of the input sequence.

**Transformer-based Diffusion Models** Recent works applied transformers to diffusion models, both for generating low-dimensional embeddings as part of a diffusion prior (Ramesh et al., 2022) and for generating compressed image latents (Peebles & Xie, 2023a; Bao et al., 2023a; Zheng et al., 2023; Gao et al., 2023; Bao et al., 2023b; Chen et al., 2023a;b) in a latent diffusion setup (Rombach et al., 2022), leading to state-of-the-art performance. Other works (Hoogeboom et al., 2023; Jing et al., 2023) also applied transformer-based architectures at the lowest level of a U-Net (Ronneberger et al., 2015), or hybridized the two architectures (Cao et al., 2022), going beyond the common practice of putting self-attention blocks into the lower levels of diffusion U-Nets (Ho et al., 2020). However, most transformer architectures for diffusion models are applied with latent diffusion and not directly in pixel space as the quadratic computational complexity of the attention mechanism makes it difficult to apply diffusion transformers for high-resolution pixel-space image synthesis, as found in (Yang et al., 2022).
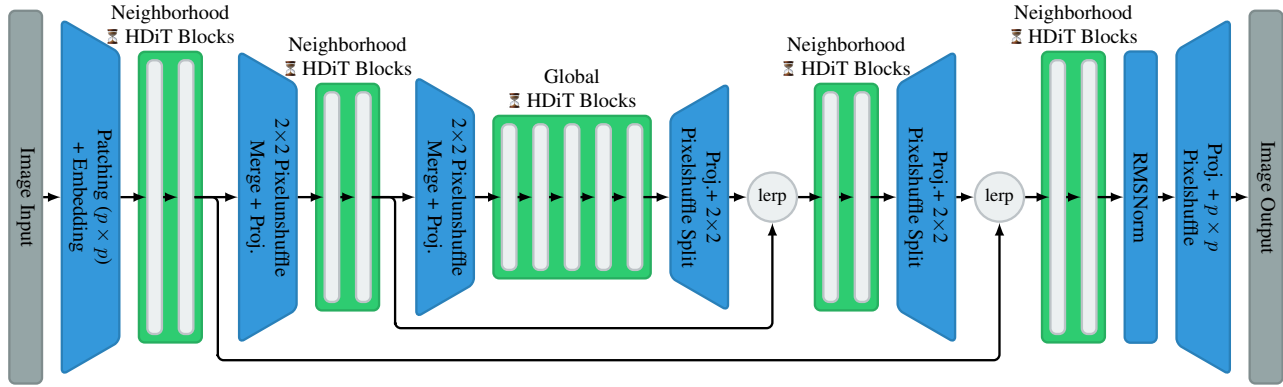
Figure 3: High-level overview of our ⌛ HDiT architecture, specifically the version for ImageNet at input resolutions of $256^2$ at patch size $p = 4$, which has three levels. For any doubling in target resolution, another neighborhood attention block is added. "lerp" denotes a linear interpolation with learnable interpolation weight. All ⌛ HDiT blocks have the noise level and the conditioning (embedded jointly using a mapping network) as additional inputs.

Based on the Diffusion Transformers (DiT) architecture (Peebles & Xie, 2023a), two works (Gao et al., 2023; Zheng et al., 2023) also explored changing the diffusion training process, adding a masking operation to it to incentivize the model to learn better relations between object parts. We consider these additional changes to be orthogonal to the goals pursued in this work.

**Transformer Improvements** As self-attention's computational complexity scales quadratically with the sequence length, many works (Liu et al., 2021; 2022a; Hassani et al., 2023) explored only applying attention to a local set of tokens in vision transformers, leading to linear computational complexity regarding the number of tokens in these local attention mechanisms, at the cost of reducing the receptive field.

Recently, the typical absolute additive, frequency-based positional embedding has also come under scrutiny, with improvements being proposed that effectively encode relative position instead of absolute position. Rotary position embeddings(Su et al., 2022) is one such example, allowing transformers to flexibly adapt to varying sequence lengths and improving performance.

Despite these developments in improving the transformer architecture, especially ViTs, these modifications have been minimally explored for diffusion transformers.

**Hourglass Transformers** The Hourglass architecture (Nawrot et al., 2022) is a hierarchical implementation of transformers that has been demonstrated to be significantly more efficient for language modeling than standard Transformer models both for training and inference. This is done by, over the course of applying the Transformer's layers, iteratively shortening and then iteratively re-expanding the sequence. Additionally, some skip connections reintroduce higher-resolution information near the expansion steps. Generally, this architecture resembles a U-Net (Ronneberger

et al., 2015) but does not use any convolutional layers. Relatedly, (Wang et al., 2022) also showed great performance of a similar structure on image restoration tasks, which can be considered closely related to the denoising diffusion objective.

## 2.2. High-Resolution Image Synthesis with Diffusion Models

There have been extensive investigations into enabling high-resolution image synthesis with diffusion models, a task they typically struggle with out of the box. The most popular approaches have been separating the generation process into multiple steps by either learning multi-stage diffusion models, where a diffusion model generates an initial low-resolution representation – either a downsampled image (Ho et al., 2021) or a learned spatially downsampled "latent" representation (Rombach et al., 2022) – from which a high-resolution image is then generated by a convolutional decoder (Rombach et al., 2022), another diffusion model (Ho et al., 2021; Li et al., 2022), or other generative models (Betker et al., 2023; Fischer et al., 2023). This approach is also used by the vast majority of transformer-based diffusion models (see Section 2.1). Recent works have also explored high-resolution image synthesis in pixel space to simplify the overall architecture, exploring fundamentally new backbone architectures (Jabri et al., 2023), transforming the image data using a discrete wavelet transform to reduce its spatial dimensions (Hoogeboom et al., 2023), and various modifications to the diffusion (training) process, including self-conditioning across sampling steps (Jabri et al., 2023), multiresolution training (Gu et al., 2023), and multiresolution losses (Hoogeboom et al., 2023). Simpler approaches that use neither multi-stage approaches nor the aforementioned adaptations of the diffusion setup (Song et al., 2021) typically struggle with producing samples that fully utilize the available resolution and are globally coherent.

## 3. Preliminaries

### 3.1. Diffusion Models

Diffusion Models generate data by learning to reverse a diffusion process. This diffusion process is most commonly defined to be a Gaussian noising process. Given a data distribution $p_{\text{data}}(\mathbf{x})$, we define a *forward* noising process with the family of distributions $p(\mathbf{x}_{\sigma_t}; \sigma_t)$ that is obtained by adding i.i.d. Gaussian noise of standard deviation $\sigma_t$ which is provided by a predefined monotonically increasing noise level schedule. Therefore, $\mathbf{x}_{\sigma_t} = \mathbf{x}_0 + \sigma_t \epsilon$ where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. A denoising neural network $D_\theta(\mathbf{x}_{\sigma_t}, \sigma_t)$ is trained to predict $\mathbf{x}_0$ given $\mathbf{x}_{\sigma_t}$. Sampling is done by starting at $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \sigma_{\max}^2 \mathbf{I})$ and sequentially denoising at each of the noise levels before resulting in the sample $\mathbf{x}$. The denoiser neural network is trained with a mean-squared error loss:

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \mathbb{E}_{\epsilon, \sigma_t \sim p(\epsilon, \sigma_t)} \left[ \lambda_{\sigma_t} \| D_\theta(\mathbf{x}_{\sigma_t}, \sigma_t) - \mathbf{x} \|_2^2 \right], \quad (1)$$

where $\lambda_{\sigma_t}$ is a weighting function. Often the denoiser is parameterized as a noise predictor:

$$\epsilon_\theta(\mathbf{x}_{\sigma_t}, \sigma_t) = \frac{\mathbf{x}_{\sigma_t} - D_\theta(\mathbf{x}_{\sigma_t}, \sigma_t)}{\sigma_t}. \quad (2)$$

This enables the formulation of a loss which predicts $\epsilon$:

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \mathbb{E}_{\epsilon, \sigma_t \sim p(\epsilon, \sigma_t)} \left[ \lambda_{\sigma_t} \| \epsilon_\theta(\mathbf{x}_{\sigma_t}, \sigma_t) - \epsilon \|_2^2 \right]. \quad (3)$$

Previous work has connected the diffusion model formulation with score-based generative models by observing that the noise prediction objective is closely related to learning the score via denoising score matching.

**Diffusion Improvements** We describe here notable recent improvements to diffusion practices adopted by our model. In EDM (Karras et al., 2022), several modifications to the diffusion framework were shown to improve performance. Most notably, preconditioning is applied to the input and output of the denoiser neural network such that the input and output magnitudes remain constant over noise levels. Specifically, we rewrite the denoiser neural network as:

$$D_\theta(\mathbf{x}_{\sigma_t}, \sigma_t) = c_{\text{out}}(\sigma_t) F_\theta(c_{\text{in}}(\sigma_t) \mathbf{x}_{\sigma_t}, c_{\text{noise}}(\sigma_t)) + c_{\text{skip}}(\sigma_t) \mathbf{x}_{\sigma_t}. \quad (4)$$

The modulation functions are given in (Karras et al., 2022).

Another recent approach demonstrated in (Hang et al., 2023) adapts the loss weighting at different noise levels based on clamped signal-to-noise ratios (SNR) in order to improve model convergence. In the EDM formulation, the loss weighting used is:

$$w(\sigma) = \frac{\min\{\text{SNR}(\sigma), \gamma\}}{c_{\text{out}}^2(\sigma)}$$
$$= \frac{\min\{\text{SNR}(\sigma), \gamma\} \cdot (\sigma^2 \cdot \sigma_{\text{data}}^2)}{\sigma_{\text{data}}^2 + \sigma^2} \quad (5)$$

Since the Min-SNR loss weighting is applied for $\mathbf{x}_0$-parameterization, the $c_{\text{out}}^{-2}(\sigma)$ factor is incorporated to account for the EDM preconditioner parameterization.

Another improvement has been the adaption of noise schedules for high resolutions. It was previously observed (Hoogeboom et al., 2023) that the commonly used noise schedules that were originally designed for low resolutions (32x32 or 64x64) fail to add enough noise at high resolutions. Therefore, the noise schedules can be shifted and interpolated from a reference low-resolution noise schedule in order to add appropriate noise at higher resolutions.

## 4. Hourglass Diffusion Transformers

Diffusion Transformers (Peebles & Xie, 2023a) and other similar works (see Section 2.1) have demonstrated impressive performance as denoising diffusion autoencoders in latent diffusion (Rombach et al., 2022) setups, surpassing prior works in terms of generative quality (Gao et al., 2023; Zheng et al., 2023). However, their scalability to high resolutions is limited by the fact that the computational complexity increases quadratically ($\mathcal{O}(n^2)$ for images of shape $h \times w \times$ channels, with $n = w \cdot h$), making them prohibitively expensive to both train and run on high-resolution inputs, effectively limiting transformers to spatially compressed latents at sufficiently small dimensions, unless very large patch sizes are used (Cao et al., 2022), which have been found to be detrimental to the quality of generated samples (Peebles & Xie, 2023a).

We propose a new, improved hierarchical architecture based on Diffusion Transformers (Peebles & Xie, 2023a), and Hourglass Transformers (Nawrot et al., 2022) – Hourglass Diffusion Transformers (⏳ HDiT) – that enables high-quality pixel-space image generation and can be efficiently adapted to higher resolutions with a computational complexity scaling of $\mathcal{O}(n)$ instead of $\mathcal{O}(n^2)$. This means that even scaling up these models to direct pixel-space generation at megapixel resolutions becomes viable, which we demonstrate for models at resolutions of up to $1024 \times 1024$ in Section 5.

### 4.1. Leveraging the Hierarchical Nature of Images

Natural images exhibit hierarchies (Saremi & Sejnowski, 2013). This makes mapping the image generation process into a hierarchical model an intuitive choice, which has previously been successfully applied in the U-Net architecture (Ronneberger et al., 2015) commonly used in diffusion models but is not commonly used by diffusion transformers (Peebles & Xie, 2023a; Bao et al., 2023a). To leverage this hierarchical nature of images for our transformer backbone, we apply the hourglass structure (Nawrot et al., 2022), which has been shown to be effective for a range of different

modalities, including images, for the high-level structure of our transformer backbone. Based on the model's primary resolution, we choose the number of levels in the hierarchy, such that the innermost level has $16 \times 16$ tokens. As lower-resolution levels have to process both low-resolution information and information that is relevant for following higher-resolution levels, we choose a larger hidden dimension for them. For every level on the encoder side, we merge $2 \times 2$ tokens into one spatially using PixelUnShuffle (Shi et al., 2016) and do the inverse on the decoder side.

**Skip Merging Mechanism** One important consideration in such architectures is the merging mechanisms of skip connections, as it can influence the final performance significantly (Bao et al., 2023a). While the previous non-hierarchical U-ViT (Bao et al., 2023a) uses a concatenation-based skip implementation, similar to the standard U-Net (Ronneberger et al., 2015), and found this to be significantly better than other options, we find additive skips to perform better for this hierarchical architecture. As the usefulness of the information provided by the skips can differ significantly, especially in very deep hierarchies, we additionally enable the model to learn the relative importance of the skip and the upsampled branch by learning a linear interpolation (lerp) coefficient $f$ between the two for each skip and implement them as

$$\mathbf{x}_{\text{merged}}^{(\text{l. lerp})} = f \cdot \mathbf{x}_{\text{skip}} + (1 - f) \cdot \mathbf{x}_{\text{upsampled}}. \quad (6)$$

### 4.2. Hourglass Diffusion Transformer Block Design



(a) ⧖ HDiT Block Architecture.  (b) DiT Block Architecture.
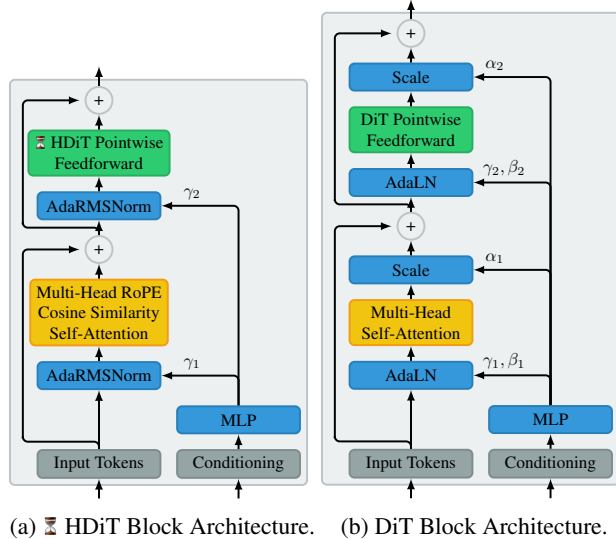
Figure 4: A comparison of our transformer block architecture and that used by DiT (Peebles & Xie, 2023a).

Our basic transformer block design (shown in comparison with that of DiT in Figure 4) is generally inspired by the blocks used by LLaMA (Touvron et al., 2023), a transformer architecture that has recently been shown to be very capable

of high-quality generation of language. To enable conditioning, we make the output scale used by the RMSNorm operations adaptive and have the mapping network, which is conditioned on the class and diffusion time step, predict them. Unlike DiT, we do not employ an (adaptive) output gate, but initialize the output projections of both self-attention and FFN blocks to zeros. To make positional information accessible to the transformer model, common diffusion transformer architectures like DiT and U-ViT use a learnable additive positional encoding. (Peebles & Xie, 2023a; Bao et al., 2023a) As it is known to improve models' generalization and their capability of extrapolating to new sequence lengths, we replace this with an adaptation of rotary positional embeddings (RoPE) (Su et al., 2022) for 2D image data: we follow an approach similar to (Ho et al., 2019) and split the encoding to operate on each axis separately, applying RoPE for each spatial axis to distinct parts of query and key respectively. We also found that applying this encoding scheme to only half of the query and key vectors and not modifying the rest to be beneficial for performance. Overall, we find empirically, that replacing the normal additive positional embedding with our adapted RoPE improves convergence and helps remove patch artifacts. Additionally to applying RoPE, we use a cosine similarity-based attention mechanism that has previously been used in (Liu et al., 2022a)[1]. We note that a similar approach has been proven at the multi-billion parameter scale for vision transformers (Dehghani et al., 2023).

For the feedforward block (see Figure 5 for a comparison with DiT), instead of having an output gate like DiT, we use GEGLU (Shazeer, 2020), where the modulation signal comes from the data itself instead of the conditioning and is applied on the first instead of the second layer of the FFN.
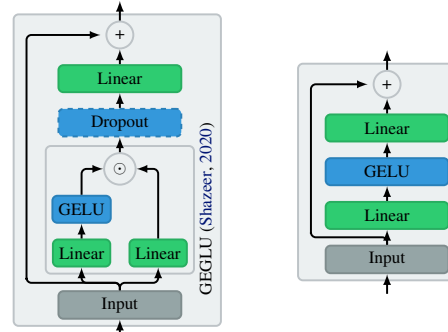


(a) ⧖ HDiT FFN Block.  (b) DiT FFN Block.

Figure 5: A comparison of our pointwise feedforward block architecture and that used by DiT (Peebles & Xie, 2023a).

---

[1]We implement a slight adaptation of their parametrization: instead of parametrizing the per-head scale in logarithmic space, we learn it in linear space, which we find improves stability. See Appendix C for details.

## 4.3. Efficient Scaling to High Resolutions

The hourglass structure enables us to process an image at a variety of resolutions. We use global self-attention at low resolutions to achieve coherence, and local self-attention (Liu et al., 2021; 2022a; Hassani et al., 2023) at all higher resolutions to enhance detail. This limits the need for quadratic-complexity global attention to a manageable amount, and enjoys linear-complexity scaling for any further increase in resolution. Asymptotically, the complexity is $\mathcal{O}(n)$ (see Appendix A) w.r.t pixel count $n$.

A typical choice for localized self-attention would be Shifted Window attention (Liu et al., 2021; 2022a) as used by previous diffusion models (Cao et al., 2022; Li et al., 2022). We find, however, that Neighborhood attention (Hassani et al., 2023) performs significantly better in practice.

The maximum resolution at which to apply global self-attention[2] is a choice determined by dataset (the size at which small features requiring long-distance coherence become large enough for attention to reason about) and by task (the smallest feature whose long-distance relationships need to be preserved in order to be acceptable). At particularly low resolutions (e.g. $256^2$), some datasets permit coherent generation with fewer levels of global attention.

## 5. Experiments

We evaluate the proposed ♟ HDiT architecture on conditional and unconditional image generation, ablating over architectural choices (Section 5.2), and evaluating both megapixel pixel-space image generation (Section 5.3) and large-scale pixel-space image generation (Section 5.4).

## 5.1. Experimental Setup

**Training** Unless mentioned otherwise, we train class-conditional models on ImageNet (Deng et al., 2009) at a resolution of $128 \times 128$ directly on RGB pixels without any kind of latent representation. We train all models with AdamW (Loshchilov & Hutter, 2019) using a constant learning rate of $5 \times 10^{-4}$ and a weight decay of $\lambda = 0.01$. We generally train at a batch size of 256 for 400k steps (following (Peebles & Xie, 2023a)) with stratified diffusion timestep sampling and do not use Dropout unless noted otherwise. For small-scale ImageNet trainings at $128 \times 128$, we do not apply any augmentation. For runs on small datasets, we apply a non-leaking augmentation scheme akin to (Karras et al., 2020a). Following

---

[2]For our FFHQ-$1024^2$ experiment, we apply two levels of global attention – one at $16^2$ and one at $32^2$. Whereas for ImageNet-$128^2$ and $256^2$, we found like prior works (Ho et al., 2020; Hoogeboom et al., 2023; Nichol & Dhariwal, 2021) that a single level of $16^2$ global attention suffices, due to the low resolutions at which images were generated.

common diffusion model training practice and (Peebles & Xie, 2023a), we also compute the exponential moving average (EMA) of the model weights with a decay of 0.9999. We use this EMA version of the model for all evaluations and generated samples, and perform our sampling using 50 steps of DPM++(3M) (Lu et al., 2023; Crowson, 2023) SDE sampling. For further details, see Table 6.

**Diffusion** We adapt our general training setup from (Karras et al., 2022), including their preconditioner, and use a continuous-time diffusion formulation. To enable classifier-free guidance (Ho & Salimans, 2021) during inference, we drop out the class conditioning information $10\%$ of the time during training on class-conditional datasets.

**Evaluation** Following common practice for generative image models, we report the Fréchet Inception Distance (FID) (Heusel et al., 2017) computed on 50k samples. To compute FID, we use the commonly used implementation from (Dhariwal & Nichol, 2021). We also report both the absolute and asymptotic computational complexity for our main ablation study, also including FLOPs for higher-resolution versions of the architecture.

## 5.2. Effect of the Architecture

To evaluate the effect of our architectural choices, we perform an ablation study where we start with a basic implementation of the hourglass architecture for diffusion and iteratively add the changes that enable our final architecture to efficiently perform high-quality megapixel image synthesis. We denote the ablation steps as **A**, **B1**, ..., **E**, and show their feature composition and experimental results in Table 1. We also provide a set of baselines **R1**-**R4**, where we trained DiT (Peebles & Xie, 2023a) models in various settings to enable a fair comparison.

We generally use DiT-B-scale models for this comparison (approx. 130M parameters for DiT, approx 105M to 120M for ♟ HDiT depending on the ablation step), due to their relatively low training cost, and train them on pixel-space ImageNet (Deng et al., 2009) at a resolution of $128^2$ and patch size of 4.

**Baselines** We train 4 versions of DiT in different setups to provide fair comparisons with it as baselines in Table 1. **R1** directly uses the official DiT implementation (Peebles & Xie, 2023b), but omits the VAE latent computation step and adjusts the scaling and variance to fit the data. No other changes were made, as DiT can be directly applied to pixel space (Peebles & Xie, 2023a). To evaluate the influence of our trainer and our loss weighting scheme, we implement a wrapper that directly wraps the original DiT model and

train it with our trainer[3]. The results of this experiment are shown as **R2**. **R3** replaces the wrapped DiT model with a hyperparameter-matched single-level version of ablation step **A**, and matches the performance of the original DiT trained with the original codebase. On top of this setup, we also add soft-min-snr loss weighting to **R4** as in ablation step **E** to enable a fair comparison with our final model. The computational cost for the same architecture at resolutions of $256 \times 256$ and $512 \times 512$ is also reported. In the case of our models, every doubling in resolution involves adding one local attention block (except for ablation step **A**, where it is global) as per Section 4.1.

**Base Hourglass Structure** Configuration **A** is a simple hourglass structure with lower-resolution levels and our linear skip interpolations, and the basic implementation of our blocks with RMSNorm, but without GEGLU, and with full global self-attention at every level. A simple additive positional encoding is used here. Even this simple architecture, without any of our additional changes, is already substantially cheaper (30% of the FLOPs per forward pass) than similarly-sized DiT (Peebles & Xie, 2023a) models operating in pixel space due to the hourglass structure. This comes at the cost of increased FID compared to the DiT baselines at this step in the ablation.

**Local Attention Mechanism** Next, we add local attention to all levels except for the lowest-resolution one. We evaluate two options – Shifted-Window (SWin) (Liu et al., 2021; 2022a) attention (**B1**, a common choice in vision transformers and previously also used in diffusion models (Cao et al., 2022; Li et al., 2022)) and Neighborhood (Hassani et al., 2023) attention (**B2**). Both result in a small reduction in FLOPs even at the low-resolution scale of $128 \times 128$ but, most importantly, reduce the computational complexity w.r.t. the base resolution from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$, enabling practical scaling to significantly higher resolutions. Both variants suffer from increased FID due to this reduced expressiveness of local attention. Still, this change is significantly less pronounced for Neighborhood attention, making it a clearly superior choice in this case compared to the common choice of SWin attention.

**Feedforward Activation** As the third step, we ablate over using GEGLU (Shazeer, 2020), where the data itself affects the modulation of the outputs of the feedforward block, compared to the standard GeLU for the feedforward network. Similar to previous work (Touvron et al., 2023), to account for the effective change of the hidden size due to the GEGLU operation, we decrease the hidden dimension

from $4 \cdot d_{\mathrm{model}}$ to $3 \cdot d_{\mathrm{model}}$. We find that this change significantly improves FID at the cost of a slight increase in computational cost, as the width of the linear projections in the feedforward block has to be increased to account for the halving in output width.

**Positional Encoding** Next, we replace the standard additive positional embedding with our 2d axial adaptation of RoPE (Su et al., 2022) in **D**, completing our Hourglass DiT backbone architecture. This further improves FID. As an additional benefit, RoPE should enable significantly better extrapolation to other resolutions than additive positional embeddings, although our ablation study does not test for that.

**Loss Weighting** Finally, we also ablate over replacing the standard $\frac{1}{\sigma^2}$ loss weighting (Ho et al., 2020; Song et al., 2021) with our adapted min-snr (Hang et al., 2023) loss weighting method that we call soft-min-snr (see Appendix B), which reduces the loss weight compared to SNR weighting for low noise levels. This substantially improves FID further, demonstrating the effectiveness of ⏳ HDiT when coupled with an appropriate training setup for pixel-space diffusion.

**Skip Implementation** Additionally to the main ablation study, we also ablate over different skip implementations based on ablation step **E**. We compare our learnable linear interpolation (lerp), which we empirically found to be especially helpful when training deep hierarchies, with both a standard additive skip, where the upsampled and skip data are directly added, and a concatenation version, where the data is first concatenated and then projected to the original channel count using a pointwise convolution. The results of this ablation are shown in Table 2. We find that, even for shallow hierarchies as used for ImageNet-$128^2$ generation in our ablations, the learnable linear interpolation outperforms the addition slightly, with both the learnable lerp and addition substantially outperforming the commonly used concatenation.

Table 2: Skip Information Merging Mechanism Ablation

| Skip Implementation | FID↓ |
|---|---|
| Concatenation (U-Net (Ronneberger et al., 2015)) | 33.75 |
| Addition (Original Hourglass (Nawrot et al., 2022)) | 28.37 |
| Learnable Linear Interpolation (**Ours**) | **27.74** |

### 5.3. High-Resolution Pixel-Space Image Synthesis

In this section, we train our model for high-resolution pixel-space image synthesis. Following previous works, we train on FFHQ-$1024^2$ (Karras et al., 2021), the standard benchmark dataset for image generation at such high resolutions.

Previous works require tricks such as self-conditioning

---

[3]The pixel-space DiT **R2** was trained with an identical setup to the rest of our ablations except for the optimizer parameters: we initially tried training this model with our optimizer parameters but found it to both be unstable and worse than with the original parameters, so we used the original parameters from (Peebles & Xie, 2023a) for the comparison.

Table 1: Ablation of our architectural choices, starting from a stripped-down implementation of our hourglass diffusion transformer that is similar to DiT-B/4 (Peebles & Xie, 2023a). We also ablate over our additional choice of using soft-min-snr loss weighting, which we use to train our full models but do not consider part of our architecture. We also present results for various DiT-B/4-based models to act as baselines. In addition to training results, we report computational cost per forward pass at multiple resolutions, including standard resolution-dependent model adaptations.

| | Configuration | FID↓ | GFLOP@$128^2$↓ | Complexity↓ | GFLOP@$256^2$ | GFLOP@$512^2$ |
|---|---|---|---|---|---|---|
| **Baselines** (**R1** uses 250 DDPM sampling steps with learned $\sigma(t)$ as in the original publication instead of 50-step DPM++ sampling) | | | | | | |
| **R1** | DiT-B/4 (Peebles & Xie, 2023a) | 42.03 | 106 | $\mathcal{O}(n^2)$ | 657 | 6,341 |
| **R2** | **R1** + our trainer (no soft-min-snr) | 69.86 | 106 | $\mathcal{O}(n^2)$ | 657 | 6,341 |
| **R3** | **R2** + our basic blocks & mapping network | 42.49 | 106 | $\mathcal{O}(n^2)$ | 657 | 6,341 |
| **R4** | **R3** + soft-min-snr | <u>30.71</u> | 106 | $\mathcal{O}(n^2)$ | 657 | 6,341 |
| **Ablation Steps** | | | | | | |
| **A** | Global Attention Diffusion Hourglass (Section 4.1) | 50.76 | 32 | $\mathcal{O}(n^2)$ | 114 | 1,060 |
| **B1** | **A** + Swin Attn. (Liu et al., 2021) | 55.93 | **29** | $\mathcal{O}(n)$ | 60 | 185 |
| **B2** | **A** + Neighborhood Attn. (Hassani et al., 2023) | 51.07 | **29** | $\mathcal{O}(n)$ | 60 | 184 |
| **C** | **B2** + GeGLU (Shazeer, 2020) | 44.36 | <u>31</u> | $\mathcal{O}(n)$ | 65 | 198 |
| **D** | **C** + Axial RoPE (Section 4.2) | 41.41 | <u>31</u> | $\mathcal{O}(n)$ | 65 | 198 |
| **E** | **D** + soft-min-snr (Appendix B) | **27.74** | <u>31</u> | $\mathcal{O}(n)$ | 65 | 198 |

(Jabri et al., 2023), multi-scale model architectures (Gu et al., 2023), or multi-scale losses (Hoogeboom et al., 2023) to enable high-quality generation at such high resolutions. We find that our model does not require such tricks to enable high-quality generation (although we expect them to further increase the quality of generated samples) and, therefore, train our model without them, with the exception of adapting the SNR at each step according to the increase in the images' redundancy (Hoogeboom et al., 2023). As seen in samples from our model in Figure 6, our model can generate high-quality, globally coherent samples that properly utilize the available resolution to produce sharp pictures with fine details, even without classifier-free guidance.



Figure 6: Samples from our 85M-parameter FFHQ-$1024^2$ model. Best viewed zoomed in.

We benchmark our models against state-of-the-at counterparts in Table 3 for a quantitative comparison. Notably, as precomputed metrics for the NCSN++ (Song et al., 2021)

baseline are unavailable, we independently compute them using the provided checkpoint[4]. We find that our model substantially outperforms this baseline both quantitatively and qualitatively (see Figure 10 and Figure 11 for uncurated samples from both our model and the NCSN++ baseline). Notably, our model excels in generating faces with symmetric features, while NCSN++ exhibits noticeable asymmetry. Moreover, ⧖ HDiT effectively leverages the available resolution, producing sharp and finely detailed images, a notable improvement over the NCSN++ model, which often yields blurry samples. We find that our model is competitive regarding FID with high-resolution transformer GANs such as HiT (Zhao et al., 2021) or StyleSwin (Zhang et al., 2022a), but does not reach the same FID as state-of-the-art GANs such as StyleGAN-XL (Sauer et al., 2022). It is worth noting that the FID metric, known for its bias towards samples generated by GANs over those from diffusion models as highlighted in (Stein et al., 2023), underscores the impressive performance of our model, suggesting that the achieved closeness might be approaching the lower limit for this specific metric for diffusion models.

---

[4]Given resource constraints and the prohibitive sampling cost associated with NCSN++ – drawing 50k samples would demand resources equivalent to training our model – we report quantitative metrics for NCSN++ based on 5k samples, and also provide 5k sample-based metrics for ⧖ HDiT.

Table 3: Comparison of our results on FFHQ 1024 × 1024 to other models in the literature. 50k samples are used for FID computation unless specified otherwise.

| Method | FID↓ |
|---|---|
| *Diffusion Models* | |
| NCSN++ (Song et al., 2021) (5k samples) | 53.52 |
| ⚡ HDiT-85M (**Ours**, 5k samples) | 8.48 |
| ⚡ HDiT-85M (**Ours**) | 5.23 |
| *Generative Adversarial Networks* | |
| HiT-B (Zhao et al., 2021) | 6.37 |
| StyleSwin (Zhang et al., 2022a) | 5.07 |
| StyleGAN2 (Karras et al., 2020b) | 2.70 |
| StyleGAN-XL (Sauer et al., 2022) | 2.02 |

### 5.4. Large-Scale ImageNet Image Synthesis

As seen in earlier experiments (see Section 5.3), ⚡ HDiT shows good performance in generating high-fidelity high-resolution samples. To also evaluate its large-scale generation capabilities, we also train a class-conditional pixel-space ImageNet-$256^2$ model. We note that we have not performed any hyperparameter tuning for this task and that this model, at 557M parameters, is significantly smaller than many state-of-the-art models. In alignment with our methodology from high-resolution experiments, we refrain from applying non-standard training tricks or diffusion modifications, and, consistent with (Hoogeboom et al., 2023), we compare results without the application of classifier-free guidance, emphasizing an out-of-the-box comparison.

We show samples in Figure 7 and compare quantitatively with state-of-the-art diffusion models in Table 4. We find that, qualitatively, our model is readily capable of generating high-fidelity samples on this task. Compared to the baseline model DiT, our model achieves a substantially lower FID and higher IS despite operating on pixel-space instead of lower-resolution latents. Compared to other single-stage pixel-space diffusion models, our model outperforms simple U-Net-based models such as ADM but is outperformed by models that use self-conditioning during sampling (RIN) or are substantially larger (simple diffusion, VDM++).

## 6. Conclusion

This work presents ⚡ HDiT, a hierarchical pure transformer backbone for image generation with diffusion models that scales to high resolutions more efficiently than previous transformer-based backbones. Instead of treating images the same regardless of resolution, this architecture adapts to the target resolution, processing local phenomena locally at high resolutions and separately processing global phenomena in low-resolution parts of the hierarchy. This yields an



Figure 7: Samples from our class-conditional 557M-parameter ImageNet-$256^2$ model without classifier-free guidance.

Table 4: Comparison of our results on ImageNet-$256^2$ to other models in the literature. Following (Hoogeboom et al., 2023), we report results without classifier-free guidance. Besides FID@50k and IS@50k, we also report trainable parameter count, samples seen (training iterations times batch size), and sampling steps.

| Method | Params | It.×BS | Steps | FID↓ | IS↑ |
|---|---|---|---|---|---|
| *Latent Diffusion Models* | | | | | |
| LDM-4 (Rombach et al., 2022) | 400M | 214M | 250 | 10.56 | 209.5 |
| DiT-XL/2 (Peebles & Xie, 2023a) | 675M | 1.8B | 250 | 9.62 | 121.5 |
| U-ViT-H/2 (Bao et al., 2023a) | 501M | 512M | 50·2 | 6.58 | - |
| MDT-XL/2 (Gao et al., 2023) | 676M | 1.7B | 250 | 6.23 | 143.0 |
| MaskDiT/2 (Zheng et al., 2023) | 736M | 2B | 40·2 | 5.69 | 178.0 |
| *Single-Stage Pixel-Space Diffusion Models* | | | | | |
| iDDPM (Nichol & Dhariwal, 2021) | - | - | 250 | 32.50 | - |
| ADM (Dhariwal & Nichol, 2021) | 554M | 507M | 1000 | 10.94 | 101.0 |
| RIN (Jabri et al., 2023) | 410M | 614M | 1000 | 4.51 | 161.0 |
| simple diffusion (Hoogeboom et al., 2023) | 2B | 1B | 512 | 2.77 | 211.8 |
| VDM++ (Kingma & Gao, 2023) | 2B | - | 256·2 | 2.40 | 225.3 |
| ⚡ HDiT (**Ours**) | 557M | 742M | 50·2 | 6.92 | 135.2 |

architecture whose computational complexity scales with $\mathcal{O}(n)$ when used at higher resolutions instead of $\mathcal{O}(n^2)$, bridging the gap between the excellent scaling properties of transformer models and the efficiency of U-Nets. We demonstrate that this architecture enables megapixel-scale pixel-space diffusion models without requiring tricks such as self-conditioning or multiresolution architectures and that it is competitive with other transformer diffusion backbones even at small resolutions, both in fairly matched pixel-space settings, where it is substantially more efficient, and when compared to transformers in latent diffusion setups.

Given the promising results in this paper, we believe that ⚡ HDiT can provide a basis for further research into efficient high-resolution image synthesis. While we only focus on unconditional and class-conditional image synthesis, ⚡ HDiT is likely well-suited to provide efficiency and performance gains in other generative tasks like super-resolution, text-to-image generation and synthesis of other modalities such as audio and video, especially with architecture scaling.

# 7. Future Work

⧖ HDiT was studied in the context of pixel-space diffusion models but future works could investigate applying ⧖ HDiT in a latent diffusion setup to increase efficiency further and achieve multi-megapixel image resolutions, or apply orthogonal tricks such as self-conditioning (Jabri et al., 2023) or progressive training (Sauer et al., 2022) to improve the quality of generated samples further.

While the results for our large-scale ImageNet training presented in Section 5.4 are promising and perform competitively to many state-of-the-art architectures, we expect that substantial further improvements are possible with hyperparameter tuning and architecture scaling. Future work could explore how to fully realize the potential of this architecture.

Our architecture with local attention blocks could also be useful for efficient diffusion superresolution and diffusion VAE feature decoding models: if all levels are set to perform local attention only (global attention blocks should not be necessary as the global structure is already present in the samples for these applications), one can train efficient transformer-based models that can scale to arbitrary resolutions.

## Acknowledgements

## References

Balaji, Y., Nah, S., Huang, X., Vahdat, A., Song, J., Zhang, Q., Kreis, K., Aittala, M., Aila, T., Laine, S., Catanzaro, B., Karras, T., and Liu, M.-Y. eDiff-I: Text-to-Image Diffusion Models with an Ensemble of Expert Denoisers, 2023.

Bao, F., Nie, S., Xue, K., Cao, Y., Li, C., Su, H., and Zhu, J. All are Worth Words: A ViT Backbone for Diffusion Models. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023a.

Bao, F., Nie, S., Xue, K., Li, C., Pu, S., Wang, Y., Yue, G., Cao, Y., Su, H., and Zhu, J. One Transformer Fits All Distributions in Multi-Modal Diffusion at Scale. In International Conference on Machine Learning (ICML). JMLR.org, 2023b.

Betker, J., Goh, G., Jing, L., Brooks, T., Wang, J., Li, L., Ouyang, L., Zhuang, J., Lee, J., Guo, Y., Manassra, W., Dhariwal, P., Chu, C., Jiao, Y., and Ramesh, A. Improving Image Generation with Better Captions. Technical report, 2023.

Blattmann, A., Rombach, R., Ling, H., Dockhorn, T., Kim, S. W., Fidler, S., and Kreis, K. Align your Latents: High-Resolution Video Synthesis with Latent Diffusion Models. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023.

Cao, H., Wang, J., Ren, T., Qi, X., Chen, Y., Yao, Y., and Zhang, L. Exploring Vision Transformers as Diffusion Learners, 2022.

Chen, J., Yu, J., Ge, C., Yao, L., Xie, E., Wu, Y., Wang, Z., Kwok, J., Luo, P., Lu, H., and Li, Z. PixArt-$\alpha$: Fast Training of Diffusion Transformer for Photorealistic Text-to-Image Synthesis, 2023a.

Chen, S., Xu, M., Ren, J., Cong, Y., He, S., Xie, Y., Sinha, A., Luo, P., Xiang, T., and Perez-Rua, J.-M. GenTron: Delving Deep into Diffusion Transformers for Image and Video Generation, 2023b.

Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., Schuh, P., Shi, K., Tsvyashchenko, S., Maynez, J., Rao, A., Barnes, P., Tay, Y., Shazeer, N., Prabhakaran, V., Reif, E., Du, N., Hutchinson, B., Pope, R., Bradbury, J., Austin, J., Isard, M., Gur-Ari, G., Yin, P., Duke, T., Levskaya, A., Ghemawat, S., Dev, S., Michalewski, H., Garcia, X., Misra, V., Robinson, K., Fedus, L., Zhou, D., Ippolito, D., Luan, D., Lim, H., Zoph, B., Spiridonov, A., Sepassi, R., Dohan, D., Agrawal, S., Omernick, M., Dai, A. M., Pillai, T. S., Pellat, M., Lewkowycz, A., Moreira, E., Child, R., Polozov, O., Lee, K., Zhou, Z., Wang, X., Saeta, B., Diaz, M., Firat, O., Catasta, M., Wei, J., Meier-Hellstern, K., Eck, D., Dean, J., Petrov, S., and Fiedel, N. PaLM: Scaling Language Modeling with Pathways. Journal of Machine Learning Research (JMLR), 2023.

Crowson, K. DPM-Solver++(3M) SDE, 2023. URL https://github.com/crowsonkb/k-diffusion/blob/cc49cf6182284e577e896943f8e2/k_diffusion/sampling.py#L656.

Dai, X., Hou, J., Ma, C.-Y., Tsai, S., Wang, J., Wang, R., Zhang, P., Vandenhende, S., Wang, X., Dubey, A., Yu, M., Kadian, A., Radenovic, F., Mahajan, D., Li, K., Zhao, Y., Petrovic, V., Singh, M. K., Motwani, S., Wen, Y., Song, Y., Sumbaly, R., Ramanathan, V., He, Z., Vajda, P., and Parikh, D. Emu: Enhancing Image Generation Models Using Photogenic Needles in a Haystack, 2023.

Dehghani, M., Djolonga, J., Mustafa, B., Padlewski, P., Heek, J., Gilmer, J., Steiner, A., Caron, M., Geirhos, R., Alabdulmohsin, I., Jenatton, R., Beyer, L., Tschannen, M., Arnab, A., Wang, X., Riquelme, C., Minderer, M., Puigcerver, J., Evci, U., Kumar, M., Van Steenkiste, S., Elsayed, G. F., Mahendran, A., Yu, F., Oliver, A., Huot, F., Bastings, J., Collier, M. P., Gritsenko, A. A., Birodkar, V., Vasconcelos, C., Tay, Y., Mensink, T., Kolesnikov, A., Pavetić, F., Tran, D., Kipf, T., Lučić, M., Zhai, X., Keysers, D., Harmsen, J., and Houlsby, N. Scaling Vision Transformers to 22 Billion Parameters. In International Conference on Machine Learning (ICML). JMLR.org, 2023.

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. ImageNet: A large-scale hierarchical image database. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2009.

Dhariwal, P. and Nichol, A. Q. Diffusion Models Beat GANs on Image Synthesis. In Conference on Neural Information Processing Systems (NeurIPS), 2021.

Fedus, W., Zoph, B., and Shazeer, N. Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity. Journal of Machine Learning Research (JMLR), 2022.

Fischer, J. S., Gui, M., Ma, P., Stracke, N., Baumann, S. A., and Ommer, B. Boosting Latent Diffusion with Flow Matching, 2023.

Gao, S., Zhou, P., Cheng, M.-M., and Yan, S. Masked Diffusion Transformer is a Strong Image Synthesizer. In IEEE/CVF International Conference on Computer Vision (ICCV), October 2023.

Gu, J., Zhai, S., Zhang, Y., Susskind, J., and Jaitly, N. Matryoshka Diffusion Models, 2023.

Hang, T., Gu, S., Li, C., Bao, J., Chen, D., Hu, H., Geng, X., and Guo, B. Efficient Diffusion Training via Min-SNR Weighting Strategy. In IEEE/CVF International Conference on Computer Vision (ICCV), October 2023.

Hassani, A., Walton, S., Li, J., Li, S., and Shi, H. Neighborhood Attention Transformer. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), June 2023.

Henry, A., Dachapally, P. R., Pawar, S. S., and Chen, Y. Query-key normalization for transformers. In Cohn, T., He, Y., and Liu, Y. (eds.), Findings of the Association for Computational Linguistics: EMNLP 2020, November 2020.

Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), Conference on Neural Information Processing Systems (NeurIPS), 2017.

Ho, J. and Salimans, T. Classifier-Free Diffusion Guidance. In NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications, 2021.

Ho, J., Kalchbrenner, N., Weissenborn, D., and Salimans, T. Axial Attention in Multidimensional Transformers, 2019.

Ho, J., Jain, A., and Abbeel, P. Denoising Diffusion Probabilistic Models. In Conference on Neural Information Processing Systems (NeurIPS), 2020.

Ho, J., Saharia, C., Chan, W., Fleet, D. J., Norouzi, M., and Salimans, T. Cascaded Diffusion Models for High Fidelity Image Generation, 2021.

Hoogeboom, E., Heek, J., and Salimans, T. Simple Diffusion: End-to-End Diffusion for High Resolution Images. In International Conference on Machine Learning (ICML). JMLR.org, 2023.

Jabri, A., Fleet, D., and Chen, T. Scalable Adaptive Computation for Iterative Generation, 2023.

Jing, X., Chang, Y., Yang, Z., Xie, J., Triantafyllopoulos, A., and Schuller, B. W. U-DiT TTS: U-Diffusion Vision Transformer for Text-to-Speech, 2023.

Karras, T., Aittala, M., Hellsten, J., Laine, S., Lehtinen, J., and Aila, T. Training Generative Adversarial Networks with Limited Data. In Conference on Neural Information Processing Systems (NeurIPS), 2020a.

Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., and Aila, T. Analyzing and Improving the Image Quality of StyleGAN. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020b.

Karras, T., Laine, S., and Aila, T. A Style-Based Generator Architecture for Generative Adversarial Networks. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), dec 2021.

Karras, T., Aittala, M., Aila, T., and Laine, S. Elucidating the Design Space of Diffusion-Based Generative Models. In Conference on Neural Information Processing Systems (NeurIPS), 2022.

Kingma, D. P. and Gao, R. Understanding Diffusion Objectives as the ELBO with Simple Data Augmentation, 2023.

Kong, Z., Ping, W., Huang, J., Zhao, K., and Catanzaro, B. DiffWave: A Versatile Diffusion Model for Audio Synthesis. In International Conference on Learning Representations (ICLR), 2021.

Li, R., Li, W., Yang, Y., Wei, H., Jiang, J., and Bai, Q. Swinv2-Imagen: Hierarchical Vision Transformer Diffusion Models for Text-to-Image Generation, 2022.

Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. In IEEE/CVF International Conference on Computer Vision (ICCV), 2021.

Liu, Z., Hu, H., Lin, Y., Yao, Z., Xie, Z., Wei, Y., Ning, J., Cao, Y., Zhang, Z., Dong, L., Wei, F., and Guo, B. Swin Transformer V2: Scaling Up Capacity and Resolution. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022a.

Liu, Z., Hu, H., Lin, Y., Yao, Z., Xie, Z., Wei, Y., Ning, J., Cao, Y., Zhang, Z., Dong, L., Wei, F., and Guo, B. SWin Transformer v2, 2022b. URL https://github.com/microsoft/Swin-Transformer/blob/2cb103f2de145ff43bb9f6fc2ae8800c24/models/swin_transformer_v2.py#L156.

Loshchilov, I. and Hutter, F. Decoupled Weight Decay Regularization. In International Conference on Learning Representations (ICLR), 2019.

Lu, C., Zhou, Y., Bao, F., Chen, J., Li, C., and Zhu, J. DPM-Solver++: Fast Solver for Guided Sampling of Diffusion Probabilistic Models, 2023.

Nawrot, P., Tworkowski, S., Tyrolski, M., Kaiser, L., Wu, Y., Szegedy, C., and Michalewski, H. Hierarchical Transformers Are More Efficient Language Models. In Findings of the Association for Computational Linguistics: NAACL 2022, July 2022.

Nichol, A. Q. and Dhariwal, P. Improved denoising diffusion probabilistic models. In International Conference on Machine Learning (ICML). PMLR, 2021.

OpenAI. GPT-4 Technical Report. Technical report, 2023.

Peebles, W. and Xie, S. Scalable Diffusion Models with Transformers. In IEEE/CVF International Conference on Computer Vision (ICCV), October 2023a.

Peebles, W. and Xie, S. facebookresearch/dit, 2023b. URL https://github.com/facebookresearch/DiT/tree/ed81ce2229091fd4ecc9a22364.

Piergiovanni, A., Kuo, W., and Angelova, A. Rethinking Video ViTs: Sparse Video Tubes for Joint Image and Video Learning. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023.

Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. Hierarchical Text-Conditional Image Generation with CLIP Latents, 2022.

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-Resolution Image Synthesis With Latent Diffusion Models. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), June 2022.

Ronneberger, O., Fischer, P., and Brox, T. U-Net: Convolutional Networks for Biomedical Image Segmentation. In Medical Image Computing and Computer-Assisted Intervention (MICCAI), 2015.

Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S. K. S., Gontijo-Lopes, R., Ayan, B. K., Salimans, T., Ho, J., Fleet, D. J., and Norouzi, M. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), Conference on Neural Information Processing Systems (NeurIPS), 2022.

Saremi, S. and Sejnowski, T. J. Hierarchical model of natural images and the origin of scale invariance. Proceedings of the National Academy of Sciences, 110 (8):3071–3076, February 2013. ISSN 1091-6490. doi: 10.1073/pnas.1222618110. URL http://dx.doi.org/10.1073/pnas.1222618110.

Sauer, A., Schwarz, K., and Geiger, A. StyleGAN-XL: Scaling StyleGAN to Large Diverse Datasets. In ACM SIGGRAPH 2022 Conference Proceedings. Association for Computing Machinery, 2022.

Shazeer, N. GLU Variants Improve Transformer, 2020.

Shi, W., Caballero, J., Huszar, F., Totz, J., Aitken, A. P., Bishop, R., Rueckert, D., and Wang, Z. Real-Time Single Image and Video Super-Resolution Using an Efficient Sub-Pixel Convolutional Neural Network. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), jun 2016.

Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-Based Generative Modeling through Stochastic Differential Equations. In International Conference on Learning Representations (ICLR), 2021.

Stein, G., Cresswell, J. C., Hosseinzadeh, R., Sui, Y., Ross, B. L., Villecroze, V., Liu, Z., Caterini, A. L., Taylor, J. E. T., and Loaiza-Ganem, G. Exposing flaws of generative model evaluation metrics and their unfair treatment of diffusion models, 2023.

Su, J., Lu, Y., Pan, S., Murtadha, A., Wen, B., and Liu, Y. RoFormer: Enhanced Transformer with Rotary Position Embedding, 2022.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. LLaMA: Open and Efficient Foundation Language Models. Technical report, 2023.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. Attention is All you Need. In Conference on Neural Information Processing Systems (NeurIPS), 2017.

Wang, P. Flash Cosine Similarity Attention, 2022. URL https://github.com/lucidrains/flash-cosine-sim-attention/tree/6f17f29a979a8bcab2479c65b7740523.

Wang, Z., Cun, X., Bao, J., Zhou, W., Liu, J., and Li, H. Uformer: A General U-Shaped Transformer for Image Restoration. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022.

Yan, J. N., Gu, J., and Rush, A. M. Diffusion Models Without Attention, 2023.

Yang, X., Shih, S.-M., Fu, Y., Zhao, X., and Ji, S. Your ViT is Secretly a Hybrid Discriminative-Generative Diffusion Model, 2022.

Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., and Wu, Y. CoCa: Contrastive Captioners are Image-Text Foundation Models. Transactions on Machine Learning Research (TMLR), 2022.

Zhang, B., Gu, S., Zhang, B., Bao, J., Chen, D., Wen, F., Wang, Y., and Guo, B. StyleSwin: Transformer-Based GAN for High-Resolution Image Generation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 11304–11314, June 2022a.

Zhang, Y., Qin, J., Park, D. S., Han, W., Chiu, C.-C., Pang, R., Le, Q. V., and Wu, Y. Pushing the Limits of Semi-Supervised Learning for Automatic Speech Recognition, 2022b.

Zhao, L., Zhang, Z., Chen, T., Metaxas, D., and Zhang, H. Improved Transformer for High-Resolution GANs. In Conference on Neural Information Processing Systems (NeurIPS), 2021.

Zheng, H., Nie, W., Vahdat, A., and Anandkumar, A. Fast Training of Diffusion Models with Masked Transformers, 2023.

Zong, Z., Song, G., and Liu, Y. DETRs with Collaborative Hybrid Assignments Training. In IEEE/CVF International Conference on Computer Vision (ICCV), 2022.

# A. Computational Complexity of HDiT

In a traditional vision transformer, including those for diffusion models (Peebles & Xie, 2023a; Bao et al., 2023a), the asymptotic computational complexity with regard to image size is dominated by the self-attention mechanism, which scales as $\mathcal{O}(n^2 d)$ with token/pixel count $n$ and embedding dimension $d$. The feedforward blocks and the attention projection heads, in turn, scale as $\mathcal{O}(nd^2)$.

For our Hourglass Diffusion Transformer architecture, we adjust the architecture for different target resolutions, similarly to previous approaches used with U-Nets (Ronneberger et al., 2015). Our architecture is generally divided into multiple hierarchical levels, where the outermost level operates at full patch resolution, and each additional level operates at half of the spatial resolution per axis. For simplicity, we will first cover the cost at square resolutions of powers of two.

When designing the architecture for a specific resolution, we start with a dataset-dependent *core* architecture, which, for natural images, typically includes one or two global-attention hierarchy levels that operate at $16^2$ or $16^2$ and $32^2$, respectively. Around that are a number of local attention levels. As this core only operates on a fixed resolution, it does not influence the asymptotic computational complexity of the overall model.

**Asymptotic Complexity Scaling** When this architecture is adapted to a higher resolution, additional local attention levels with shared parameters are added to keep the innermost level operating at $16^2$. This means that the number of levels in our hierarchy scales with the number of image tokens as $\mathcal{O}(\log(n))$. While this might intuitively lead one to the conclusion of the overall complexity being $\mathcal{O}(n \log(n))$, as local attention layers' complexity is $\mathcal{O}(nd)$, the reduction in resolution at each level in the hierarchy has to be considered: due to the spatial downsampling, the number of tokens decreases by a factor of four at every level in the hierarchy, making the cost of the self-attention – the only part of our model whose complexity does not scale linearly with token count – of the additional levels

$$\sum_{l=1}^{\log_4(n)-\log_4(\mathrm{res_{core}})} \frac{nd}{4^{l-1}}.$$

Factoring out $n$ and defining $m = l - 1$ yields

$$n \cdot \sum_{m=0}^{\log_4(n)-\log_4(\mathrm{res_{core}})-1} d \cdot \left(\frac{1}{4}\right)^m,$$

a (cut-off) geometric series with a common ratio of less than one, which means that, as the geometric series converges, it does not affect the asymptotic complexity, making the cumulative complexity of the local self-attention of the additional levels $\mathcal{O}(n)$. Thus, as no other parts of the scale worse

than $\mathcal{O}(n)$ either, the overall complexity of the Hourglass Diffusion Transformer architecture, as the target resolution is increased, is $\mathcal{O}(n)$.

**Local Complexity Scaling at Arbitrary Resolutions** When the target resolution is increased by a factor smaller than a power of two per axis, the architecture is not adapted. This means that, for these intermediate resolutions, a different scaling behavior prevails. Here, the cost of the local attention levels, whose number does not change in this case, scales with $\mathcal{O}(n)$ as before, but the global attention levels incur a quadratic increase in cost with the resolution. As the resolution is increased further, however, new levels are added, which reduce the resolution the global attention blocks operate at to their original values, and retaining the overall asymptotic scaling behavior of $\mathcal{O}(n)$.

# B. Soft-Min-SNR Loss Weighting

Min-SNR loss weighting (Hang et al., 2023) is a recently introduced training loss weighting scheme that improves diffusion model training. It adapts the SNR weighting scheme (for image data scaled to $\mathbf{x} \in [-1, 1]^{h \times w \times c}$)

$$w_{\mathrm{SNR}}(\sigma) = \frac{1}{\sigma^2} \tag{7}$$

by clipping it at an SNR of $\gamma = 5$:

$$w_{\mathrm{Min\text{-}SNR}}(\sigma) = \min\left\{\frac{1}{\sigma^2}, \gamma\right\}. \tag{8}$$

We utilize a slightly modified version that smoothes out the transition between the normal SNR weighting and the clipped section:

$$w_{\mathrm{Soft\text{-}Min\text{-}SNR}}(\sigma) = \frac{1}{\sigma^2 + \gamma^{-1}}. \tag{9}$$

For $\sigma \ll \gamma$ and $\sigma \gg \gamma$, this matches Min-SNR, while providing a smooth transition between both sections.

In practice, we also change the hyperparameter $\gamma$ from $\gamma = 5$ to $\gamma = 4$.

Plotting the resulting loss weight for both min-snr and our soft-min-snr as shown in Figure 8 shows that our loss weighting is identical to min-snr, except for the transition, where it is significantly smoother. An ablation of our soft-min-snr compared to min-snr also shows that our loss weighting scheme leads to an improved FID score (see Table 5) for our model.
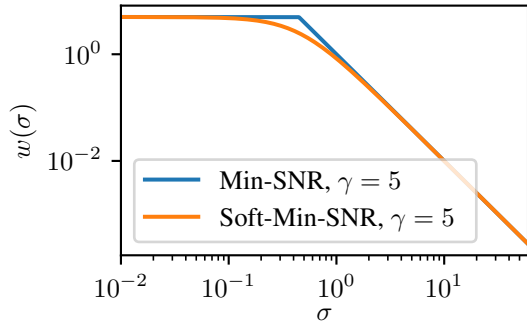
Figure 8: The resulting loss weighting over $\sigma$ for our soft-min-snr weighting (orange) and min-snr weighting (blue) with $\gamma = 5$.

Table 5: Soft-Min-SNR ablation on RGB ImageNet-$128^2$.

| Loss Weighting | FID↓ |
|---|---|
| SNR (Table 1 step **D**) | 41.41 |
| Min-SNR (Hang et al., 2023) ($\gamma = 5$) | 36.65 |
| Min-SNR (Hang et al., 2023) ($\gamma = 4$) | 35.62 |
| Soft-Min-SNR (**Ours**, $\gamma = 4$, Table 1 step **E**) | **27.74** |

## C. Scaled Cosine Similarity Attention

For the attention mechanism in ⏳ HDiT, we use a slight variation of the cosine similarity-based attention introduced in (Liu et al., 2022a) they dub *Scaled Cosine Attention*: instead of computing the self-attention as

$$\text{SA}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_{\text{head}}}}\right)V, \quad (10)$$

they compute it as

$$\text{SCA}(Q, K, V) = \text{softmax}\left(\frac{\text{sim}_{\cos}(Q, K)}{\tau} + B_{ij}\right)V, \quad (11)$$

with $\tau$ being a per-head per-layer learnable scalar, and $B_{ij}$ being the relative positional bias between pixel $i$ and $j$ (which we do not use in our models). In practice, they parametrize $\tau$ based on a learnable parameter $\theta$ in the following way (Liu et al., 2022b):

$$\frac{1}{\tau} = \exp\left(\min\left\{\theta, \log\frac{1}{0.01}\right\}\right), \quad (12)$$

with $\theta$ being initialized to $\theta = \log 10$.

### C.1. Improving Scale Learning Stability

We find that their parametrization of $\tau$ causes the learned scales to vary significantly during training, necessitating the

clamping to a maximum value of 100 before exponentiation to prevent destabilization of the training. In this setting, we find that a significant number of scale factors $\tau$ reach this maximum value and values below 1 during our trainings. We speculate that this instability might be the cause of the behaviour observed in (Wang, 2022), where using scaled cosine similarity attention was detrimental to the performance of generative models. To alleviate this problem, we find simply learning $\tau$ directly, as done for normal attention in (Henry et al., 2020), prevents this large variance of its values in our models, with our converged models' scale typically reaching a range between 5 and 50.

## D. Additional Results for ImageNet-$256^2$

In addition to the analyses in Section 5.4, which do not use classifier-free guidance (Ho & Salimans, 2021), we also analyze the FID-IS-tradeoff for difference guidance scales $w_{cfg}$ (we follow the guidance scale formulation used in (Saharia et al., 2022), where $w_{cfg} = 1$ corresponds to no classifier-free guidance being applied). The resulting curve is shown in Figure 9, with the lowest FID of 3.21 being achieved around $w_{cfg} = 1.3$, with a corresponding IS of 220.6.
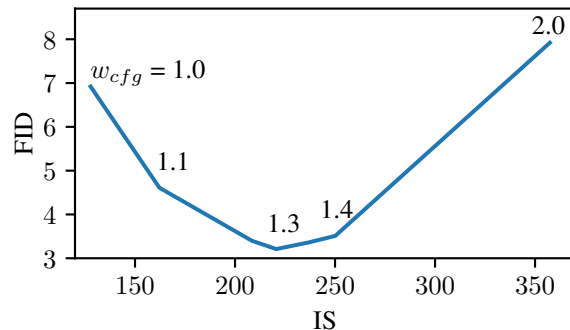


Figure 9: Inception Score vs. Fréchet Inception Distance at different classifier-free guidance weight scales (1 = no guidance) for our 557M ImageNet-$256^2$ model.

# E. Experiment Details

| Parameter | ImageNet-$128^2$ | FFHQ-$1024^2$ | ImageNet-$256^2$ |
|---|---|---|---|
| Experiment | Ablation **E**[5] (Section 5.2) | High-Res Synthesis (Section 5.3) | Large-Scale (Section 5.4) |
| Parameters | 117M | 85M | 557M |
| GFLOP/forward | 31 | 206 | 198 |
| Training Steps | 400k | 1M | 2.2M[6] |
| Batch Size | 256 | 256 | 256+[6] |
| Precision | bfloat16 | bfloat16 | bfloat16 |
| Training Hardware | 4 A100 | 64 A100 | 8 H100 |
| Training Time | 15 hours | 5 days | - |
| Patch Size | 4 | 4 | 4 |
| Levels (Local + Global Attention) | 1 + 1 | 3 + 2 | 2 + 1 |
| Depth | [2, 11] | [2, 2, 2, 2, 2] | [2, 2, 16] |
| Widths | [384, 768] | [128, 256, 384, 768, 1024] | [384, 768, 1536] |
| Attention Heads (Width / Head Dim) | [6, 12] | [2, 4, 6, 12, 16] | [6, 12, 24] |
| Attention Head Dim | 64 | 64 | 64 |
| Neighborhood Kernel Size | 7 | 7 | 7 |
| Mapping Depth | 1 | 2 | 2 |
| Mapping Width | 768 | 768 | 768 |
| Data Sigma | 0.5 | 0.5 | 0.5 |
| Sigma Range | [1e-3, 1e3] | [1e-3, 1e3] | [1e-3, 1e3] |
| Sigma Sampling Density | interpolated cosine | interpolated cosine | interpolated cosine |
| Augmentation Probability | 0 | 0.12 | 0 |
| Dropout Rate | 0 | [0, 0, 0, 0, 0.1] | 0 |
| Conditioning Dropout Rate | 0.1 | 0.1 | 0.1 |
| Optimizer | AdamW | AdamW | AdamW |
| Learning Rate | 5e-4 | 5e-4 | 5e-4 |
| Betas | [0.9, 0.95] | [0.9, 0.95] | [0.9, 0.95] |
| Eps | 1e-8 | 1e-8 | 1e-8 |
| Weight Decay | 1e-2 | 1e-2 | 1e-2 |
| EMA Decay | 0.9999 | 0.9999 | 0.9999 |
| Sampler | DPM++(3M) SDE | DPM++(3M) SDE | DPM++(3M) SDE |
| Sampling Steps | 50 | 50 | 50 |

Table 6: Details of our training and inference setup.

---

[5]The other ablation steps generally use the same parameters, except for the architectural changes indicated in the experiment description.

[6]We initially trained for 2M steps. We then experimented with progressively increasing the batch size (waiting until the loss plateaued to a new, lower level each time), training at batch size 512 for an additional 50k steps, at batch size 1024 for 100k, and at batch size 2048 for 50k steps.

# F. Our FFHQ-$1024^2$ Samples



Figure 10: Uncurated samples from our 85M ⧗ HDiT FFHQ-$1024^2$ model.

# G. NCSN++ (Song et al., 2021) FFHQ-$1024^2$ Reference Samples



Figure 11: Uncurated reference samples from the NCSN++ (Song et al., 2021) FFHQ-$1024^2$ baseline model.

# H. Our ImageNet-$256^2$ Samples



Figure 12: Uncurated random class-conditional samples from our 557M ⧗ HDiT ImageNet-$256^2$ model.

Figure 13: More uncurated random class-conditional samples from our ⌛ HDiT-557M ImageNet-256$^2$ model.