

# PixNerd: Pixel Neural Field Diffusion

Shuai Wang<sup>1</sup> Ziteng Gao<sup>3</sup> Chenhui Zhu<sup>1</sup> Weilin Huang<sup>2</sup> Limin Wang<sup>1</sup>

<sup>1</sup>Nanjing University <sup>2</sup>ByteDance Seed <sup>3</sup>National University of Singapore

<https://github.com/MCG-NJU/PixNerd>

<https://huggingface.co/spaces/MCG-NJU/PixNerd>

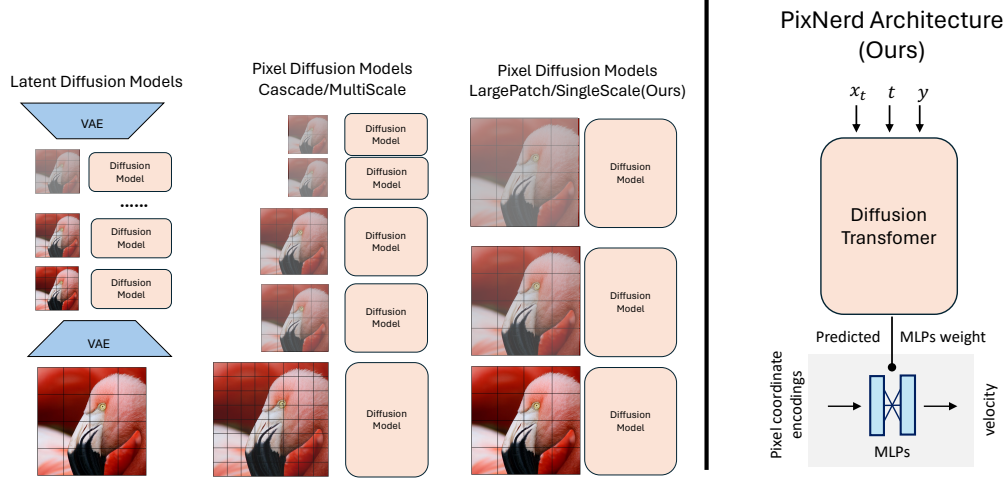


Figure 1: **Left: Comparison with other diffusion models.** Our LargePatch/SingleScale pixel space diffusion keeps consistent tokens as latent diffusion among diffusion steps. **Right: PixNerd architecture.** PixNerd follows the diffusion transformer design, replacing the final linear projection with a neural field to model the large patch details.

## Abstract

The current success of diffusion transformers heavily depends on the compressed latent space shaped by the pre-trained variational autoencoder (VAE). However, this two-stage training paradigm inevitably introduces accumulated errors and decoding artifacts. To address the aforementioned problems, researchers return to pixel space at the cost of complicated cascade pipelines and increased token complexity. In contrast to their efforts, we propose to model the patch-wise decoding with neural field and present a single-scale, single-stage, efficient, end-to-end solution, coined as pixel neural field diffusion (PixelNerd). Thanks to the efficient neural field representation in PixNerd, we directly achieved 2.15 FID on ImageNet  $256 \times 256$  and 2.84 FID on ImageNet  $512 \times 512$  without any complex cascade pipeline or VAE. We also extend our PixNerd framework to text-to-image applications. Our PixNerd-XXL/16 achieved a competitive 0.73 overall score on the GenEval benchmark and 80.9 overall score on the DPG benchmark.

## 1 Introduction

The current success of diffusion transformers largely depends on variational autoencoders (VAEs) [1, 2, 3]. VAEs significantly reduce the spatial dimension of raw pixels while providing a compact and nearly lossless latent space, substantially easing the learning difficulty for diffusion transformers. By operating in this compressed latent space, diffusion transformers can effectively learn either the score function or velocity of the diffusion process using small patch sizes. However, training high-quality

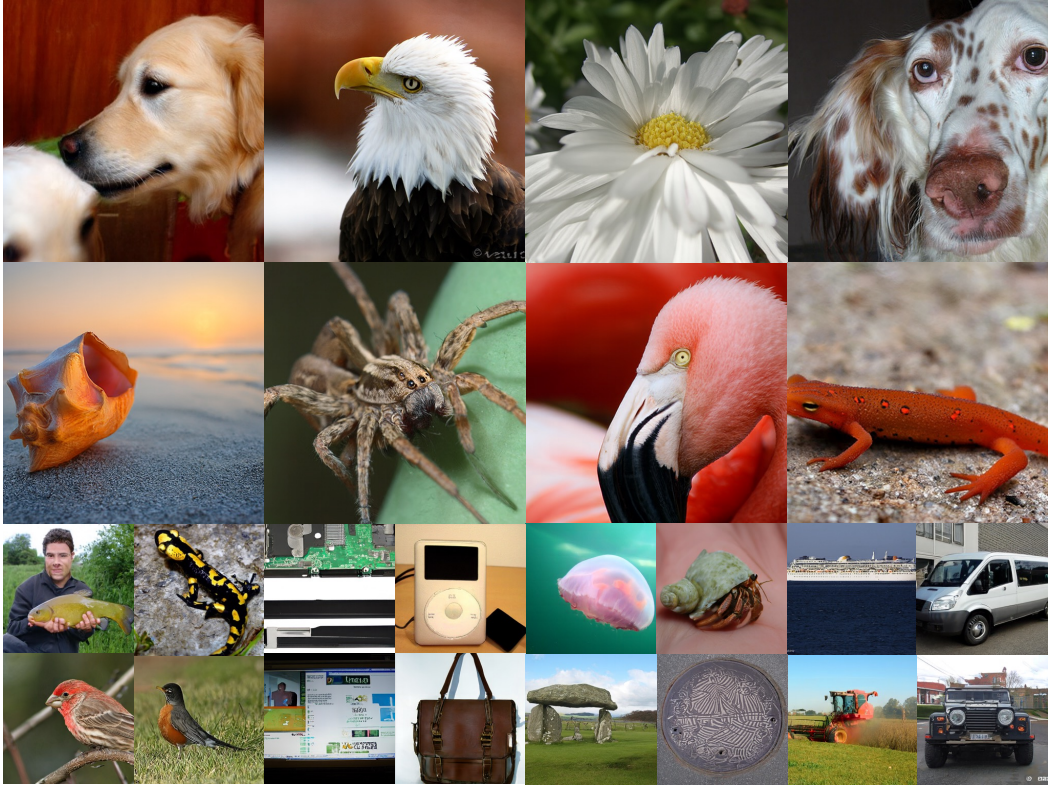


Figure 2: **Selected  $256 \times 256$  and  $512 \times 512$  resolution samples.** Generated from PixNerd-XL/16 trained on ImageNet  $256 \times 256$  resolution and ImageNet  $512 \times 512$  resolution with CFG = 3.5.

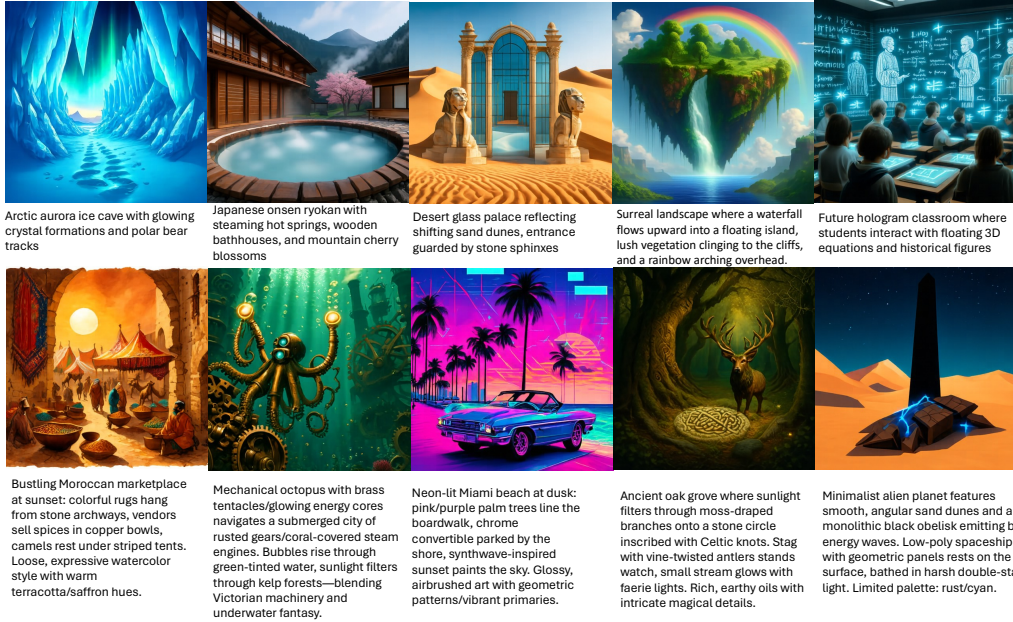


Figure 3: **The Text-to-Image  $512 \times 512$  visualization with text descriptions of different lengths and styles.** Given text descriptions of different lengths and styles, PixNerd can generate promising samples with a large patch size of 16. We used Adams-2nd solver with 25 steps and a CFG value of 4.0 for sampling.

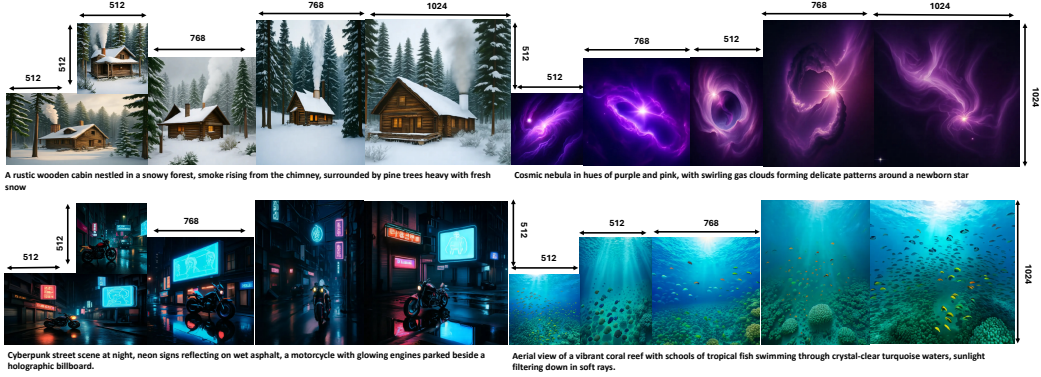


Figure 4: **Training-free arbitrary resolution generation.** We keep the amount of tokens in PixNerd as constant as pretraining resolution, we **only interpolate the neural field coordinates** for different resolutions to yield multi-resolution images.

VAEs typically requires adversarial training [4, 5, 1, 2, 3] and additional supervision [6], introducing complex optimization challenges. Moreover, this two-stage training paradigm leads to accumulated errors and inevitable decoding artifacts. To address these limitations, some researchers have explored joint training approaches, though these come with substantial computational costs [7].

An alternative approach involves implementing diffusion models directly in raw pixel space [8, 9, 10]. In contrast to the success of latent diffusion transformers, progress in pixel-space diffusion transformers proves significantly more challenging. Without the dimensionality reduction provided by VAEs [2, 1, 3], pixel diffusion transformers allocate substantially more image tokens [9], requiring impractical computational resources when using the same patch size as latent diffusion transformers. To maintain a comparable number of image tokens, pixel diffusion transformers must employ much larger patch sizes during denoising training. However, due to the vastness of raw pixel space, larger patches make diffusion learning particularly difficult. Previous methods [9, 10] have employed cascade solutions that divide the diffusion process across different scales to reduce computational costs. However, this cascade approach complicates both training and inference. In contrast to these methods, our work explores the performance upper bound using a large-patch diffusion transformer while maintaining the same token count and comparable computational requirements as latent diffusion models.

Inspired by the success of implicit neural fields in scene reconstruction [11, 12], we propose modeling large-patch decoding using an implicit neural field. We replace the traditional diffusion transformer’s linear projection with an implicit neural field, naming this novel pixel-space architecture PixelNerd (**P**ixel **N**erual **F**ield **D**iffusion). Specifically, we predict the weights of each patch’s neural field MLPs using the diffusion transformer’s last hidden states. For each pixel within a local patch, we first transform its local coordinates into coordinate encoding. This encoding, combined with the corresponding noisy pixel value, is then processed by the neural field MLPs to predict the diffusion velocity. Our approach significantly alleviates the challenge of learning fine details under large-patch configurations.

Compared to previous latent diffusion transformers and other pixel-space diffusion models, our end-to-end PixelNerd offers a simpler, more elegant, and efficient solution. For class-conditional image generation on ImageNet  $256 \times 256$ , PixelNerd-XL/16 achieves a competitive 2.15 FID and significantly better 4.55 sFID (indicating superior spatial structure) than PixelFlow [9, 10]. On ImageNet  $512 \times 512$ , PixelNerd-XL/16 maintains comparable performance with 2.84 FID. For text-to-image generation, PixelNerd-XXL/16- $512 \times 512$  achieves 73.0 on the GenEval benchmark and 80.9 average score on DPG benchmark.

Our contributions are summarized as follows:

- We propose a novel, elegant, efficient end-to-end pixel space diffusion transformer with neural field, deemed as PixNerd.
- For the class-to-image generation, on ImageNet  $256 \times 256$ , our PixNerd-XL/16 achieves a comparable 2.15 FID with similar computation demands as its latent counterpart. On



ImageNet  $512 \times 512$ , our PixNerd-XL/16 achieves a comparable 2.84 FID with similar computation demands as its latent counterpart.

- For the text-to-image generation, our PixNerd-XXL/16 achieves a 0.73 overall score on the GenEval benchmark and 80.9 average score on DPG benchmark.

## 2 Related Work

**Latent Diffusion Models** train diffusion models on a compact latent space shaped by a VAE [1]. Compared to raw pixel space, this latent space significantly reduces spatial dimensions, easing both learning difficulty and computational demands [1, 2, 3]. Thus VAE has become a core component in modern diffusion models [13, 14, 15, 16, 17, 18, 19, 20]. However, VAE training typically involves adversarial training and perceptual supervision, complicating the overall pipeline. Insufficient VAE training can lead to decoding artifacts [21], limiting the broader applicability of diffusion generative models. Earlier latent diffusion models primarily focused on U-Net-like architectures. The pioneering work of DiT [13] introduced transformers into diffusion models, replacing the traditionally dominant U-Net [22, 8]. Empirical results [13] show that, given sufficient training iterations, diffusion transformers outperform conventional approaches without relying on long residual connections. SiT [14] further validated the transformer architecture with linear flow diffusion.

**Pixel Diffusion Models** have progressed much more slowly than their latent counterparts. Due to the vastness of pixel space, learning difficulty and computational demands are typically far greater than those of latent diffusion models [8, 23, 24]. Current pixel-space diffusion models still rely on long residuals [10, 8], limiting further scaling. Early attempts primarily split the diffusion process into chunks at different resolution scales to reduce computational burdens [9, 10]. Pixelflow [9] uses the same denoising model across all scales, while Relay Diffusion [10] employs distinct models for each. However, this cascaded pipeline inevitably complicates both training and sampling. Additionally, training these models at isolated resolutions hinders end-to-end optimization and requires carefully designed strategies. FractalGen [25] constructs fractal generative models by recursively applying atomic generative modules, resulting in self-similar architectures that excel in pixel-by-pixel image generation. TarFlow [26] introduces a Transformer-based normalizing flow architecture capable of directly modeling and generating pixels.

## 3 Method

### 3.1 Preliminary

**Diffusion Models** gradually adds  $\mathbf{x}_0$  with Gaussian noise  $\epsilon$  to perturb the corresponding known data distribution  $p(\mathbf{x}_0)$  into a simple Gaussian distribution. The discrete perturbation function of each  $t$  satisfies  $\mathcal{N}(\mathbf{x}_t | \alpha_t \mathbf{x}_0, \sigma_t^2 \mathbf{I})$ , where  $\alpha_t, \sigma_t > 0$ . It can also be written as Eq. (1).

$$\mathbf{x}_t = \alpha_t \mathbf{x}_{\text{real}} + \sigma_t \epsilon. \quad (1)$$

Moreover, as shown in Eq. (2), Eq. (1) has a forward continuous-SDE description, where  $f(t) = \frac{d \log \alpha_t}{dt}$  and  $g(t) = \frac{d \sigma_t^2}{dt} - (\frac{d \log \alpha_t}{dt} \sigma_t^2)$ . [27] establishes a pivotal theorem that the forward SDE has an equivalent reverse-time diffusion process as in Eq. (3), so the generating process is equivalent to solving the diffusion SDE. Typically, diffusion models employ neural networks and distinct prediction parametrization to estimate the score function  $\nabla \log_x p_{\mathbf{x}_t}(\mathbf{x}_t)$  along the sampling trajectory [28, 29, 30].

$$d\mathbf{x}_t = f(t)\mathbf{x}_t dt + g(t)d\mathbf{w}. \quad (2)$$

$$d\mathbf{x}_t = [f(t)\mathbf{x}_t - g(t)^2 \nabla_{\mathbf{x}} \log p(\mathbf{x}_t)]dt + g(t)d\mathbf{w}. \quad (3)$$

VP [28] also shows that there exists a corresponding deterministic process Eq. (4) whose trajectories share the same marginal probability densities of Eq. (3).

$$d\mathbf{x}_t = [f(t)\mathbf{x}_t - \frac{1}{2}g(t)^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)]dt. \quad (4)$$

Rectified Flow Model simplifies diffusion model under the framework of Eq. (2) and Eq. (3). Different from [30] introduces non-linear transition scheduling, the rectified-flow model adopts linear

function to transform data to standard Gaussian noise. Instead of estimating the score function  $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ , rectified-flow models directly learn a neural network  $v_\theta(\mathbf{x}_t, t)$  to predict the velocity field  $\mathbf{v}_t = d\mathbf{x}_t = (\mathbf{x}_{\text{real}} - \epsilon)$ .

**Diffusion Transformer** was introduced into diffusion models [13] to replace the traditionally dominant UNet architecture [22, 8]. Empirical evidence demonstrates that given sufficient training iterations, diffusion transformers outperform conventional approaches even without relying on long residual connections. SiT [14] further validated the transformer architecture with linear flow diffusion. Given a noisy image latent  $\mathbf{x}_t$  as Eq. (1),  $\mathbf{y}$  is the condition,  $t$  is the timestep, we first partition it into non-overlapping patches, converting it into a 1D sequence. These noisy patches are then processed through stacked self-attention and FFN blocks, with class label conditions incorporated via AdaLN modulation. Finally, a simple linear projection decodes the feature patches into either patch-wise score or velocity estimates:

$$\mathbf{X}_t = \mathbf{X}_t + \text{AdaLN}(\mathbf{y}, t, \text{Attention}(\mathbf{X}_t)), \quad (5)$$

$$\mathbf{X}_t = \mathbf{X}_t + \text{AdaLN}(\mathbf{y}, t, \text{FFN}(\mathbf{X}_t)). \quad (6)$$

Recent architectural improvements such as SwiGLU [31, 32], RoPE [33], and RMSNorm [31, 32] have been extensively validated in the research community [34, 2, 35, 15, 36, 37, 38].

**Neural Field** is usually adopted to represent a scene through MLPs that map coordinates encodings to signals [11, 12, 39, 40, 41]. It has been widely applied to objects [42, 43] and surface reconstruction [44, 45, 46]. Specifically, recall that an MLP consists of a Linear, SiLU, and another Linear, we regard  $\mathbf{W}_1^n$  as the weight for the first linear layer in MLP, while  $\mathbf{W}_2^n$  is for the second. If we have a neural field with a single MLP  $\{\mathbf{W}_1, \mathbf{W}_2\}$  for naive 2D scene, given the query coordinate  $(i, j)$ , the coordinate will be transformed into cosine/sine encodings in Eq. (7) or DCT-basis encodings in Eq. (12):

$$\text{PE}(i, j) = \sin(2^0 \pi i), \cos(2^0 \pi i), \dots, \sin(2^L \pi i), \cos(2^L \pi i), \dots, \sin(2^L \pi j), \cos(2^L \pi j) \quad (7)$$

Then this encoding feature will be fed into neural field MLPs to extract features  $\mathbf{V}^n(i, j)$  as Eq. (8):

$$\mathbf{V}^n(i, j) = \text{MLP}(\text{PE}(i, j)) | \{\mathbf{W}_1, \mathbf{W}_2\}. \quad (8)$$

Finally,  $\mathbf{V}^n(i, j)$  can be used to regress the needed value, eg. RGB [11, 12], density and SDF [45].

### 3.2 Diffusion Transformer with Patch-wise Neural Field

While VAEs [1, 2, 3] significantly reduce spatial dimensions in the latent space, rendering a single linear projection sufficient for velocity modeling in latent diffusion transformers, pixel diffusion transformers must handle substantially larger patch sizes to maintain computational parity with their latent counterparts. Under such conditions, a simple linear projection becomes inadequate for capturing fine details.

To address the limitations of linear projection, we propose modeling patch-wise velocity decoding using an implicit neural field. Formally, given the last hidden states  $\mathbf{X}^n$  of the  $n$ -th patch in diffusion transformer, we predict neural field parameters  $\{\mathbf{W}_1^n \in \mathbb{R}^{D_2 \times D_1}, \mathbf{W}_2^n \in \mathbb{R}^{D_1 \times D_2}\}$  from  $\mathbf{X}^n$ .

$$\mathbf{W}_1^n, \mathbf{W}_2^n = \text{Linear}(\text{SiLU}(\mathbf{X}^n)). \quad (9)$$

To decode the pixel-wise velocity  $\mathbf{v}^n(i, j)$  for the pixel coordinate  $(i, j)$  in the  $n$ -th patch feature  $\mathbf{X}^n$ , where  $i, j \in (0, K)$ , we first encode the coordinates into encodings. These encodings  $(\text{PE}(i, j))$ , along with the noisy pixel value  $\mathbf{x}^n(i, j)$ , are then fed into the neural field MLP to predict the velocity. To enhance performance and stabilize training, we apply row-wise normalization to the neural field parameters. For brevity, we omit the timestep subscript here.

$$\mathbf{V}^n(i, j) = \text{MLP}(\text{Concat}([\text{PE}(i, j), \mathbf{x}^n(i, j)]) | \{ \frac{\mathbf{W}_1^n}{\|\mathbf{W}_1^n\|}, \frac{\mathbf{W}_2^n}{\|\mathbf{W}_2^n\|} \} ). \quad (10)$$

Finally, as shown in Eq. (11) the pixel velocity  $\mathbf{v}^n(i, j)$  is decoded from  $\mathbf{V}^n(i, j)$  through a linear projection:

$$\mathbf{v}^n(i, j) = \text{Linear}(\mathbf{V}^n(i, j)). \quad (11)$$

Model	Inference			Training	
	1 image	1 step	Mem (GB)	Speed (s/it)	Mem (GB)
SiT-L/2(VAE-f8)	0.51s	0.0097s	2.9	0.30	18.4
Baseline-L/16	0.48s	0.0097s	2.1	0.18	18
PixNerd-L/16	0.51s	0.010s	2.1	0.19	22
ADM-G	4.21s	0.08s	2.23	/	/
PixelFlow-XL/4	10.1s	0.084s <sup>†</sup>	4.0	/	/
PixNerd-XL/16	0.65s	0.012s	3.1	0.27	33.9

Table 1: **The resource consumption comparison.** <sup>†</sup> means the average time consumption for a single step across different stages. Our PixNerd consumes much less memory and latency (Nearly 8× faster than other pixel diffusion models).



Figure 5: **The visualization Comparison with Baseline-L/16 under 400k training steps.** With the help of neural field representation, our PixNerd-L/16 yields promising details and better structure.

## 4 Experiments

We conduct ablation studies and baseline comparison experiments on ImageNet- $256 \times 256$ . For class-to-image generation, we provide system-level comparisons on ImageNet- $256 \times 256$  and ImageNet- $512 \times 512$ , and report FID [47], sFID [48], IS [49], Precision, and Recall [50] as the main metrics. For text-to-image generation, we report results collected on the GenEval [51] and DPG [52] benchmarks.

**Training Details.** Inspired by recent advances in training schedulers [53], architecture design [34, 2, 35, 15], and representation alignment [54], we incorporate SwiGLU, RMSNorm [34, 31, 32], lognorm sampling [53], and representation alignment from DINOv2 [55, 54, 15] to enhance our PixNerd model. Specifically, we add an additional representation alignment loss with a weight of 0.5 to align the intermediate features from the 8th layer of our PixNerd model with the features extracted by DINOv2-Base.

**Resource Consumption.** We employ `torch.compile` to optimize memory allocation and reduce redundant computations for both the baseline and PixNerd. As shown in Tab. 1, compared with latent counterparts, our PixNerd-L/16 achieves much higher training throughput without VAE latency and with similar inference memory. Compared to other pixel-space diffusion models, our PixNerd consumes significantly less memory and has lower latency (nearly 8× faster than ADM-G [8] and PixelFlow [9]).

### 4.1 Comparison with Baselines

We conduct a baseline comparison on ImageNet  $256 \times 256$  with a large-size model. Both Baseline-L/16 and PixNerd-L/16 are built upon SwiGLU [32, 31], RoPE2d [33], RMSNorm, and trained with lognorm sampling. All optimizer configurations are consistently aligned. As shown in Fig. 6b and Fig. 6a, our model achieves consistently lower loss expectation. We also provide visualization comparisons with Baseline-L/16 in Fig. 5: PixNerd-L/16, trained for the same number of steps, generates better details and structures.

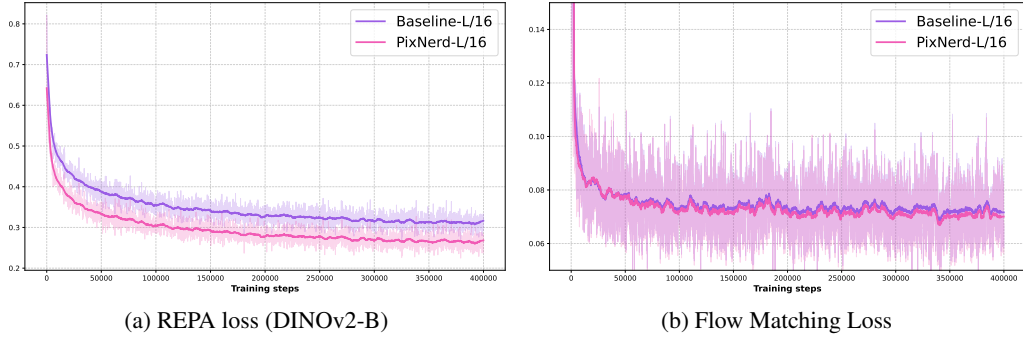


Figure 6: **Loss Comparison with Diffusion Transformer Baselines.** Our PixNerd-L/16 achieves consistently lower REPA loss and flow matching loss than its diffusion transformer counterpart.

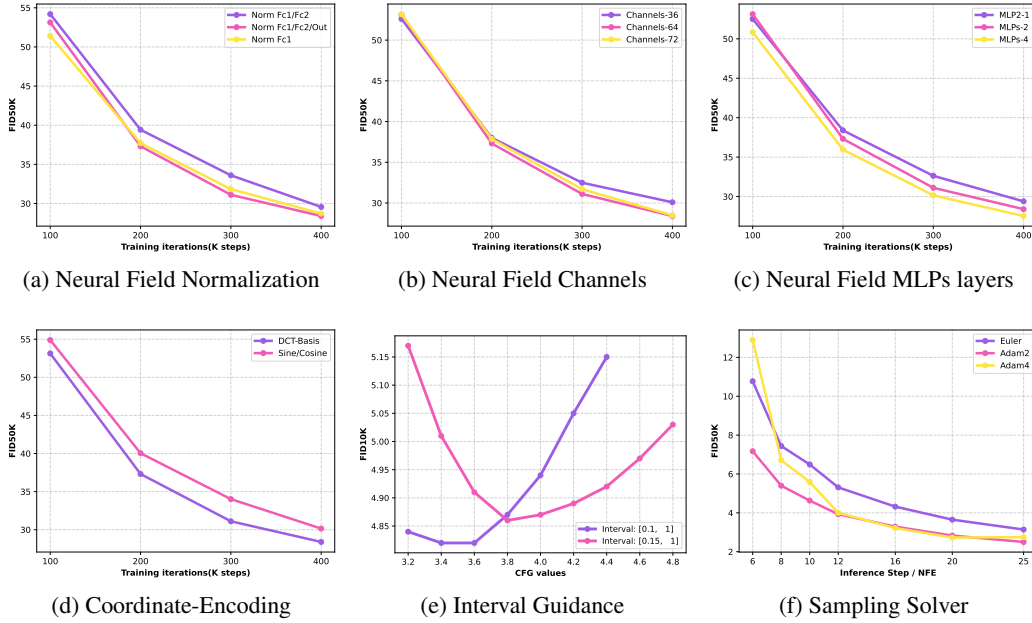


Figure 7: **Ablations studies of PixNerd.** We conduct ablation studies on class-to-image generation benchmark ImageNet $256 \times 256$  with PixNerd-L/16.

## 4.2 Neural Field Design

We conduct ablation studies on PixNerd-L/16, which comprises 22 transformer layers with 1024 channels. The Neural Field design is configured to have a computational burden comparable to that of a two-layer transformer block, so PixNerd-L/16 has inference latency similar to its counterpart, Baseline-L/16 (24 transformer layers with 1024 channels).

**Neural Field Normalization** As illustrated in Fig. 7a, we evaluate different neural field normalization strategies. Our baseline compares three approaches: (1) normalizing only the first weight (FC1), (2) normalizing both weight (FC1/FC2), and (3) our default strategy that additionally normalizes the output features. Experimental results demonstrate that the default strategy achieves optimal performance and convergence speed.

**Neural Field MLP channels** We conduct an empirical study of different MLP channel configurations (36, 64, and 72 channels) as shown in Fig. 7b. Our experiments reveal that: (1) a minimal configuration of 36 channels leads to noticeable performance degradation compared to 64 channels; (2) while the 72-channel variant achieves marginally better results, it incurs significant computational overhead, including slower training speed and increased parameter count. Based on this trade-off analysis, we select 64 channels as our default configuration.

	ImageNet 256×256						
	Params	Epochs	FID↓	sFID↓	IS↑	Pre.↑	Rec.↑
<i>Latent Generative Models</i>							
LDM-4 [1]	400M + 86M	170	3.6	-	247.7	0.87	0.48
DiT-XL [13]	675M + 86M	1400	2.27	4.60	278.2	0.83	0.57
SiT-XL [14]	675M + 86M	1400	2.06	4.50	270.3	0.82	0.59
FlowDCN [18]	618M + 86M	400	2.00	4.33	263.1	0.82	0.58
REPA [54]	675M + 86M	800	1.42	4.70	305.7	0.80	0.64
DDT-XL [15]	675M + 86M	400	1.26	4.51	310.6	0.79	0.65
MAR-L [58]	479M + 86M	800	1.78	-	296.0	0.81	0.60
CausalFusion [59]	676M + 86M	800	1.77	-	282.3	0.82	0.61
<i>Pixel Generative Models</i>							
ADM [8]	554M	400	4.59	5.13	186.7	0.82	0.52
RDM [10]	553M + 553M	/	1.99	3.99	260.4	0.81	0.58
JetFormer [60]	2.8B	/	6.64	-	-	0.69	0.56
FractalMAR-H [25]	844M	600	6.15	-	348.9	0.81	0.46
PixelFlow-XL/4 [9]	677M	320	1.98	5.83	282.1	0.81	0.60
<b>PixNerd-L/16 (Euler-50)</b>	458M	160	2.64	5.25	297	0.78	0.60
<b>PixNerd-XL/16 (Euler-50)</b>	700M	160	2.29	4.82	303	0.80	0.59
<b>PixNerd-XL/16 (Adam2-50)</b>	700M	160	2.16	4.93	291	0.78	0.60
<b>PixNerd-XL/16 (Euler-100)</b>	700M	160	2.15	4.55	297	0.79	0.59

Table 2: **System performance comparison** on ImageNet 256 × 256 class-conditioned generation. Our PixNerd-XL/16 achieves comparable results with latent diffusion models under similar computation demands while achieving much better results than other pixel space generative models. We adopt the interval guidance with interval [0.1, 1] and CFG of 3.5.

**Neural Field MLP Depth.** We investigate the impact of neural field depth by evaluating PixNerd-L/16 with 1, 2, and 4 MLP layers, as shown in Fig. 7c. Our experiments demonstrate consistent performance improvements with increasing network depth. However, considering the trade-off between computational efficiency (inference latency and training speed) and model performance, we establish the 2-layer configuration as our optimal default architecture.

**Coordinate Encodings** We compared the DCT-Basis coordinate encoding with traditional sine/cosine encoding in Fig. 7d. Our DCT-Basis encoding achieves much better results than sine/cosine encoding in terms of both convergence and final result.

$$\text{DCT-PE}(i, j) = \{\cos(k_1 i) \cos(k_2 j), \}_{k_1, k_2 \in (0, K]}. \quad (12)$$

### 4.3 Inference Scheduler design

**Interval Guidance** Classifier-free guidance [13, 14, 56] is a commonly used technique to improve the diffusion model performance. Interval guidance [57] is an improved cfg technique, and has been validated by recent works [2, 15]. We sweep different CFG values from 3.0 to 5.0 with a step of 0.2 to find the optimal CFG scheduler. As shown in Fig. 7e, our PixNerd-XL/16 achieves best result FID10k result with 3.4 or 3.6 within the interval [0.1, 1]. So we take 3.5 as the default CFG value.

**Sampling Solver** In Fig. 7f, we armed PixNerd-XL/16 with an Euler solver and an Adams-like linear multistep solver with 2/4 orders. The Adams-2 order solver consistently achieves better results than the Euler solver with limited sampling steps. Due to the learning difficulty of pixel spaces, Adams-4-order solver performs unstable compared to the Euler and Adams2 solvers. However, with sufficient sampling steps, their performance gap becomes marginal.

### 4.4 Class-to-Image Generation

**Training details** In class-to-image generation, to ensure a fair comparative analysis, we did not use gradient clipping and learning rate warm-up techniques. We adopt EMA with 0.9999 to stabilize the training. Our default training infrastructure consisted of 8 × A100 GPUs.



	ImageNet $512 \times 512$					
Model	Params	FID↓	sFID↓	IS↑	Pre.↑	Rec.↑
<i>Latent Diffusion Models</i>						
DiT-XL/2 [13]	675M + 86M	3.04	5.02	240.82	0.84	0.54
SiT-XL/2 [14]	675M + 86M	2.62	4.18	252.21	0.84	0.57
REPA-XL/2 [54]	675M + 86M	2.08	4.19	274.6	0.83	0.58
FlowDCN-XL/2 [18]	608M + 86M	2.44	4.53	252.8	0.84	0.54
EDM2 [29]	1.5B + 86M	1.81				
DDT-XL/2 [15]	675M + 86M	1.28	4.22	305.1	0.80	0.63
<i>Pixel Diffusion Models</i>						
ADM-G [8]	559M	7.72	6.57	172.71	0.87	0.42
ADM-G, ADM-U	559M	3.85	5.86	221.72	0.84	0.53
RIN [61]	320M	3.95	-	210	-	-
SimpleDiffusion[24]	2B	3.54	-	205	-	-
<b>PixNerd-XL/16 (Euler50)</b>	700M	3.41	6.43	246.45	0.80	0.58
<b>PixNerd-XL/16 (Euler100)</b>	700M	2.84	5.95	245.62	0.80	0.59

Table 3: **Benchmarking class-conditional image generation on ImageNet  $512 \times 512$ .** Our PixNerd-XL/16( $512 \times 512$ ) is fine-tuned from the same model trained on  $256 \times 256$  resolution. We adopt the interval guidance with interval  $[0.1, 1]$  and CFG of 3.5.

**Visualizations.** We placed the selected visual examples from PixNerd-XL/16 trained on ImageNet $256 \times 256$  and ImageNet $512 \times 512$  at Fig. 2. Our PixNerd-XL/16 can generate images with promising details. We generate these images with a CFG of 3.5 and the Euler-50 solver.

**ImageNet  $256 \times 256$  Benchmark.** We report the metrics of PixNerd-L/16 and PixNerd-XL/16 in Tab. 2. Under merely 160 training epochs, PixNerd-L/16 achieves 2.64 FID, significantly better than other pixel generative models like Jetformer [60], FractalMAR [25]. Further, PixNerd-XL/16 achieves 2.29 FID under 50 steps with the Euler solver. When used with Adams-2-order-solver(Adam2 for brevity), PixNerd-XL/16 achieves 2.16 FID, comparable to DiT. With enough sampling steps, PixNerd-XL/16 boosts FID to 2.15 under 100 steps. We adopt the interval guidance with interval  $[0.1, 1]$  and CFG of 3.5.

**ImageNet  $512 \times 512$  Benchmark.** We provide the final metrics of PixNerd-XL/16 at Tab. 3. To validate the superiority of our PixNerd model, we take our PixNerd-XL/16 trained on ImageNet  $256 \times 256$  as the initialization, fine-tune our PixNerd-XL/16 on ImageNet  $512 \times 512$  for *abc* steps. We adopt the aforementioned interval guidance [57] and we achieved 2.84 FID, with CFG of 3.5 within the time interval  $[0.3, 1.0]$ . PixNerd-XL/16 achieves comparable performance to other diffusion models.

#### 4.5 Text-to-Image Generation

**Data preprocess details** For text-to-image generation, we trained our model on a mixed dataset containing approximately 45M images from open-sourced datasets, e.g., SAM [66], JourneyDB [67], ImageNet-1K [68]. We recaption all the images with Qwen2.5-VL-7B [69] to yield English caption descriptions of various lengths. Note that our caption results only contain English descriptions. All the images are cropped into a square shape of  $256 \times 256$  or  $512 \times 512$ , we do not adopt various aspect ratio training. We leave the native resolution [70] or native aspect training [36, 37, 38] as future works.

**Training details** We adopt Qwen3-1.7B<sup>1</sup> as the text encoder. To improve the alignment of frozen text features [63], we jointly train several transformer layers on the frozen text features similar to Fluid [63]. To further enhance the generation quality, we adopt an SFT stage at resolution  $512 \times 512$  with the dataset released by BLIP-3o [71]. The total batch size is 1536 for  $256 \times 256$  resolution pretraining and 512 for  $512 \times 512$  resolution pretraining. We trained PixNerd on  $256 \times 256$  resolution for 200K steps and trained on  $512 \times 512$  resolution for 80K steps. The default training infrastructure consisted of  $16 \times$  A100 GPUs. We adopt the gradient clip to stabilize training. We

<sup>1</sup><https://huggingface.co/Qwen/Qwen3-1.7B>

	#Params	Sin.Obj.	Two.Obj.	Counting	Colors	Pos	Color.Attr.	Overall
<i>autoregressive model</i>								
Show-o [62]	1.3B	0.95	0.52	0.49	0.82	0.11	0.28	0.53
MAR [63]	1.1B+4.7B+86M	0.96	0.77	0.61	0.78	0.34	0.53	0.67
SimpleAR [64]	0.5B	-	0.82	-	-	0.26	0.38	0.59
SimpleAR [64]	1.5B	-	0.90	-	-	0.26	0.45	0.63
<i>latent diffusion model</i>								
LDM[1]	1.4B	0.92	0.29	0.23	0.70	0.02	0.05	0.37
DALL-E 2	4.2B	0.94	0.66	0.49	0.77	0.10	0.19	0.52
DALL-E 3	-	0.96	0.87	0.47	0.83	0.43	0.45	0.67
Imagen	3B	-	-	-	-	-	-	-
SD3 [53]	8B	0.98	0.84	0.66	0.74	0.40	0.43	0.68
Transfusion [65]	7.3B	-	-	-	-	-	-	0.63
<i>pixel diffusion model</i>								
PixelFlow-XL/4 [9]	882M + 3B	-	-	-	-	-	-	0.60
PixelFlow-XL/4 <sup>†</sup> [9]	882M + 3B	-	-	-	-	-	-	0.64
PixelNerd-XXL/16	1.2B + 1.7B	0.97	0.86	0.44	0.83	0.71	0.53	0.73

Table 4: **Comparison with other text-to-image models on GenEval Benchmark.**<sup>†</sup> indicates prompt rewriting. Parameters consist of denoiser+text encoder+vae. Our PixNerd-XXL/16 achieves competitive performance compared with others under a much-limited data scale (45M images).

Model	#Params	DPG Benchmark					
		Global	Entity	Attribute	Relation	Other	Average
<i>latent diffusion model</i>							
SD v2	0.86B + 0.7B + 86M	77.67	78.13	74.91	80.72	80.66	68.09
PixArt- $\alpha$	0.61B + 4.7B + 86M	74.97	79.32	78.60	82.57	76.96	71.11
Playground v2	2.61B + 0.7B + 86M	83.61	79.91	82.67	80.62	81.22	74.54
DALL-E 3	-	90.97	89.61	88.39	90.58	89.83	83.50
SD v1.5	0.86B + 0.7B + 86M	74.63	74.23	75.39	73.49	67.81	63.18
SDXL	2.61B + 1.4B + 86M	83.27	82.43	80.91	86.76	80.41	74.65
<i>pixel diffusion model</i>							
PixelFlow-XL/4	882M + 3B	-	-	-	-	-	77.90
PixelNerd-XXL/16	1.2B + 1.7B	80.5	87.9	87.2	91.3	72.8	80.9

Table 5: **Comparison with other text-to-image models on DPG Benchmark.** Parameters consist of denoiser+text encoder+vae. Our PixNerd-XXL/16 achieves competitive performance compared with others under a much-limited data scale (45M images).

adopted *torch.compile* to optimize the computation graph to reduce the memory and computation overhead. We use the Adams-2nd solver with 25 steps as the default choice for sampling.

**Visualizations.** We provided  $512 \times 512$  visualizations with prompts provided by DouBao-1.5-Pro<sup>2</sup> in Fig. 3. As illustrated in Fig. 3, our PixNerd-XXL/16 is capable of generating visually compelling images from complex text prompts. Overall, the atmosphere and color tones are largely accurate. Noted that we only trained PixNerd with English prompts. As shown in Fig. 8, we can generate images even with other languages like Chinese and Japanese thanks to the powerful embedding space of Qwen3 models [69]. Nevertheless, occasional blurry or unnatural artifacts appear in certain scenarios (e.g., steampunk lab image). We posit that appropriate post-training processing could mitigate such artifacts [72], and we intend to explore pixel-space post-training as future work.

**Sampling solvers.** We provide denoising trajectories in  $x_1$  space (clean data) of different solvers in Fig. 8, including Euler, Adams-2nd, and Adams-3rd solvers. Adams-2nd solver achieves stable and fast sampling results.

<sup>2</sup><https://www.volcengine.com/product/doubao>

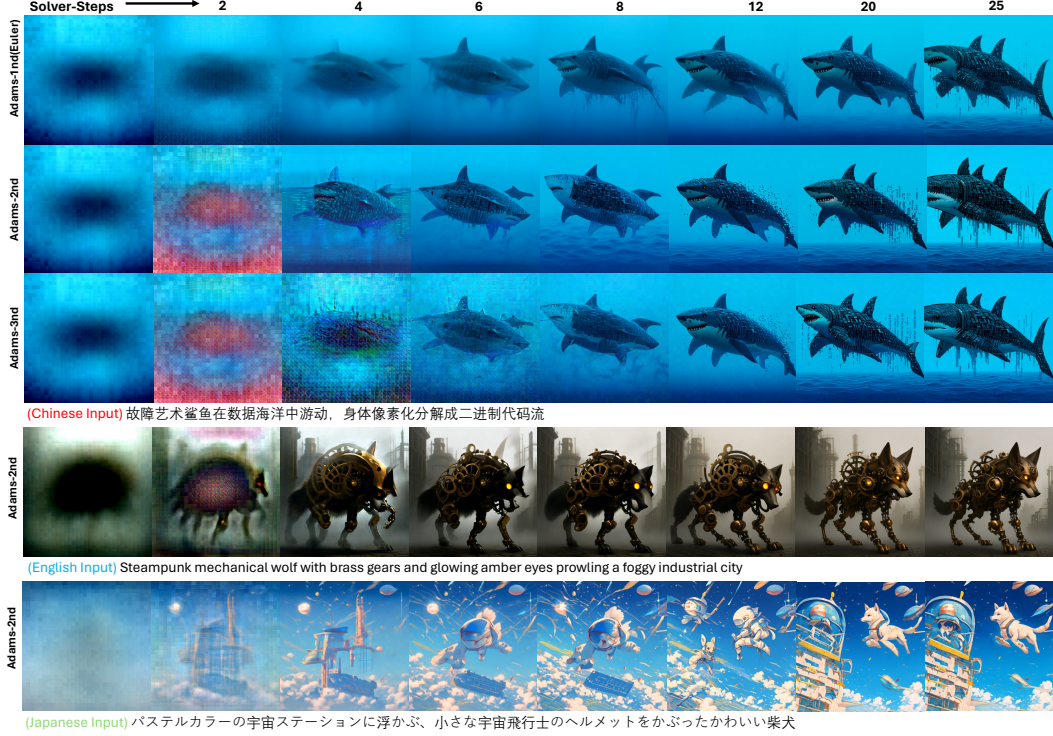


Figure 8: **The Text-to-Image  $512 \times 512$  visualization with different solvers.** We armed PixNerd with different ODE solvers, eg, Euler, Adams-2nd, Adams-3rd. Adams solver achieves better visual quality than the naive Euler solver. Also, thanks to the powerful text embedding in Qwen3 models, though we only trained PixNerd with English captions, PixNerd can generate samples of promising quality with other languages.

**Training-free arbitrary resolution generation.** As shown in Fig. 4, without any resolution adaptation fine-tuning, we can achieve arbitrary resolution generation through the coordinate interpolation while keeping the amount of tokens as constant as the pretraining resolution. Specifically, we sampled pretraining resolution from the given noisy image, then fed this sampled version into the transformer. This keeps the token amount in our PixNerd as consistent as in the pretraining stage. To match the velocity field with the desired resolution of the given noisy image, we then interpolate the coordinates for neural field decoding accordingly.

**GenEval Benchmark** We provided quantity comparison on Geneval [51] benchmark in Tab. 4. Our PixNerd-XXL/16 achieves comparable results under enormous patch sizes and limited data scales. As shown in Tab. 4, our PixNerd-XXL/16 achieves 0.73 overall score, beating pixelflow [9] with a significant margin.

**DPG Benchmark** We provided quantity comparison on DPG [52] benchmark in Tab. 5. Our PixNerd-XXL/16 achieves competitive results compared to its latent counterparts. As shown in Tab. 5, our PixNerd-XXL/16 achieves 0.82 overall score, beating other pixel generation models with a significant margin.

## 5 Discussion

Several related works [41, 39, 73] also combine diffusion with neural fields. In practice, however, they differ fundamentally from PixNerd. For example, INFd [41] and DDMI [39] leverage neural fields to enhance VAEs rather than diffusion models, and their generative capacity still stems from a latent diffusion model. DenoisedWeights [74] trains independent neural weights for each image before training a generative model on these pre-collected weights. This remains a two-stage framework and poses non-trivial challenges for large-scale training. INRFLOW [73] and PatchDiffusion[75] utilize coordinate encodings to enhance diffusion model performance. Beyond diffusion-based generative

models, GAN-based methods [76, 77, 40, 78] also utilize neural fields or coordinate encodings. PixNerd is a simple yet elegant single-stage pixel-space generative model that does not rely on a VAE. Current latent generative models inevitably cascade errors due to their two-stage configurations. Further, a high-quality VAE usually demands numerous losses supervisions, e.g., adversarial loss, LPIPS loss. In particular, adversarial loss is unstable in training and tends to introduce artifacts. Pixel generative model has more potential in the future, and PixNerd is a simple yet elegant solution for a pixel-space generative model.

## 6 Conclusion

In this paper, we return to pixel space diffusion with neural field. We present a single-scale, single-stage, efficient, end-to-end solution, the pixel neural field diffusion (PixelNerd). We achieved 2.15 FID on ImageNet  $256 \times 256$  and 2.84 FID on ImageNet  $512 \times 512$  without any complex cascade pipeline. Our PixNerd-XXL/16 achieved a competitive 0.73 overall score on the GenEval benchmark and 80.9 average score on DPG benchmark. However, current PixNerd shows unclear details in some cases, as Fig. 3 and still has gaps with its latent counterparts.

## References

- [1] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 3, 4, 5, 8, 10
- [2] Jingfeng Yao and Xinggang Wang. Reconstruction vs. generation: Taming optimization dilemma in latent diffusion models. *arXiv preprint arXiv:2501.01423*, 2025. 1, 3, 4, 5, 6, 8
- [3] Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang, Haotian Tang, Muyang Li, Yao Lu, and Song Han. Deep compression autoencoder for efficient high-resolution diffusion models. *arXiv preprint arXiv:2410.10733*, 2024. 1, 3, 4, 5
- [4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 3
- [5] Axel Sauer, Katja Schwarz, and Andreas Geiger. Stylegan-xl: Scaling stylegan to large diverse datasets. In *ACM SIGGRAPH 2022 conference proceedings*, pages 1–10, 2022. 3
- [6] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 3
- [7] Xingjian Leng, Jaskirat Singh, Yunzhong Hou, Zhenchang Xing, Saining Xie, and Liang Zheng. Repa-e: Unlocking vae for end-to-end tuning with latent diffusion transformers. *arXiv preprint arXiv:2504.10483*, 2025. 3
- [8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 3, 4, 5, 6, 8, 9
- [9] Shoufa Chen, Chongjian Ge, Shilong Zhang, Peize Sun, and Ping Luo. Pixelflow: Pixel-space generative models with flow. *arXiv preprint arXiv:2504.07963*, 2025. 3, 4, 6, 8, 10, 11
- [10] Jiayan Teng, Wendi Zheng, Ming Ding, Wenyi Hong, Jianqiao Wangni, Zhuoyi Yang, and Jie Tang. Relay diffusion: Unifying diffusion process across resolutions for image synthesis. *arXiv preprint arXiv:2309.03350*, 2023. 3, 4, 8
- [11] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 3, 5
- [12] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in neural information processing systems*, 33:7462–7473, 2020. 3, 5
- [13] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 4, 5, 8, 9



- [14] Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. *arXiv preprint arXiv:2401.08740*, 2024. 4, 5, 8, 9
- [15] Shuai Wang, Zhi Tian, Weilin Huang, and Limin Wang. Decoupled diffusion transformer. *arXiv preprint arXiv:2504.05741*, 2025. 4, 5, 6, 8, 9
- [16] Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing and improving the training dynamics of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24174–24184, 2024. 4
- [17] Xiaoyu Yue, Zidong Wang, Zeyu Lu, Shuyang Sun, Meng Wei, Wanli Ouyang, Lei Bai, and Luping Zhou. Diffusion models need visual priors for image generation. *arXiv preprint arXiv:2410.08531*, 2024. 4
- [18] Shuai Wang, Zexian Li, Tianhui Song, Xubin Li, Tiezheng Ge, Bo Zheng, and Limin Wang. Exploring dcn-like architecture for fast image generation with arbitrary resolution. *Advances in Neural Information Processing Systems*, 37:87959–87977, 2024. 4, 8, 9
- [19] Yao Teng, Yue Wu, Han Shi, Xuefei Ning, Guohao Dai, Yu Wang, Zhenguo Li, and Xihui Liu. Dim: Diffusion mamba for efficient high-resolution image synthesis. *arXiv preprint arXiv:2405.14224*, 2024. 4
- [20] Tianhui Song, Weixin Feng, Shuai Wang, Xubin Li, Tiezheng Ge, Bo Zheng, and Limin Wang. Dmm: Building a versatile image generation model via distillation-based model merging. *arXiv preprint arXiv:2504.12364*, 2025. 4
- [21] Mingyuan Zhou, Huangjie Zheng, Zhendong Wang, Mingzhang Yin, and Hai Huang. Score identity distillation: Exponentially fast distillation of pretrained diffusion models for one-step generation. In *Forty-first International Conference on Machine Learning*, 2024. 4
- [22] Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22669–22679, 2023. 4, 5
- [23] Diederik Kingma and Ruiqi Gao. Understanding diffusion objectives as the elbo with simple data augmentation. *Advances in Neural Information Processing Systems*, 36:65484–65516, 2023. 4
- [24] Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for high resolution images. In *International Conference on Machine Learning*, pages 13213–13232. PMLR, 2023. 4, 9
- [25] Tianhong Li, Qinyi Sun, Lijie Fan, and Kaiming He. Fractal generative models. *arXiv preprint arXiv:2502.17437*, 2025. 4, 8, 9
- [26] Shuangfei Zhai, Ruixiang Zhang, Preetum Nakkiran, David Berthelot, Jiatao Gu, Huangjie Zheng, Tianrong Chen, Miguel Angel Bautista, Navdeep Jaitly, and Josh Susskind. Normalizing flows are capable generative models. *arXiv preprint arXiv:2412.06329*, 2024. 4
- [27] Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982. 4
- [28] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 4
- [29] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577, 2022. 4, 9
- [30] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 4
- [31] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 5, 6
- [32] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 5, 6

- [33] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. 5, 6
- [34] Xiangxiang Chu, Jianlin Su, Bo Zhang, and Chunhua Shen. Visionllama: A unified llama interface for vision tasks. *arXiv preprint arXiv:2403.00522*, 2024. 5, 6
- [35] Zeyu Lu, Zidong Wang, Di Huang, Chengyue Wu, Xihui Liu, Wanli Ouyang, and Lei Bai. Fit: Flexible vision transformer for diffusion model. *arXiv preprint arXiv:2402.12376*, 2024. 5, 6
- [36] Lixue Gong, Xiaoxia Hou, Fanshi Li, Liang Li, Xiaochen Lian, Fei Liu, Liyang Liu, Wei Liu, Wei Lu, Yichun Shi, et al. Seedream 2.0: A native chinese-english bilingual image generation foundation model. *arXiv preprint arXiv:2503.07703*, 2025. 5, 9
- [37] Yu Gao, Lixue Gong, Qiushan Guo, Xiaoxia Hou, Zhichao Lai, Fanshi Li, Liang Li, Xiaochen Lian, Chao Liao, Liyang Liu, et al. Seedream 3.0 technical report. *arXiv preprint arXiv:2504.11346*, 2025. 5, 9
- [38] Chao Liao, Liyang Liu, Xun Wang, Zhengxiong Luo, Xinyu Zhang, Wenliang Zhao, Jie Wu, Liang Li, Zhi Tian, and Weilin Huang. Mogao: An omni foundation model for interleaved multi-modal generation. *arXiv preprint arXiv:2505.05472*, 2025. 5, 9
- [39] Dogyun Park, Sihyeon Kim, Sojin Lee, and Hyunwoo J Kim. Ddmi: Domain-agnostic latent diffusion models for synthesizing high-quality implicit neural representations. *arXiv preprint arXiv:2401.12517*, 2024. 5, 11
- [40] Chieh Hubert Lin, Chia-Che Chang, Yu-Sheng Chen, Da-Cheng Juan, Wei Wei, and Hwann-Tzong Chen. Coco-gan: Generation by parts via conditional coordinating. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4512–4521, 2019. 5, 12
- [41] Yinbo Chen, Oliver Wang, Richard Zhang, Eli Shechtman, Xiaolong Wang, and Michael Gharbi. Image neural field diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8007–8017, 2024. 5, 11
- [42] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5855–5864, 2021. 5
- [43] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751*, 2019. 5
- [44] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021. 5
- [45] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural information processing systems*, 35:25018–25032, 2022. 5
- [46] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3295–3306, 2023. 5
- [47] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6
- [48] Charlie Nash, Jacob Menick, Sander Dieleman, and Peter W Battaglia. Generating images with sparse representations. *arXiv preprint arXiv:2103.03841*, 2021. 6
- [49] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016. 6
- [50] Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *Advances in neural information processing systems*, 32, 2019. 6
- [51] Dhruva Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023. 6, 11

- [52] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024. 6, 11
- [53] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. *arXiv preprint arXiv:2403.03206*, 2024. 6, 10
- [54] Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and Saining Xie. Representation alignment for generation: Training diffusion transformers is easier than you think. *arXiv preprint arXiv:2410.06940*, 2024. 6, 8, 9
- [55] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 6
- [56] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 8
- [57] Tuomas Kynkäänniemi, Miika Aittala, Tero Karras, Samuli Laine, Timo Aila, and Jaakko Lehtinen. Applying guidance in a limited interval improves sample and distribution quality in diffusion models. *arXiv preprint arXiv:2404.07724*, 2024. 8, 9
- [58] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *Advances in Neural Information Processing Systems*, 37:56424–56445, 2025. 8
- [59] Chaorui Deng, Deyao Zh, Kunchang Li, Shi Guan, and Haoqi Fan. Causal diffusion transformers for generative modeling. *arXiv preprint arXiv:2412.12095*, 2024. 8
- [60] Michael Tschannen, André Susano Pinto, and Alexander Kolesnikov. Jetformer: An autoregressive generative model of raw images and text. *arXiv preprint arXiv:2411.19722*, 2024. 8, 9
- [61] Allan Jabri, David Fleet, and Ting Chen. Scalable adaptive computation for iterative generation. *arXiv preprint arXiv:2212.11972*, 2022. 9
- [62] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024. 10
- [63] Lijie Fan, Tianhong Li, Siyang Qin, Yuanzhen Li, Chen Sun, Michael Rubinstein, Deqing Sun, Kaiming He, and Yonglong Tian. Fluid: Scaling autoregressive text-to-image generative models with continuous tokens. *arXiv preprint arXiv:2410.13863*, 2024. 9, 10
- [64] Junke Wang, Zhi Tian, Xun Wang, Xinyu Zhang, Weilin Huang, Zuxuan Wu, and Yu-Gang Jiang. Simplear: Pushing the frontier of autoregressive visual generation through pretraining, sft, and rl. *arXiv preprint arXiv:2504.11455*, 2025. 10
- [65] Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model. *arXiv preprint arXiv:2408.11039*, 2024. 10
- [66] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 9
- [67] Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, et al. Journeydb: A benchmark for generative image understanding. *Advances in neural information processing systems*, 36:49659–49678, 2023. 9
- [68] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015. 9
- [69] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 9, 10
- [70] Zidong Wang, Lei Bai, Xiangyu Yue, Wanli Ouyang, and Yiyuan Zhang. Native-resolution image synthesis. *arXiv preprint arXiv:2506.03131*, 2025. 9

- [71] Jiuhai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025. 9
- [72] Yuxi Ren, Xin Xia, Yanzuo Lu, Jiacheng Zhang, Jie Wu, Pan Xie, Xing Wang, and Xuefeng Xiao. Hyper-sd: Trajectory segmented consistency model for efficient image synthesis. *arXiv preprint arXiv:2404.13686*, 2024. 10
- [73] Yuyang Wang, Anurag Ranjan, Josh Susskind, and Miguel Angel Bautista. Inrflow: Flow matching for inrs in ambient space. *arXiv preprint arXiv:2412.03791*, 2024. 11
- [74] Yifan Gong, Zheng Zhan, Yanyu Li, Yerlan Idelbayev, Andrey Zharkov, Kfir Aberman, Sergey Tulyakov, Yanzhi Wang, and Jian Ren. Efficient training with denoised neural weights. In *European Conference on Computer Vision*, pages 18–34. Springer, 2024. 11
- [75] Zhendong Wang, Yifan Jiang, Huangjie Zheng, Peihao Wang, Pengcheng He, Zhangyang Wang, Weizhu Chen, Mingyuan Zhou, et al. Patch diffusion: Faster and more data-efficient training of diffusion models. *Advances in neural information processing systems*, 36:72137–72154, 2023. 11
- [76] Evangelos Ntavelis, Mohamad Shahbazi, Iason Kastanis, Radu Timofte, Martin Danelljan, and Luc Van Gool. Arbitrary-scale image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11533–11542, 2022. 12
- [77] Ivan Skorokhodov, Savva Ignatyev, and Mohamed Elhoseiny. Adversarial generation of continuous images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10753–10764, 2021. 12
- [78] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in neural information processing systems*, 34:852–863, 2021. 12