

All are Worth Words: A ViT Backbone for Diffusion Models

Fan Bao¹, Shen Nie², Kaiwen Xue², Yue Cao³, Chongxuan Li^{2*}, Hang Su¹, Jun Zhu^{1*}

¹Dept. of Comp. Sci. & Tech., Institute for AI, BNRist Center

¹Tsinghua-Bosch Joint ML Center, THBI Lab, Tsinghua University, Beijing, 100084 China

²Gaoling School of Artificial Intelligence, Renmin University of China,

²Beijing Key Laboratory of Big Data Management and Analysis Methods, Beijing, China

³Beijing Academy of Artificial Intelligence

bf19@mails.tsinghua.edu.cn; nieshen@ruc.edu.cn; {kevin.kaiwenxue, caoyue10}@gmail.com

chongxuanli@ruc.edu.cn; {suhangss, dcszj}@tsinghua.edu.cn

Abstract

Vision transformers (ViT) have shown promise in various vision tasks while the U-Net based on a convolutional neural network (CNN) remains dominant in diffusion models. We design a simple and general ViT-based architecture (named U-ViT) for image generation with diffusion models. U-ViT is characterized by treating all inputs including the time, condition and noisy image patches as tokens and employing long skip connections between shallow and deep layers. We evaluate U-ViT in unconditional and class-conditional image generation, as well as text-to-image generation tasks, where U-ViT is comparable if not superior to a CNN-based U-Net of a similar size. In particular, latent diffusion models with U-ViT achieve record-breaking FID scores of 2.29 in class-conditional image generation on ImageNet 256×256 , and 5.48 in text-to-image generation on MS-COCO, among methods without accessing large external datasets during the training of generative models.

Our results suggest that, for diffusion-based image modeling, the long skip connection is crucial while the down-sampling and up-sampling operators in CNN-based U-Net are not always necessary. We believe that U-ViT can provide insights for future research on backbones in diffusion models and benefit generative modeling on large scale cross-modality datasets. Our code is available at <https://github.com/baofff/U-ViT>.

1. Introduction

Diffusion models [25, 62, 67] are powerful deep generative models that emerge recently for high quality image generation [13, 26, 54]. They grow rapidly and find applications in text-to-image generation [52, 54, 56], image-

*Corresponding to C. Li and J. Zhu.

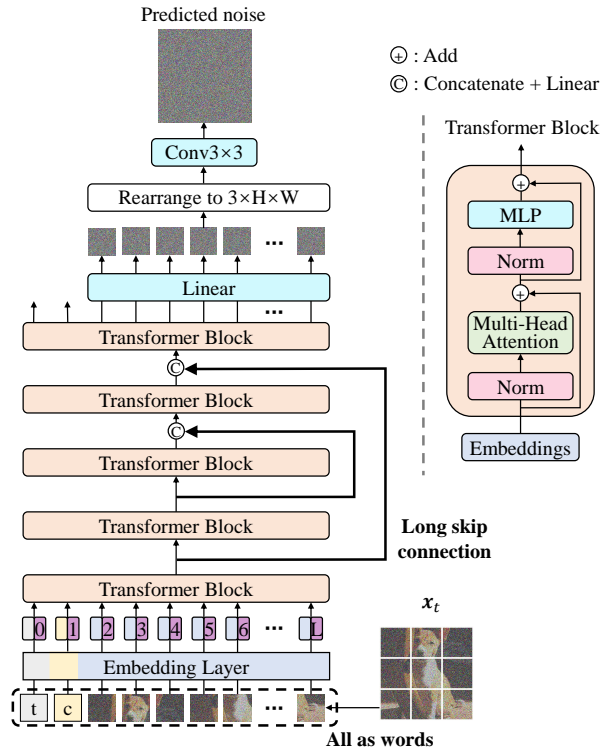


Figure 1. The U-ViT architecture for diffusion models, which is characterized by treating **all** inputs including the time, condition and noisy image patches as **tokens** and employing $(\#Blocks-1)/2$ **long skip connections** between shallow and deep layers.

to-image generation [10, 45, 80], video generation [24, 28], speech synthesis [6, 34], and 3D synthesis [50].

Along with the development of algorithms [2, 3, 15, 25, 33, 43, 44, 49, 63, 64, 67, 71], the revolution of backbones plays a central role in diffusion models. A representative example is U-Net based on a convolutional neural network

(CNN) employed in prior work [25, 65]. The CNN-based U-Net is characterized by a group of down-sampling blocks, a group of up-sampling blocks, and long skip connections between the two groups, which dominates diffusion models for image generation tasks [13, 52, 54, 56]. On the other hand, vision transformers (ViT) [16] have shown promise in various vision tasks, where ViT is comparable or even superior to CNN based approaches [9, 21, 38, 68, 81]. Therefore, a very natural question arises: *whether the reliance of the CNN-based U-Net is necessary in diffusion models?*

In this paper, we design a simple and general ViT-based architecture called U-ViT (Figure 1). Following the design methodology of transformers, U-ViT treats all inputs including the time, condition and noisy image patches as tokens. Crucially, U-ViT employs long skip connections between shallow and deep layers inspired by U-Net. Intuitively, low-level features are important to the pixel-level prediction objective in diffusion models and such connections can ease the training of the corresponding prediction network. Besides, U-ViT optionally adds an extra 3×3 convolutional block before output for better visual quality. See a systematical ablation study for all elements in Figure 2.

We evaluate U-ViT in three popular tasks: unconditional image generation, class-conditional image generation and text-to-image generation. In all settings, U-ViT is comparable if not superior to a CNN-based U-Net of a similar size. In particular, latent diffusion models with U-ViT achieve record-breaking FID scores of 2.29 in class-conditional image generation on ImageNet 256×256 , and 5.48 in text-to-image generation on MS-COCO, among methods without accessing large external datasets during the training of generative models.

Our results suggest that the long skip connection is crucial while the down/up-sampling operators in CNN-based U-Net are not always necessary for image diffusion models. We believe that U-ViT can provide insights for future research on diffusion model backbones and benefit generative modeling on large scale cross-modality datasets.

2. Background

Diffusion models [25, 62, 67] gradually inject noise to data, and then reverse this process to generate data from noise. The noise-injection process, also called the forward process, is formalized as a Markov chain:

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}).$$

Here \mathbf{x}_0 is the data, $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t|\sqrt{\alpha_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$, and α_t and β_t represent the noise schedule such that $\alpha_t + \beta_t = 1$. To reverse this process, a Gaussian model $p(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\mu}_t(\mathbf{x}_t), \sigma_t^2\mathbf{I})$ is adopted to approximate the ground truth reverse transition $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$, and

the optimal mean [3] is

$$\boldsymbol{\mu}_t^*(\mathbf{x}_t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \mathbb{E}[\boldsymbol{\epsilon}|\mathbf{x}_t] \right).$$

Here $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$, and $\boldsymbol{\epsilon}$ is the standard Gaussian noises injected to \mathbf{x}_t . Thus, the learning is equivalent to a noise prediction task. Formally, a noise prediction network $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ is adopted to learn $\mathbb{E}[\boldsymbol{\epsilon}|\mathbf{x}_t]$ by minimizing a noise prediction objective, i.e., $\min_{\theta} \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|_2^2$, where t is uniform between 1 and T . To learn conditional diffusion models, e.g., class-conditional [13] or text-to-image [52] models, the condition information is further fed into the noise prediction objective:

$$\min_{\theta} \mathbb{E}_{t, \mathbf{x}_0, c, \boldsymbol{\epsilon}} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, c)\|_2^2, \quad (1)$$

where c is the condition or its continuous embedding. In prior work on image modeling, the success of diffusion models heavily rely on CNN-based U-Net [55, 65], which is a convolutional backbone characterized by a group of down-sampling blocks, a group of up-sampling blocks and long skip connections between the two groups, and c is fed into U-Net by mechanisms such as adaptive group normalization [13] and cross attention [54].

Vision Transformer (ViT) [16] is a pure transformer architecture that treats an image as a sequence of tokens (words). ViT rearranges an image into a sequence of flattened patches. Then, ViT adds learnable 1D position embeddings to linear embeddings of these patches before feeding them into a transformer encoder [72]. ViT has shown promise in various vision tasks but it is not clear whether it is suitable for diffusion-based image modeling yet.

3. Method

U-ViT is a simple and general backbone for diffusion models in image generation (Figure 1). In particular, U-ViT parameterizes the noise prediction network¹ $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, c)$ in Eq. (1). It takes the time t , the condition c and the noisy image \mathbf{x}_t as inputs and predicts the noise injected into \mathbf{x}_t . Following the design methodology of ViT, the image is split into patches, and U-ViT treats all inputs including the time, condition and image patches as tokens (words).

Inspired by the success of the CNN-based U-Net in diffusion models [65], U-ViT also employs similar long skip connections between shallow and deep layers. Intuitively, the objective in Eq. (1) is a pixel-level prediction task and is sensitive to low-level features. The long skip connections provide shortcuts for the low-level features and therefore ease the training of the noise prediction network.

¹U-ViT can also parameterize other types of prediction, e.g., \mathbf{x}_0 -prediction [25].

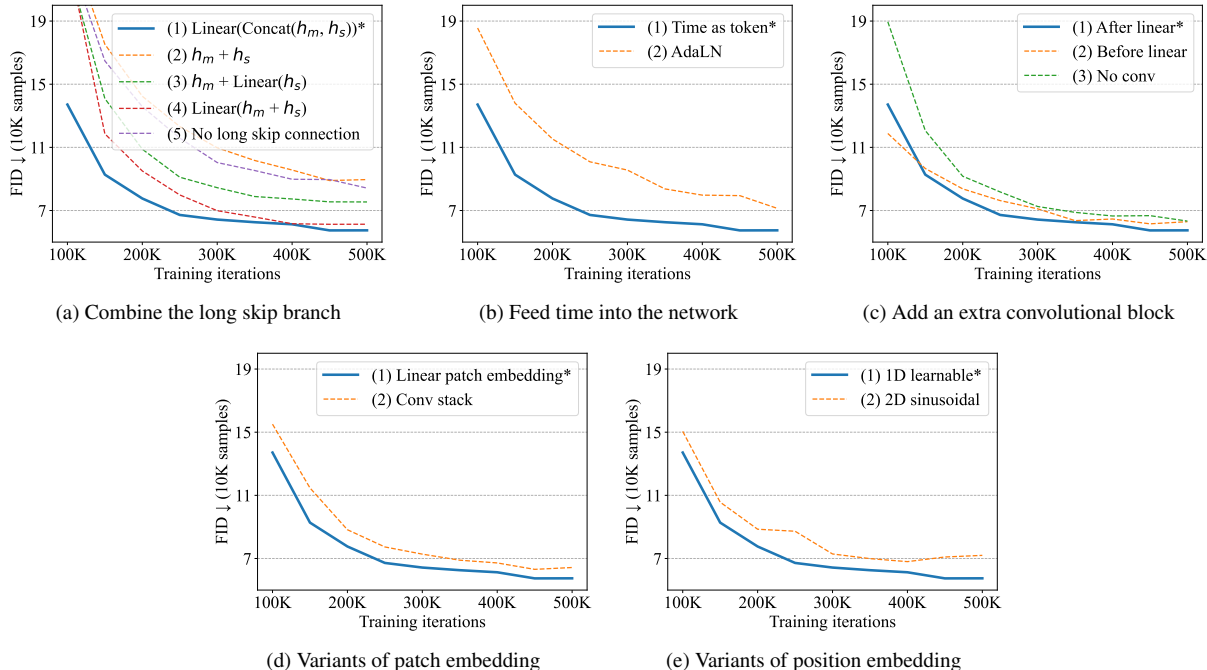


Figure 2. Ablate design choices. The one marked with * is the adopted choice of U-ViT illustrated in Figure 1. Since this ablation aims to determine implementation details, we evaluate FID on 10K generated samples (instead of 50K samples for efficiency).

Additionally, U-ViT optionally adds a 3×3 convolutional block before output. This is intended to prevent the potential artifacts in images produced by transformers [78]. The block improves the visual quality of the samples generated by U-ViT according to our experiments.

In Section 3.1, we present the implementation details of U-ViT. In Section 3.2, we present the scaling properties of U-ViT by studying the effect of depth, width and patch size.

3.1. Implementation Details

Although U-ViT is conceptually simple, we carefully design its implementation. To this end, we perform a systematical empirical study on key elements in U-ViT. In particular, we ablate on CIFAR10 [36], evaluate the FID score [23] every 50K training iterations on 10K generated samples (instead of 50K samples for efficiency), and determine default implementation details.

The way to combine the long skip branch. Let $\mathbf{h}_m, \mathbf{h}_s \in \mathbb{R}^{L \times D}$ be the embeddings from the main branch and the long skip branch respectively. We consider several ways to combine them before feeding them to the next transformer block: (1) concatenating them and then performing a linear projection as illustrated in Figure 1, i.e., $\text{Linear}(\text{Concat}(\mathbf{h}_m, \mathbf{h}_s))$; (2) directly adding them, i.e., $\mathbf{h}_m + \mathbf{h}_s$; (3) performing a linear projection to \mathbf{h}_s and then adding them, i.e., $\mathbf{h}_m + \text{Linear}(\mathbf{h}_s)$; (4) adding them and then performing a linear projection, i.e., $\text{Linear}(\mathbf{h}_m + \mathbf{h}_s)$. (5) We also compare with the case where the long skip

connection is dropped. As shown in Figure 2 (a), directly adding $\mathbf{h}_m, \mathbf{h}_s$ does not provide benefits. Since a transformer block has skip connections via the adding operator inside it, \mathbf{h}_m already contains information of \mathbf{h}_s in a linear form. As a result, the only effect of $\mathbf{h}_m + \mathbf{h}_s$ is to increase the coefficient of \mathbf{h}_s in the linear form, which does not change the nature of the network. In contrast, all other ways to combine \mathbf{h}_s perform a linear projection on \mathbf{h}_s and improve the performance compared to no long skip connection. Among them, the first way with concatenation performs the best. In Appendix D, we visualize the similarity between representations in a network, and we find the first way with concatenation significantly changes the representations, which validates its effectiveness.

The way to feed the time into the network. We consider two ways to feed t into the network. (1) The first way is to treat it as a token as illustrated in Figure 1. (2) The second way is to incorporate the time after the layer normalization in the transformer block [20], which is similar to the adaptive group normalization [13] used in U-Net. The second way is referred to as adaptive layer normalization (AdaLN). Formally, $\text{AdaLN}(h, y) = y_s \text{LayerNorm}(h) + y_b$, where h is an embedding inside a transformer block, and y_s, y_b are obtained from a linear projection of the time embedding. As shown in Figure 2 (b), while simple, the first way that treats time as a token performs better than AdaLN.

The way to add an extra convolutional block after the transformer. We consider two ways to add an extra convo-

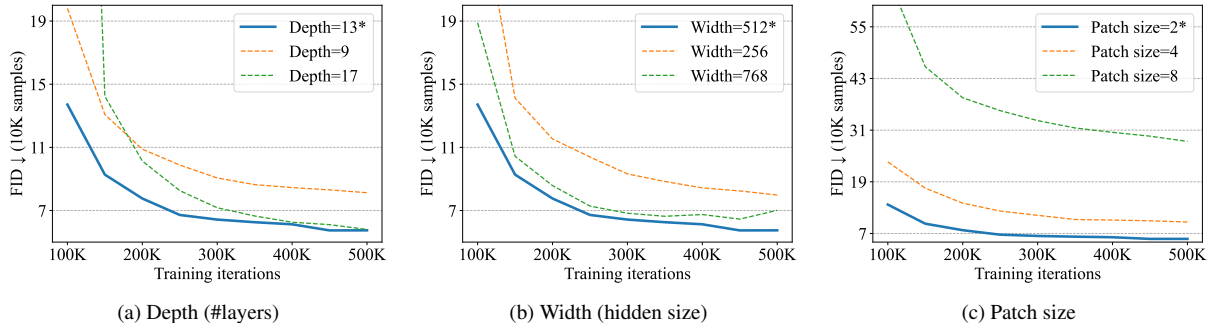


Figure 3. Effect of depth, width and patch size. The one marked with * corresponds to the setting of U-ViT-S/2 (see Table 2).

lutional block after the transformer. (1) The first way is to add a 3×3 convolutional block after the linear projection that maps the token embeddings to image patches, as illustrated in Figure 1. (2) The second way is to add a 3×3 convolutional block before this linear projection, which needs to first rearrange the 1D sequence of token embeddings $\mathbf{h} \in \mathbb{R}^{L \times D}$ to a 2D feature of shape $H/P \times W/P \times D$, where P is the patch size. (3) We also compare with the case where we drop the extra convolutional block. As shown in Figure 2 (c), the first way that adds a 3×3 convolutional block after the linear projection performs slightly better than other two choices.

Variants of the patch embedding. We consider two variants of the patch embedding. (1) The original patch embedding adopts a linear projection that maps a patch to a token embedding, as illustrated in Figure 1. (2) Alternatively, [73] use a stack of 3×3 convolutional blocks followed by a 1×1 convolutional block to map an image to token embeddings. We compare them in Figure 2 (d), and the original patch embedding performs better.

Variants of the position embedding. We consider two variants of the position embedding. (1) The first one is the 1-dimensional learnable position embedding proposed in the original ViT [16], which is the default setting in this paper. (2) The second one is the 2-dimensional sinusoidal position embedding, which is obtained by concatenating the sinusoidal embeddings [72] of i and j for a patch at position (i, j) . As shown in Figure 2 (e), the 1-dimensional learnable position embedding performs better. We also try not use any position embedding, and find the model fails to generate meaningful images, which implies the position information is critical in image generation.

3.2. Effect of Depth, Width and Patch Size

We present scaling properties of U-ViT by studying the effect of the depth (i.e., the number of layers), width (i.e., the hidden size D) and patch size on CIFAR10. As shown in Figure 3, the performance improves as the depth (i.e., the number of layers) increases from 9 to 13. Nevertheless, U-ViT does not gain from a larger depth like 17 in 50K training iterations. Similarly, increasing the width (i.e., the hidden

size) from 256 to 512 improves the performance, and further increase to 768 brings no gain; decreasing the patch size from 8 to 2 improves the performance, and further decrease to 1 brings no gain. Note that a small patch size like 2 is required for a good performance. We hypothesize it is because that the noise prediction task in diffusion models is low-level and requires small patches, differing from high-level tasks (e.g., classification). Since using a small patch size is costly for high resolution images, we firstly convert them to low-dimensional latent representations [54] and model these latent representations using U-ViT.

4. Related Work

Transformers in diffusion models. A related work is GenViT [76]. GenViT employs a smaller ViT that does not employ long skip connections and the 3×3 convolutional block, and incorporates time before normalization layers for image diffusion models. Empirically, our U-ViT performs much better than GenViT (see Table 1) by carefully designing implementation details. Another related work is VQ-Diffusion [20] and its variants [61, 69]. VQ-Diffusion firstly obtains a sequence of discrete image tokens via a VQ-GAN [17], and then models these tokens using a discrete diffusion model [1, 62] with a transformer as its backbone. Time and condition are fed into the transformer through cross attention or adaptive layer normalization. In contrast, our U-ViT simply treats all inputs as tokens, and employs long skip connections between shallow and deep layers, which achieves a better FID (see Table 1 and Table 4). In addition to images, transformers in diffusion models are also employed to encode texts [48, 52, 54, 56], decode texts [7, 29, 39, 46] and generate CLIP embeddings [52].

U-Net in diffusion models. [65, 66] initially introduce CNN-based U-Net to model the gradient of log-likelihood function for continuous image data. After that, improvements on the CNN-based U-Net for (continuous) image diffusion models are made, including using group normalization [25], multi-head attention [49], improved residual block [13] and cross attention [54]. In contrast, our U-ViT is a ViT-based backbone with conceptually simple design,



Figure 4. Image generation results of U-ViT: selected samples on ImageNet 512×512 and ImageNet 256×256, and random samples on CIFAR10, CelebA 64×64, and ImageNet 64×64.

and meanwhile has a comparable performance if not superior to a CNN-based U-Net of a similar size (see Table 1 and Table 4).

Improvements of diffusion models. In addition to the backbone, there are also improvements on other aspects, such as fast sampling [3, 44, 58, 63, 74], improved training methodology [2, 15, 30, 32, 33, 43, 49, 64, 71] and controllable generation [4, 10, 13, 22, 27, 45, 60, 80].

5. Experiments

We evaluate the proposed U-ViT in unconditional and class-conditional image generation (Section 5.2), as well as text-to-image generation (Section 5.3). Before presenting these results, we list main experimental setup below, and more details such as the sampling hyperparameters are provided in Appendix A.

5.1. Experimental Setup

Datasets. For unconditional learning, we consider CIFAR10 [36], which contain 50K training images, and

CelebA 64×64 [41], which contain 162,770 training images of human faces. For class-conditional learning, we consider ImageNet [12] at 64×64, 256×256 and 512×512 resolutions, which contains 1,281,167 training images from 1K different classes. For text-to-image learning, we consider MS-COCO [40] at 256×256 resolution, which contains 82,783 training images and 40,504 validation images. Each image is annotated with 5 captions.

High resolution image generation. We follow latent diffusion models (LDM) [54] for images at 256×256 and 512×512 resolutions. We firstly convert them to latent representations at 32×32 and 64×64 resolutions respectively, using a pretrained image autoencoder provided by Stable Diffusion² [54]. Then we model these latent representations using the proposed U-ViT.

Text-to-image learning. On MS-COCO, we convert discrete texts to a sequence of embeddings using a CLIP text encoder following Stable Diffusion. Then these embeddings are fed into U-ViT as a sequence of tokens.

²<https://github.com/CompVis/stable-diffusion>

Model on CIFAR10			FID ↓
GAN			
StyleGAN2-ADA [31]			2.92
Diff. based on U-Net		#Params	
DDPM [25]	36M		3.17
IDDPM [49]	53M		2.90
DDPM++ cont. [67]	62M		2.55
EDM [†] [30]	56M		1.97
Diff. based on ViT		#Params	
GenViT [76]	11M		20.20
U-ViT-S/2	44M		3.11
Model on CelebA 64×64			FID ↓
GAN			
DDIM [63]	79M		3.26
Soft Truncation [†] [32]	62M		1.90
Diff. model based on ViT		#Params	
U-ViT-S/4	44M		2.87
Model on ImageNet 64×64			FID ↓
GAN			
BigGAN-deep [5]			4.06
StyleGAN-XL [59]			1.51
Diff. based on U-Net		#Params	
IDDPM (small) [49]	100M		6.92
IDDPM (large) [49]	270M		2.92
CDM [26]	Unknown		1.48
ADM [13]	296M		2.07
EDM [†] [30]	296M		1.36
Diff. based on ViT		#Params	
U-ViT-M/4	131M		5.85
U-ViT-L/4	287M		4.26
Model on ImageNet 256×256			FID ↓
GAN			
BigGAN-deep [5]			6.95
StyleGAN-XL [59]			2.30
Discrete diff. based on transformer			
VQ-Diffusion [20]			11.89
VQ-Diffusion (acc0.05) [20]			5.32
Diff. based on U-Net		#Params	
IDDPM [49]	270M + 280M (SR)		12.26
CDM [26]	Unknown		4.88
ADM [13]	554M		10.94
ADM-U [13]	296M + 312M (SR)		7.49
ADM-G [13]	554M + 54M (Cls)		4.59
ADM-G, ADM-U [13]	296M + 65M (Cls) + 312M (SR)		3.94
LDM [‡] [54]	400M + 55M (AE)		3.60
Diff. based on ViT		#Params	
U-ViT-H/2 [‡]	501M + 84M (AE)		2.29
Model on ImageNet 512×512			FID ↓
GAN			
BigGAN-deep [5]			8.43
StyleGAN-XL [59]			2.41
Diff. based on U-Net		#Params	
ADM [13]	559M		23.24
ADM-U [13]	422M + 309M (SR)		9.96
ADM-G [13]	559M + 54M (Cls)		7.72
ADM-G, ADM-U [13]	422M + 43M (Cls) + 309M (SR)		3.85
Diff. based on ViT		#Params	
U-ViT-H/4 [‡]	501M + 84M (AE)		4.05

Table 1. FID results of unconditional image generation on CIFAR10 and CelebA 64×64, and class-conditional image generation on ImageNet 64×64, 256×256 and 512×512. We mark the best results *among diffusion models* in **bold**. We also include GAN results (gray) for completeness. Methods marked with [†] use advanced training techniques for diffusion models. Methods marked with [‡] model latent representations of images [54] and use classifier-free guidance [27]. We also present the number of parameters of auxiliary components for diffusion models, where SR represents a super-resolution module, AE represents an image autoencoder, and Cls represents a classifier.

U-ViT configurations. We identify several configurations of U-ViT in Table 2. In the rest of the paper, we use brief notation to represent the U-ViT configuration and the input patch size (for instance, U-ViT-H/2 means the U-ViT-Huge configuration with an input patch size of 2×2).

Model	#Layers	Hidden size D	MLP size	#Heads	#Params
U-ViT-Small	13	512	2048	8	44M
U-ViT-Small (Deep)	17	512	2048	8	58M
U-ViT-Mid	17	768	3072	12	131M
U-ViT-Large	21	1024	4096	16	287M
U-ViT-Huge	29	1152	4608	16	501M

Table 2. Configurations of U-ViT.

Training. We use the AdamW optimizer [42] with a

weight decay of 0.3 for all datasets. We use a learning rate of 2e-4 for most datasets, except ImageNet 64×64 where we use 3e-4. We train 500K iterations on CIFAR10 and CelebA 64×64 with a batch size of 128. We train 300K iterations on ImageNet 64×64 and ImageNet 256×256, and 500K iterations on ImageNet 512×512, with a batch size of 1024. We train 1M iterations on MS-COCO with a batch size of 256. On ImageNet 256×256, ImageNet 512×512 and MS-COCO, we adopt classifier-free guidance [27] following [54]. We provide more details, such as the training time and how we choose hyperparameters in Appendix A.

	4	5	10	15	20
LDM (trained 178K)	34.48	12.73	4.51	3.87	3.68
U-ViT-H/2 (trained 200K)	16.48	4.94	3.87	3.54	2.91
U-ViT-H/2 (trained 500K)	15.44	4.64	3.18	2.92	2.53

Table 3. FID results on ImageNet 256×256 under different number of sampling steps using the DPM-Solver sampler [44].

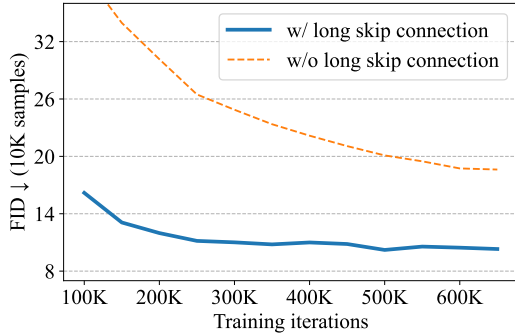


Figure 5. Ablate the long skip connection on ImageNet 256×256 (w/o classifier-free guidance).

5.2. Unconditional and Class-Conditional Image Generation

We compare U-ViT with prior diffusion models based on U-Net. We also compare with GenViT [76], a smaller ViT which does not employ long skip connections, and incorporates time before normalization layers. Consistent with previous literature, we report the FID score [23] on 50K generated samples to measure the image quality.

As shown in Table 1, U-ViT is comparable to U-Net on unconditional CIFAR10 and CelebA 64×64 , and meanwhile performs much better than GenViT.

On class-conditional ImageNet 64×64 , we initially try the U-ViT-M configuration with 131M parameters. As shown in Table 1, it gets a FID of 5.85, which is better than 6.92 of IDDPm that employs a U-Net with 100M parameters. To further improve the performance, we employ the U-ViT-L configuration with 287M parameters, and the FID improves from 5.85 to 4.26.

Meanwhile, we find that our U-ViT performs especially well in the latent space [54], where images are firstly converted to their latent representations before applying diffusion models. On class-conditional ImageNet 256×256 , our U-ViT obtains a state-of-the-art FID of 2.29, which outperforms all prior diffusion models. Table 3 further demonstrates that our U-ViT outperforms LDM under different number of sampling steps using the same sampler. Note that our U-ViT also outperforms VQ-Diffusion, which is a discrete diffusion model [1] that employs a transformer as its backbone. We also try replace our U-ViT with a U-Net with a similar amount of parameters and computational

cost, where our U-ViT still outperforms U-Net (see details in Appendix E). On class-conditional ImageNet 512×512 , our U-ViT outperforms ADM-G that directly models the pixels of images. In Figure 4, we provide selected samples on ImageNet 256×256 and ImageNet 512×512 , and random samples on other datasets, which have good quality and clear semantics. We provide more generated samples including class-conditional and random ones in Appendix F.

In Section 3.1 we have demonstrated the importance of long skip connection on small-scale dataset (i.e., CIFAR10). Figure 5 further shows it is also critical for large-scale dataset such as ImageNet.

In Appendix C, we present results of other metrics (e.g., sFID, inception score, precision and recall) as well as the computational cost (GFLOPs) with more U-ViT configurations on ImageNet. Our U-ViT is still comparable to state-of-the-art diffusion models on other metrics, and meanwhile has comparable if not smaller GFLOPs.

5.3. Text-to-Image Generation on MS-COCO

We evaluate U-ViT for text-to-image generation on the standard benchmark dataset MS-COCO. We train our U-ViT in the latent space of images [54] as detailed in Section 5.1. We also train another latent diffusion model that employs a U-Net of comparable model size to U-ViT-S, and leave other parts unchanged. Its hyperparameters and training details are provided in Appendix B. We report the FID score [23] to measure the image quality. Consistent with previous literature, we randomly draw 30K prompts from the MS-COCO validation set, and generate samples on these prompts to compute FID.

As shown in Table 4, our U-ViT-S already achieves a state-of-the-art FID among methods without accessing large external datasets during the training of generative models. By further increasing the number of layers from 13 to 17, our U-ViT-S (Deep) can even achieve a better FID of 5.48. Figure 6 shows generated samples of U-Net and U-ViT using the same random seed for a fair comparison. We find U-ViT generates more high quality samples, and meanwhile the semantics matches the text better. For example, given the text ``a baseball player swinging a bat at a ball``, U-Net generates neither the bat nor the ball. In contrast, our U-ViT-S generates the ball with even a smaller number of parameters, and our U-ViT-S (Deep) further generates the bat. We hypothesize this is because texts and images interact at every layer in our U-ViT, which is more frequent than U-Net that only interact at cross attention layer. We provide more samples in Appendix F.

6. Conclusion

This work presents U-ViT, a simple and general ViT-based architecture for image generation with diffusion models. U-ViT treats all inputs including the time, condition and

Model	FID	Type	Training datasets	#Params
Generative model trained on external large dataset (zero-shot)				
DALL-E [53]	~ 28	Autoregressive	DALL-E dataset (250M)	12B
CogView [14]	27.1	Autoregressive	Internal dataset (30M)	4B
LAFITE [82]	26.94	GAN	CC3M (3M)	75M + 151M (TE)
GLIDE [48]	12.24	Diffusion	DALL-E dataset (250M)	3.5B + 1.5B (SR)
Make-A-Scene [19]	11.84	Autoregressive	Union datasets (without MS-COCO) (35M)	4B
DALL-E 2 [52]	10.39	Diffusion	DALL-E dataset (250M)	4.5B + 700M (SR)
Imagen [56]	7.27	Diffusion	Internal dataset (460M) + LAION (400M)	2B + 4.6B (TE) + 600M (SR)
Parti [77]	7.23	Autoregressive	LAION (400M) + FIT (400M) + JFT (4B)	20B + 630M (AE)
Re-Imagen [8]	6.88	Diffusion	KNN-ImageText (50M)	2.5B + 750M (SR)
Generative model trained on external large dataset with access to MS-COCO				
VQ-Diffusion [†] [20]	13.86	Discrete diffusion	Conceptual Caption Subset (7M)	370M
Make-A-Scene [19]	7.55	Autoregressive	Union datasets (with MS-COCO) (35M)	4B
Re-Imagen [‡] [8]	5.25	Diffusion	KNN-ImageText (50M)	2.5B + 750M (SR)
Parti [†] [77]	3.22	Autoregressive	LAION (400M) + FIT (400M) + JFT (4B)	20B + 630M (AE)
Generative model trained on MS-COCO				
AttnGAN [75]	35.49	GAN	MS-COCO (83K)	230M
DM-GAN [83]	32.64	GAN	MS-COCO (83K)	46M
VQ-Diffusion [20]	19.75	Discrete diffusion	MS-COCO (83K)	370M
DF-GAN [70]	19.32	GAN	MS-COCO (83K)	19M
XMC-GAN [79]	9.33	GAN	MS-COCO (83K)	166M
Fririo [18]	8.97	Diffusion	MS-COCO (83K)	512M + 186M (TE) + 68M (AE)
LAFITE [82]	8.12	GAN	MS-COCO (83K)	75M + 151M (TE)
U-Net*	7.32	Latent diffusion	MS-COCO (83K)	53M + 123M (TE) + 84M (AE)
U-ViT-S/2	5.95	Latent diffusion	MS-COCO (83K)	45M + 123M (TE) + 84M (AE)
U-ViT-S/2 (Deep)	5.48	Latent diffusion	MS-COCO (83K)	58M + 123M (TE) + 84M (AE)

Table 4. FID results of different models on MS-COCO validation (256×256). U-ViT-S (Deep) increases the number of layers from 13 to 17 compared to U-ViT-S. We also present the number of parameters of auxiliary components for a model when it is reported in the corresponding paper, where SR represents a super-resolution module, AE represents an image autoencoder, and TE represents a text encoder. Methods marked with [†] finetune on MS-COCO. Methods marked with [‡] use MS-COCO as a knowledge base for retrieval. The U-Net* is trained by ourselves to serve as a direct baseline of U-ViT, where we leave other parts unchanged except for the backbone.

noisy image patches as tokens and employs long skip connections between shallow and deep layers. We evaluate U-ViT in tasks including unconditional and class-conditional image generation, as well as text-to-image generation. Experiments demonstrate U-ViT is comparable if not superior to a CNN-based U-Net of a similar size. These results suggest that, for diffusion-based image modeling, the long skip connection is crucial while the down/up-sampling operators in CNN-based U-Net are not always necessary. We believe that U-ViT can provide insights for future research on backbones in diffusion models and benefit generative modeling on large scale cross-modality datasets.

Acknowledgments

This work was supported by NSF of China Projects (Nos. 62061136001, 61620106010, 62076145, U19B2034, U1811461, U19A2081, 6197222); Beijing Outstanding Young Scientist Program NO. BJJWZYJH012019100020098; a grant from Tsinghua Institute for Guo Qiang; the High Performance Computing

Center, Tsinghua University; the Fundamental Research Funds for the Central Universities, and the Research Funds of Renmin University of China (22XNKJ13). C. Li was also sponsored by Beijing Nova Program. J.Z was also supported by the XPlorer Prize.

References

- [1] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. In *Advances in Neural Information Processing Systems*, 2021. 4, 7
- [2] Fan Bao, Chongxuan Li, Jiacheng Sun, Jun Zhu, and Bo Zhang. Estimating the optimal covariance with imperfect mean in diffusion probabilistic models. *ArXiv preprint*, 2022. 1, 5
- [3] Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-dpm: an analytic estimate of the optimal reverse variance in diffusion probabilistic models. *ArXiv preprint*, 2022. 1, 2, 5
- [4] Fan Bao, Min Zhao, Zhongkai Hao, Peiyao Li, Chongxuan Li, and Jun Zhu. Equivariant energy-guided sde for inverse molecular design. *ArXiv preprint*, 2022. 5



Figure 6. Text-to-image generation on MS-COCO. All the other settings except the backbone are the same. U-Net and U-ViT generate samples using the same random seed for a fair comparison. The random seed is unselected.

- [5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *7th International Conference on Learning Representations*, 2019. 6
- [6] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation. In *9th International Conference on Learning Representations*, 2021. 1
- [7] Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. Analog bits: Generating discrete data using diffusion models with self-conditioning. *ArXiv preprint*, 2022. 4
- [8] Wenhui Chen, Hexiang Hu, Chitwan Saharia, and William W Cohen. Re-imagen: Retrieval-augmented text-to-image generator. *ArXiv preprint*, 2022. 8
- [9] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. 2
- [10] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. *ArXiv preprint*, 2021. 1, 5
- [11] Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. Algorithms for learning kernels based on centered alignment. *The Journal of Machine Learning Research*, 13(1):795–828, 2012. 15
- [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2009. 5
- [13] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. *ArXiv preprint*, 2021. 1, 2, 3, 4, 5, 6, 14, 15
- [14] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. In *Advances in Neural Information Processing Systems*, 2021. 8
- [15] Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Score-based generative modeling with critically-damped langevin diffusion. *ArXiv preprint*, 2021. 1, 5
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *9th International Conference on Learning Representations*, 2021. 2, 4
- [17] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021. 4
- [18] Wan-Cyuan Fan, Yen-Chun Chen, DongDong Chen, Yu Cheng, Lu Yuan, and Yu-Chiang Frank Wang. Frido: Feature pyramid diffusion for complex scene image synthesis. *ArXiv preprint*, 2022. 8
- [19] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. *ArXiv preprint*, 2022. 8
- [20] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 3, 4, 6, 8

- [21] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 2
- [22] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *ArXiv preprint*, 2022. 5
- [23] 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. In *Advances in Neural Information Processing Systems*, 2017. 3, 7, 14
- [24] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *ArXiv preprint*, 2022. 1
- [25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020. 1, 2, 4, 6
- [26] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*, 2022. 1, 6
- [27] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *ArXiv preprint*, 2022. 5, 6
- [28] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *ArXiv preprint*, 2022. 1
- [29] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. In *Advances in Neural Information Processing Systems*, 2021. 4
- [30] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *ArXiv preprint*, 2022. 5, 6
- [31] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. In *Advances in Neural Information Processing Systems*, 2020. 6
- [32] Dongjun Kim, Seungjae Shin, Kyungwoo Song, Wanmo Kang, and Il-Chul Moon. Soft truncation: A universal training technique of score-based diffusion model for high precision score estimation. In *International Conference on Machine Learning*, 2022. 5, 6
- [33] Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *ArXiv preprint*, 2021. 1, 5
- [34] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. In *9th International Conference on Learning Representations*, 2021. 1
- [35] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International Conference on Machine Learning*, pages 3519–3529. PMLR, 2019. 15
- [36] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 3, 5
- [37] 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. 14
- [38] Kwonjoon Lee, Huiwen Chang, Lu Jiang, Han Zhang, Zhuowen Tu, and Ce Liu. Vitgan: Training gans with vision transformers. *ArXiv preprint*, 2021. 2
- [39] Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *ArXiv preprint*, 2022. 4
- [40] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, 2014. 5
- [41] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *IEEE International Conference on Computer Vision*, 2015. 5
- [42] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations*, 2019. 6
- [43] Cheng Lu, Kaiwen Zheng, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Maximum likelihood training for score-based diffusion odes by high order denoising score matching. In *International Conference on Machine Learning*, 2022. 1, 5
- [44] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *ArXiv preprint*, 2022. 1, 5, 7, 13
- [45] Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. *ArXiv preprint*, 2021. 1, 5
- [46] Eliya Nachmani and Shaked Dovrat. Zero-shot translation using diffusion models. *ArXiv preprint*, 2021. 4
- [47] Charlie Nash, Jacob Menick, Sander Dieleman, and Peter W Battaglia. Generating images with sparse representations. *arXiv preprint arXiv:2103.03841*, 2021. 14
- [48] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *ArXiv preprint*, 2021. 4, 8
- [49] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *Proceedings of the 38th International Conference on Machine Learning*, 2021. 1, 4, 5, 6
- [50] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *ArXiv preprint*, 2022. 1
- [51] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *Advances in Neural Information Processing Systems*, 34:12116–12128, 2021. 15
- [52] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *ArXiv preprint*, 2022. 1, 2, 4, 8

- [53] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *Proceedings of the 38th International Conference on Machine Learning*, 2021. 8
- [54] 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*, 2022. 1, 2, 4, 5, 6, 7, 13, 14
- [55] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, 2015. 2
- [56] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *ArXiv preprint*, 2022. 1, 2, 4, 8
- [57] 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. 14
- [58] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *ArXiv preprint*, 2022. 5
- [59] Axel Sauer, Katja Schwarz, and Andreas Geiger. Stylegan-xl: Scaling stylegan to large diverse datasets. In *ACM SIGGRAPH 2022 Conference Proceedings*, 2022. 6
- [60] Vikash Sehwal, Caner Hazirbas, Albert Gordo, Firat Ozgenel, and Cristian Canton. Generating high fidelity data from low-density regions using diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 5
- [61] Shelly Sheynin, Oron Ashual, Adam Polyak, Uriel Singer, Oran Gafni, Eliya Nachmani, and Yaniv Taigman. Knn-diffusion: Image generation via large-scale retrieval. *ArXiv preprint*, 2022. 4
- [62] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Conference on Machine Learning*, 2015. 1, 2, 4
- [63] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *9th International Conference on Learning Representations*, 2021. 1, 5, 6
- [64] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *arXiv e-prints*, 2021. 1, 5
- [65] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, 2019. 2, 4
- [66] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. In *Advances in Neural Information Processing Systems*, 2020. 4
- [67] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *9th International Conference on Learning Representations*, 2021. 1, 2, 6, 13
- [68] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. 2
- [69] Zhicong Tang, Shuyang Gu, Jianmin Bao, Dong Chen, and Fang Wen. Improved vector quantized diffusion models. *ArXiv preprint*, 2022. 4
- [70] Ming Tao, Hao Tang, Fei Wu, Xiao-Yuan Jing, Bing-Kun Bao, and Changsheng Xu. Df-gan: A simple and effective baseline for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 8
- [71] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *ArXiv preprint*, 2021. 1, 5
- [72] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017. 2, 4
- [73] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. Early convolutions help transformers see better. In *Advances in Neural Information Processing Systems*, 2021. 4
- [74] Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with denoising diffusion gans. *ArXiv preprint*, 2021. 5
- [75] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 8
- [76] Xiulong Yang, Sheng-Min Shih, Yinlin Fu, Xiaoting Zhao, and Shihao Ji. Your vit is secretly a hybrid discriminative-generative diffusion model. *ArXiv preprint*, 2022. 4, 6, 7
- [77] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *ArXiv preprint*, 2022. 8
- [78] Bowen Zhang, Shuyang Gu, Bo Zhang, Jianmin Bao, Dong Chen, Fang Wen, Yong Wang, and Baining Guo. Styleswin: Transformer-based gan for high-resolution image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 3
- [79] Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, and Yinfei Yang. Cross-modal contrastive learning for text-to-image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021. 8
- [80] Min Zhao, Fan Bao, Chongxuan Li, and Jun Zhu. Egsde: Unpaired image-to-image translation via energy-guided stochastic differential equations. *ArXiv preprint*, 2022. 1, 5
- [81] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao

Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021. 2

- [82] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. Lafite: Towards language-free training for text-to-image generation. *ArXiv preprint*, 2021. 8
- [83] Minfeng Zhu, Pingbo Pan, Wei Chen, and Yi Yang. DM-GAN: dynamic memory generative adversarial networks for text-to-image synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 8

A. Experimental Setup

We list the experimental setup for U-ViT presented in the main paper in Table 5.

Dataset	CIFAR10	CelebA 64×64	ImageNet 64×64	ImageNet 256×256	ImageNet 512×512	MS-COCO
Latent space	×	×	×	✓	✓	✓
Latent shape	-	-	-	32×32×4	64×64×4	32×32×4
Image decoder	-	-	-	ft-EMA	ft-EMA	original
Batch size	128	128	1024	1024	1024	256
Training iterations	500K	500K	300K	500K	500K	1M
Warm-up steps	2.5K	5K	5K	5K	5K	5K
Optimizer	AdamW	AdamW	AdamW	AdamW	AdamW	AdamW
Learning rate	2e-4	2e-4	3e-4	2e-4	2e-4	2e-4
Weight decay	0.03	0.03	0.03	0.03	0.03	0.03
Betas	(0.99, 0.999)	(0.99, 0.99)	(0.99, 0.99)	(0.99, 0.99)	(0.99, 0.99)	(0.9, 0.9)
Noise schedule	VP	VP	VP	SD	SD	SD
Sampler	EM	EM	DPM-Solver	DPM-Solver	DPM-Solver	DPM-Solver
Sampling steps	1K	1K	50	50	50	50
CFG	×	×	×	✓	✓	✓
p_{uncond}	-	-	-	0.1	0.1	0.1
Guidance strength	-	-	-	0.4	0.7	1
Convolution	✓	✓	✓	×	×	✓

Table 5. The experimental setup for U-ViT in the main paper. “ft-EMA” and “original” correspond to different weights of the image decoder provided in <https://huggingface.co/stabilityai/sd-vae-ft-ema>. “VP” represents the continuous-time variance preserving noise schedule [67]. “SD” represents the discrete-time noise schedule used in Stable Diffusion [54]. “EM” represents the Euler-Maruyama SDE sampler [67]. “DPM-Solver” represents the DPM-Solver ODE sampler [44]. “ p_{uncond} ” represents the unconditional training probability in classifier free guidance (CFG). “Convolution” represents whether to add a 3×3 convolutional block before output.

In our early experiments, we try learning rates between $1e-4$ and $5e-4$, and find that a learning rate of $2e-4$ performs well for all datasets. On ImageNet 64×64 , a learning rate of $3e-4$ could further improve the performance. We try weight decay between 0.01 and 0.05, and find that a weight decay of 0.03 performs well for all datasets. We try the running coefficients β_1, β_2 of AdamW among $\{0.9, 0.99, 0.999\}$, and find that $(\beta_1, \beta_2) = (0.99, 0.99)$ performs well for all datasets. On CIFAR10, $\beta_2 = 0.999$ could further improve the performance. On MS-COCO, $(\beta_1, \beta_2) = (0.9, 0.9)$ could further improve the performance. We train with mixed precision for efficiency, and the training time and devices are listed in Table 6. Besides, the training memory of U-ViT can be greatly reduced with the gradient checkpointing trick. For example, the memory for forward and backward on a single A100 can be reduced from 53GB to 10GB when training U-ViT-L/2 with a batch size of 128 on ImageNet 256×256 .

During inference, with 1 A100, generating 500 samples with DPM-Solver takes around 19 seconds, 34 seconds, 59 seconds, 89 seconds, with U-ViT-S, U-ViT-M, U-ViT-L, U-ViT-H respectively. The time would double if classifier-free guidance is used.

Dataset	Model	Training devices	Training time	Training iterations
CIFAR10	U-ViT-S/2	4 GeForce RTX 2080 Ti	24 hours	500K
CelebA	U-ViT-S/4	4 GeForce RTX 2080 Ti	24 hours	500K
ImageNet 64×64	U-ViT-M/4	8 A100	59 hours	300K
ImageNet 64×64	U-ViT-L/4	8 A100	100 hours	300K
ImageNet 256×256	U-ViT-L/2	8 A100	100 hours	300K
ImageNet 256×256	U-ViT-H/2	8 A100	208 hours	500K
ImageNet 512×512	U-ViT-L/4	8 A100	166 hours	500K
ImageNet 512×512	U-ViT-H/4	8 A100	208 hours	500K
MS-COCO	U-ViT-S/2	4 A100	60 hours	1M
MS-COCO	U-ViT-S/2 (deep)	4 A100	74 hours	1M

Table 6. The training devices and time.

B. Details of the U-Net Baseline on MS-COCO

We employ the U-Net with cross attention provided by LDM [54] for the baseline. The U-Net is performed on the 32×32 resolution latent representation, and down-samples it to 16×16 , 8×8 and 4×4 resolution. The number of channels is 128 at 32×32 resolution, and 256 at other resolutions. Each resolution has 2 residual blocks. The U-Net performs self attention and cross attention at 16×16 and 8×8 resolution. Such a configuration leads to a total of 53M parameters, which is comparable to 45M of U-ViT-Small for a fair comparison. We use the AdamW optimizer with weight decay set to 0.01 and running coefficients β_1, β_2 set to (0.9, 0.999), which are the setting used across LDM [54]. We tune the learning rate of U-Net and find $2e-4$ performs the best. The training iterations and the batch size of U-Net are the same to U-ViT for a fair comparison.

C. Results of Other Metrics and Configurations on ImageNet

We present results of FID [23], sFID [47], inception score (IS) [57], precision and recall [37] on ImageNet in Table 7. Our U-ViT is still comparable to state-of-the-art diffusion models based on U-Net on these metrics, and meanwhile has comparable if not smaller GFLOPs.

ImageNet 64×64	#Params	GFLOPs	FID↓	sFID↓	IS↑	Precision↑	Recall↑
ADM [13]	296M	110	2.07	4.29	-	0.74	0.63
U-ViT-M/4 (VP, trained 300K, w/ conv)	131M	35	5.85	4.09	33.71	0.69	0.61
U-ViT-L/4 (VP, trained 300K, w/ conv)	287M	77	4.26	3.77	40.66	0.71	0.62
ImageNet 256×256	#Params	GFLOPs	FID↓	sFID↓	IS↑	Precision↑	Recall↑
ADM-G, ADM-U [13]	296M + 65M (Cls) + 312M (SR)	110 + 19 (Cls) + 632 (SR)	3.94	6.14	215.84	0.83	0.53
LDM [54]	400M + 55M (AE)	104 + 336 (AE)	3.60	-	247.67	0.87	0.48
U-ViT-L/2 (VP, trained 300K, w/ conv, original, $p_{uncond}=0.15$)	287M + 84M (AE)	77 + 312 (AE)	3.40	6.63	219.94	0.83	0.52
U-ViT-H/2 (VP, trained 300K, w/ conv, original, $p_{uncond}=0.1$)	501M + 84M (AE)	133 + 312 (AE)	3.10	6.70	250.82	0.84	0.53
U-ViT-H/2 (VP, trained 300K, w/o conv, original, $p_{uncond}=0.1$)	501M + 84M (AE)	133 + 312 (AE)	3.74	8.04	244.47	0.84	0.51
U-ViT-H/2 (SD, trained 300K, w/ conv, original, $p_{uncond}=0.15$)	501M + 84M (AE)	133 + 312 (AE)	3.14	7.81	229.03	0.82	0.55
U-ViT-H/2 (SD, trained 300K, w/o conv, original, $p_{uncond}=0.15$)	501M + 84M (AE)	133 + 312 (AE)	2.90	7.70	242.59	0.81	0.56
U-ViT-H/2 (SD, trained 300K, w/o conv, original, $p_{uncond}=0.1$)	501M + 84M (AE)	133 + 312 (AE)	2.78	7.55	251.83	0.82	0.56
U-ViT-H/2 (SD, trained 500K, w/o conv, original, $p_{uncond}=0.1$)	501M + 84M (AE)	133 + 312 (AE)	2.65	8.17	260.34	0.81	0.57
U-ViT-H/2 (SD, trained 500K, w/o conv, ft-EMA, $p_{uncond}=0.1$)	501M + 84M (AE)	133 + 312 (AE)	2.29	5.68	263.88	0.82	0.57
ImageNet 512×512	#Params	GFLOPs	FID↓	sFID↓	IS↑	Precision↑	Recall↑
ADM-G, ADM-U [13]	422M + 43M (Cls) + 309M (SR)	307 + 21 (Cls) + 2506 (SR)	3.85	5.86	221.72	0.84	0.53
U-ViT-L/4 (VP, trained 500K, w/ conv, original, $p_{uncond}=0.15$)	287M + 84M (AE)	77 + 1260 (AE)	4.67	5.87	213.28	0.87	0.45
U-ViT-H/4 (SD, trained 500K, w/o conv, original, $p_{uncond}=0.1$)	501M + 84M (AE)	133 + 1260 (AE)	4.34	8.44	261.13	0.84	0.48
U-ViT-H/4 (SD, trained 500K, w/o conv, ft-EMA, $p_{uncond}=0.1$)	501M + 84M (AE)	133 + 1260 (AE)	4.05	6.44	263.79	0.84	0.48

Table 7. Results of FID [23], sFID [47], inception score (IS) [57], precision and recall [37] on ImageNet. We also show the number of parameters as well as the GFLOPs.

D. CKA Analysis

Centered kernel alignment (CKA) is widely used to analyze similarity between hidden representations in deep neural networks [11, 35, 51]. In this section, we use the CKA method to analyze hidden representations of networks that employ three ways to combine long skip branches: (1) concatenation, i.e., $\text{Linear}(\text{Concat}(\mathbf{h}_m, \mathbf{h}_s))$; (2) addition, i.e., $\mathbf{h}_m + \mathbf{h}_s$; (3) no long skip connection. These three ways are elaborated in Section 3.1 in the main paper. We evaluate hidden representations after each transformer block and fix the input time as $t = 0.5$ on CIFAR10.

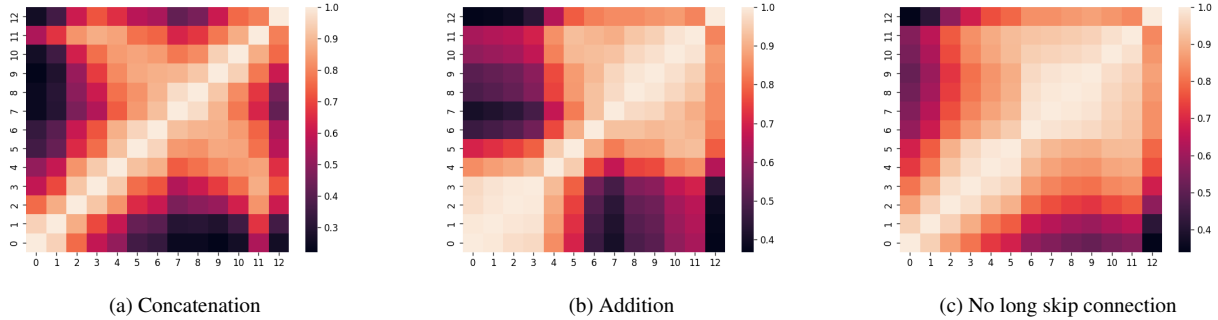


Figure 7. CKA analysis on hidden representations of networks that employ three ways to combine long skip branches. We analyze the similarity between hidden representations after each transformer block in the same network.

We find that the “addition” and “no long skip connection” settings share a similar phenomenon that neighboring blocks in the network have similar representations, e.g., blocks 0-3, 6-11 in Figure 7 (b), and blocks 0-5, 6-11 in Figure 7 (c). In contrast, the representations of neighboring blocks under the “concatenation” setting have low similarity, as shown in Figure 7 (a). Thus, the “concatenation” setting significantly changes the representations in the transformer, while the “addition” setting does not.

E. Compare with U-Net Under Similar Amount of Parameters and Computational Cost

On ImageNet 256×256 , we also try replace our U-ViT with a U-Net with a similar amount of parameters and computational cost. The U-Net employs implementation from ADM [13]. We set the model channels as 320, the channel multiplier as (2, 2, 4), the number of residual blocks as 3, and employs attention at $2 \times$ and $4 \times$ down-sampling. This leads to a U-Net of 646M parameters and 135 GFLOPs, and our U-ViT has 501M parameters and 133 GFLOPs. We use the same optimizer configuration as ADM. As shown in Figure 8, our U-ViT consistently outperforms U-Net at different training iterations without classifier-free guidance. We also evaluate FID with 50K samples at 500K training iterations. With no classifier-free guidance, U-ViT obtains a FID of 6.58 and U-Net obtains a FID of 10.69. With a classifier-free guidance scale of 0.4, U-ViT obtains a FID of 2.29 and U-Net obtains a FID of 2.66. Under both settings, our U-ViT outperforms U-Net.

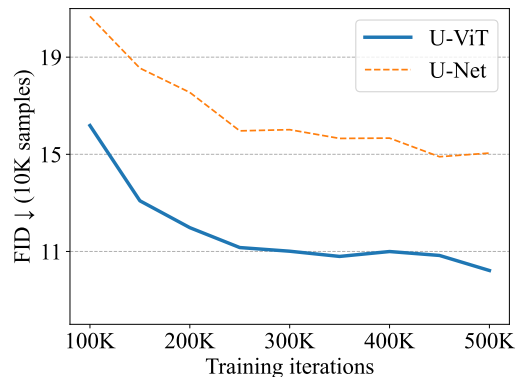


Figure 8. Compare with U-Net under similar amount of parameters and computational cost (w/o classifier-free guidance).

F. Additional Samples



Figure 9. Generated samples on ImageNet 512×512 , conditioned on goldfish (1), arctic fox (279), monarch butterfly (323), african elephant (386), flamingo (130), tennis ball (852).



Figure 10. Generated samples on ImageNet 512×512 , conditioned on cheeseburger (933), fountain (562), balloon (417), tabby cat (281), lorikeet (90), agaric (992).

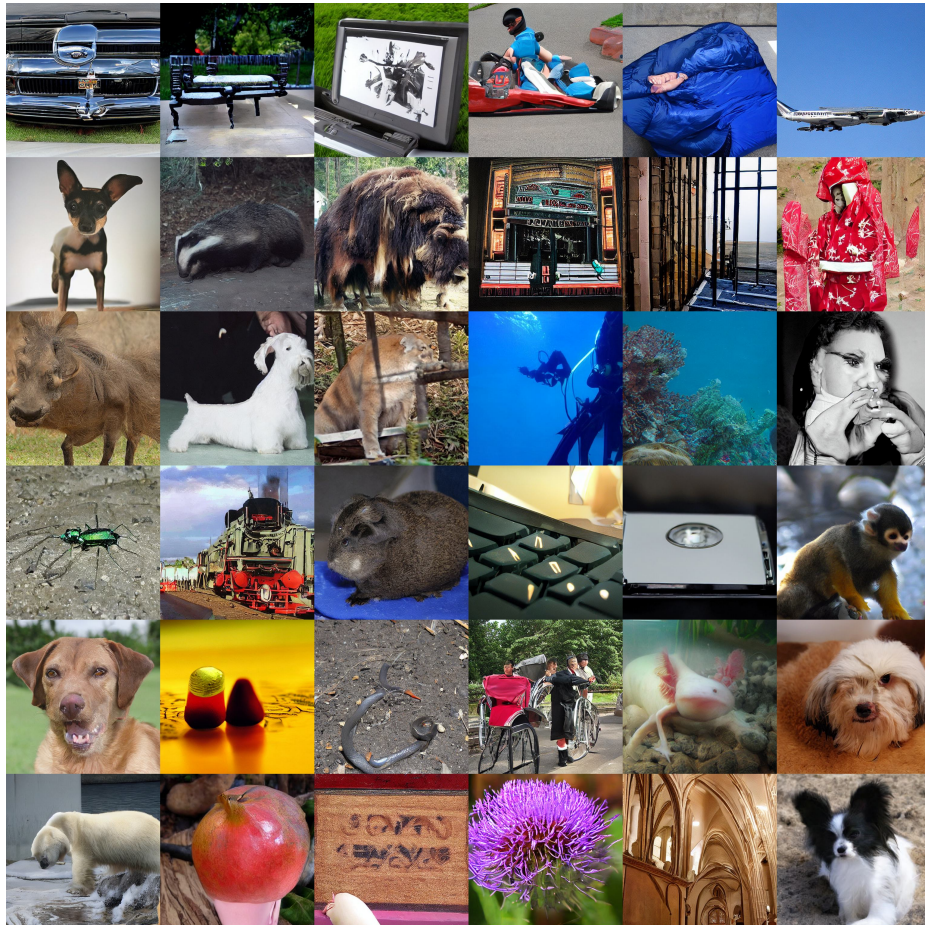


Figure 11. Random samples on ImageNet 512×512.



Figure 12. Generated samples on ImageNet 256×256 , conditioned on goldfish (1), arctic fox (279), monarch butterfly (323), african elephant (386), flamingo (130), tennis ball (852), cheeseburger (933), fountain (562), balloon (417), tabby cat (281), lorikeet (90), agaric (992).

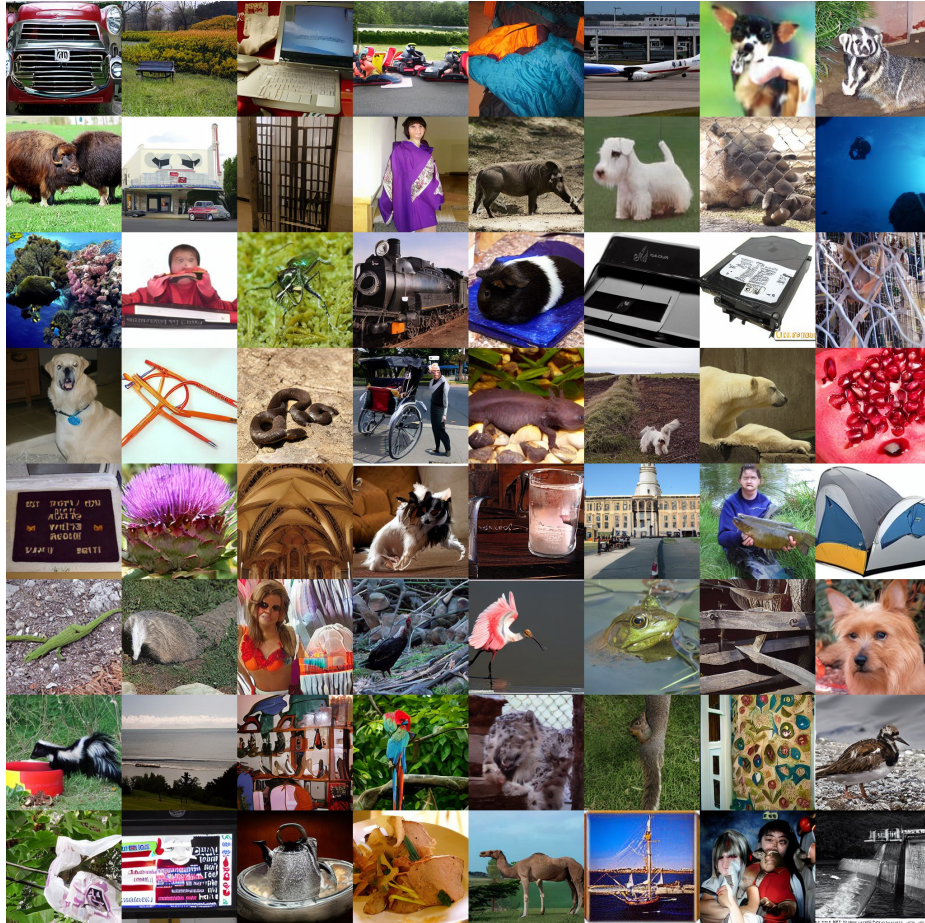


Figure 13. Random samples on ImageNet 256×256 .

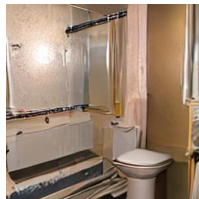
Group of whimsical, colorful artificial flowers in bottles.



A man wearing black glasses and a mustache.



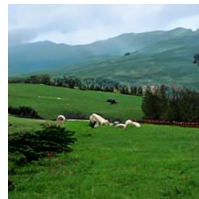
A bathroom with a toilet, sink bowl and mirror.



A black and white photo of two teddy bears posing near two cameras.



A group of sheep in a grassy area with trees in the back ground.



A bench sitting along side of river next to tree.



A group photo of a tennis team on the court.



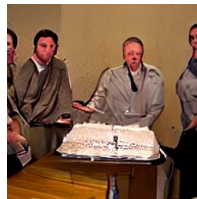
a close up of a head of broccoli in a garden



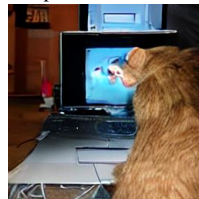
A clear bowl of broccoli and chopped nuts.



A group of people standing around a white cake on a table.



A cat watches a blonde haired man on a laptop computer screen.



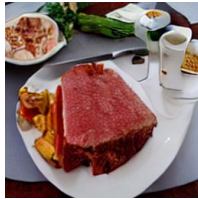
A display case displays various types of deserts.



A young boy kicking a soccer ball across a field.



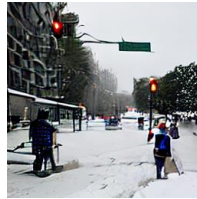
Close up of a plate with food on it.



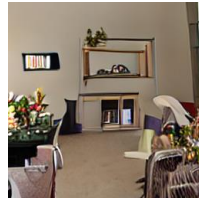
Group of whimsical, colorful artificial flowers in bottles.



People are at a stop light on a snowy street.



The furniture in the living room is decorated with flowers.



Kites fly high in the air over a park.



Figure 14. Random samples on MS-COCO. Prompts are randomly drawn from the validation set.