
All are Worth Words: a ViT Backbone for Score-based Diffusion Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Vision transformers (ViT) have shown promise in various vision tasks including low-
2 level ones while the U-Net remains dominant in score-based diffusion models. In
3 this paper, we perform a systematical empirical study on the ViT-based architectures
4 in diffusion models. Our results suggest that adding extra long skip connections
5 (like the U-Net) to ViT is crucial to diffusion models. The new ViT architecture,
6 together with other improvements, is referred to as U-ViT. On several popular
7 visual datasets, U-ViT achieves competitive generation results to SOTA U-Net
8 while requiring comparable amount of parameters and computation if not less.

9 1 Introduction

10 Along with the development of algorithms, the revolution of backbones plays a central role in
11 the success of (score-based) diffusion models. A representative example is the U-Net architecture
12 employed in prior work [15, 5], which remains dominant in diffusion models for image generation
13 tasks. A very natural question is whether the reliance of the U-Net is necessary in such models.

14 On the other hand, vision transformers (ViT) [3] have shown promise in various vision tasks [1, 4]
15 including low-level ones [17, 19]. Compared to CNN, ViT is preferable at a large scale because of
16 its scalability and efficiency [3]. Although the score-based diffusion models have been scaled up
17 dramatically [12], it is still not clear whether ViT is suitable for score modeling or not.

18 In this paper, we perform a systematical empirical study on the ViT-based architectures in diffusion
19 models. We modify the standard ViT as follows:

- 20 1. adding extra long skip connections (like the U-Net),
- 21 2. adding an extra 3x3 convolutional block before output, and
- 22 3. treating everything including the time embedding, label embedding and patches of the noisy
23 image as tokens.

24 The resulting architecture is referred to as *U-ViT*.

25 On several popular visual datasets, U-ViT achieves competitive generation results to SOTA U-Net
26 architectures while requires comparable amount of parameters and computation if not less. Our
27 results suggest that

- 28 1. ViT is promising for score-based diffusion models;
- 29 2. the long skip connections play a central role in the success of diffusion models; and
- 30 3. the down-sampling and up-sampling operators are not necessary for diffusion models.

31 We believe that future diffusion models on large scale or cross-modality datasets potentially benefit
32 from U-ViT.

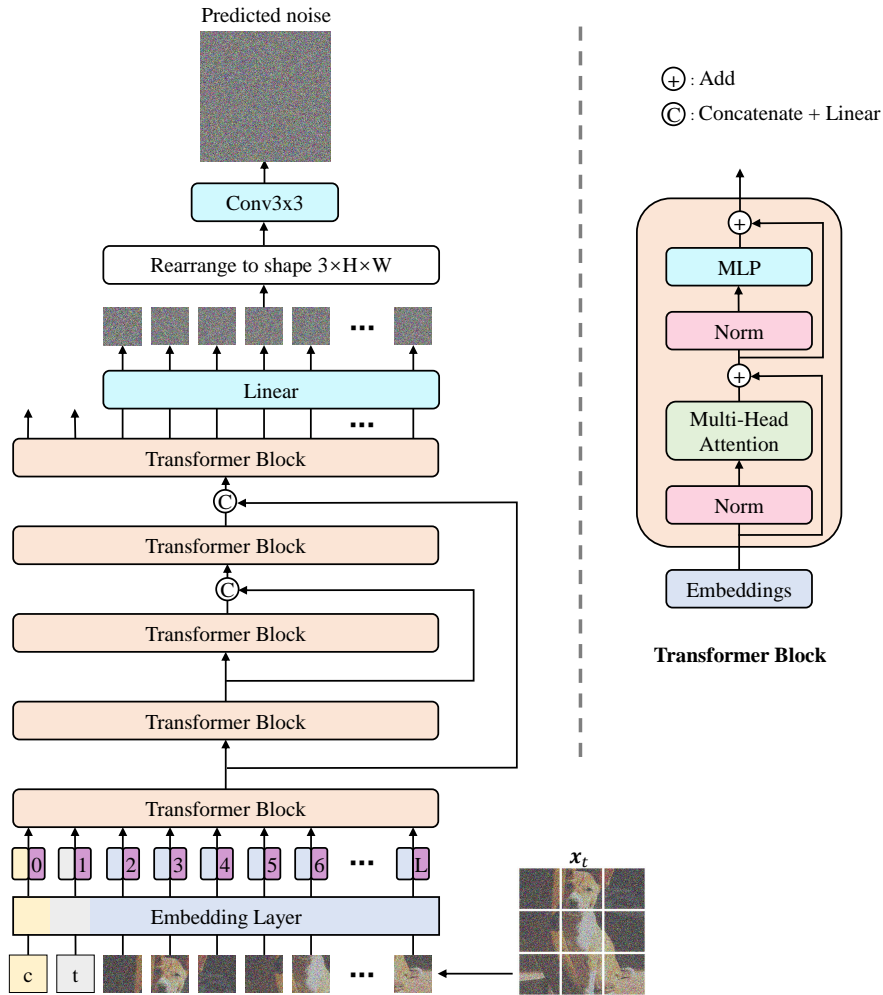


Figure 1: The U-ViT architecture.

33 2 Development of the U-ViT Architecture

34 We first attempt to train a diffusion model using a vanilla ViT [3] on CIFAR10. For simplicity, we
 35 treat everything including the time embedding, label embedding and patches of the noisy image as
 36 tokens. With carefully tuned hyperparameters, a 13-layer ViT of size 41M achieves a FID 5.97, which
 37 is significantly better than 20.20 of the prior ViT-based diffusion models [18]. We conjecture that this
 38 is mainly because our model is larger. However, this is clearly worse than 3.17 of the U-Net [5] of a
 39 similar size.

40 The importance of the skip connections in U-Net has been realized for a long time in low-level vision
 41 tasks [13]. Since all local information are also crucial in score modeling (or noise prediction), we
 42 hypothesize that the skip connections play a central role in such tasks as well. Therefore, we add
 43 extra skip connections to ViT and obtain a FID of 4.24.

44 Finally, we add a 3x3 convolutional block before the output to avoid potential artifacts between
 45 patches and obtain a FID of 3.11, which is competitive to the results of DDPM [5]. The overall
 46 architecture is illustrated in Fig. 1 and the ablation results are summarized in Table 1 for clarity.

Table 1: Ablation study on the architecture design on CIFAR10.

Skip connection	Conv3x3	FID
✓	✓	3.11
✓	×	4.24
×	✓	7.37
×	×	5.97

47 3 Experiments

48 We evaluate U-ViT on CIFAR10 [7], CelebA 64x64 [8], and ImageNet 64x64 [2]. We provide
49 detailed experimental settings in Table 2.

Dataset	CIFAR10	CelebA 64x64	ImageNet 64x64
Patch size	2	4	4
Layers	13	13	17
Hidden size	512	512	768
MLP size	2048	2048	3072
Heads	8	8	12
Params	44M	44M	131M
Noise schedule	VP SDE [16]	VP SDE	VP SDE
Batch size	128	128	1024
Training steps	500K	500K	300K
Warm-up steps	5K	5K	5K
Optimizer	AdamW [9]	AdamW	AdamW
Learning rate	2e-4	2e-4	3e-4
Weight decay	0.03	0.03	0.03
Betas	(0.99, 0.999)	(0.99, 0.99)	(0.99, 0.99)
Sampler	EM	EM	DPM-Solver [10]
Sampling steps	1K	1K	50

Table 2: The experimental settings. EM represents the Euler-Maruyama sampler.

50 We compare U-ViT with commonly used U-Net in diffusion models [5, 11, 16]. We also compare with
51 GenViT [18], a smaller ViT which does not employ long skip connections and the 3x3 convolutional
52 block, and incorporates time before normalization layers. As shown in Table 3, the FID results on
53 CIFAR10 and CelebA 64x64 are comparable to U-Net. As shown in Table 4, on ImageNet 64x64,
54 U-ViT is comparable to IDDPM U-Net (small), which has a comparable number of parameters.
55 Note that there is still a gap between U-ViT and IDDPM U-Net (large), which could potentially be
56 narrowed by further increasing the U-ViT size or increasing the batch size and training steps. We
57 provide generated samples of U-ViT in Figure 2, which have good quality and clear semantics.

Table 3: FID ↓ results on unconditional datasets.

Architecture	CIFAR10	CelebA 64x64
DDPM U-Net [5]	3.17	3.26 [14]
IDDPM U-Net [11]	2.90	-
DDPM++ U-Net [16]	2.55	1.90 [6]
GenViT [18]	20.20	-
U-ViT (ours)	3.11	3.13



Figure 2: Generated samples of U-ViT.

Table 4: FID \downarrow results on class-conditional ImageNet 64x64 and comparison of experimental setting.

Architecture	FID \downarrow	Params	Batch size	Training steps
IDDPM U-Net (small) [11]	6.92	100M	2048	1700K
IDDPM U-Net (large) [11]	2.92	270M	2048	250K
U-ViT (ours)	6.75	131M	1024	300K

58 3.1 Efficiency Comparison

59 We compare efficiency of U-Net and U-ViT on CIFAR10 in Table 5. U-ViT has fewer parameters.
 60 When the computation resource is unsaturated, e.g., using a batch size of 1, U-ViT has a much higher
 61 throughput than U-Net. When the computation resource is saturated, e.g., using a large batch size of
 62 500, U-ViT has a slightly lower throughput than U-Net. This means that U-ViT has a slightly larger
 63 computational cost, but meanwhile enjoys a better parallelism than U-Net.

Table 5: Efficiency comparison on CIFAR10 in one A40 GPU. Throughput is measured by the number of processed inputs in a second.

Method	FID \downarrow	Params	Throughput (batch size=1)	Throughput (batch size=500)
IDDPM U-Net [11]	2.90	53M	22/s	1297/s
U-ViT (ours)	3.11	44M	55/s	1125/s

64 References

- 65 [1] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised
 66 vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer
 67 Vision*, pages 9640–9649, 2021.
- 68 [2] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
 69 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern
 70 recognition*, pages 248–255. Ieee, 2009.
- 71 [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
 72 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
 73 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint
 74 arXiv:2010.11929*, 2020.

- 75 [4] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
76 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on*
77 *Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.
- 78 [5] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv*
79 *preprint arXiv:2006.11239*, 2020.
- 80 [6] Dongjun Kim, Seungjae Shin, Kyungwoo Song, Wanmo Kang, and Il-Chul Moon. Soft
81 truncation: A universal training technique of score-based diffusion model for high precision
82 score estimation. In *International Conference on Machine Learning*, pages 11201–11228.
83 PMLR, 2022.
- 84 [7] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
85 2009.
- 86 [8] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the
87 wild. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile,*
88 *December 7-13, 2015*, pages 3730–3738. IEEE Computer Society, 2015.
- 89 [9] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
90 *arXiv:1711.05101*, 2017.
- 91 [10] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A
92 fast ode solver for diffusion probabilistic model sampling in around 10 steps. *arXiv preprint*
93 *arXiv:2206.00927*, 2022.
- 94 [11] Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *arXiv*
95 *preprint arXiv:2102.09672*, 2021.
- 96 [12] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
97 text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 98 [13] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for
99 biomedical image segmentation. In *International Conference on Medical image computing and*
100 *computer-assisted intervention*, pages 234–241. Springer, 2015.
- 101 [14] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
102 *preprint arXiv:2010.02502*, 2020.
- 103 [15] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data
104 distribution. *arXiv preprint arXiv:1907.05600*, 2019.
- 105 [16] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and
106 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*
107 *preprint arXiv:2011.13456*, 2020.
- 108 [17] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer
109 for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on*
110 *Computer Vision*, pages 7262–7272, 2021.
- 111 [18] Xiulong Yang, Sheng-Min Shih, Yinlin Fu, Xiaoting Zhao, and Shihao Ji. Your vit is secretly a
112 hybrid discriminative-generative diffusion model. *arXiv preprint arXiv:2208.07791*, 2022.
- 113 [19] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei
114 Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation
115 from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF*
116 *conference on computer vision and pattern recognition*, pages 6881–6890, 2021.