# Supplementary Materials for ECCV 2022 paper PixelFolder: An Efficient Progressive Pixel Synthesis Network for Image Generation

Jing He<sup>1</sup>, Yiyi Zhou<sup>1\*</sup>, Qi Zhang<sup>2</sup>, Jun Peng<sup>1</sup>, Yunhang Shen<sup>2</sup>, Xiaoshuai Sun<sup>1,3</sup>, Chao Chen<sup>2</sup>, and Rongrong Ji<sup>1,3</sup>

<sup>1</sup>Media Analytics and Computing Lab, Department of Artificial Intelligence, School of Informatics, Xiamen University. <sup>2</sup>Youtu Lab, Tencent. <sup>3</sup>Institute of Artificial Intelligence, Xiamen University.

{blinghe, pengjun}@stu.xmu.edu.cn, {zhouyiyi, xssun, rrji}@xmu.edu.cn, {merazhang, aaronccchen}@tencent.com, shenyunhang01@gmail.com

## A More Experiments

### A.1 Comparison of additional metrics for image quality

To further validate our PixelFolder in terms of image quality, we report two new indicators, namely Inception Score (IS) and Perceptual Path Length (PPL). The IS was shown to correlate well with the human scoring for judging the realism of generated images on CIFAR-10 dataset [6], the higher IS score indicators the better image quality. The PPL indicates the smoothness of the mapping from a latent space to the output images, which is very important for high quality image generation [3]. The smooth mapping will have a small PPL. In Tab. A, the overall performance of PixelFolder is still better.

Mathad	FFHQ, $256 \times 256$		LSUN Church, $256 \times 256$		
Method	IS $\uparrow$	$\mathrm{PPL}\downarrow$	$IS\uparrow$	$\mathrm{PPL}\downarrow$	
INR-GAN	$4.51 {\pm} 0.08$	62	$2.69 \pm 0.02$	529	
CIPS	$4.75 {\pm} 0.06$	273	$2.73 {\pm} 0.02$	3132	
StyleGAN2	$4.86 {\pm} 0.04$	209	$2.78 \pm 0.02$	350	
PixelFolder	$5.04{\pm}0.07$	197	$2.79{\pm}0.03$	307	

Table A: The comparison of inception score (IS) and perceptual path length (PPL) on FFHQ and LSUN Church. The higher IS score and the lower PPL score indicate the better performance in terms of image quality. These results further validate the superior of the proposed PixelFolder.

\* Corresponding Author.

Mothod	FFHQ, $512 \times 512$		FFHQ, $1024 \times 1024$		
Method	$\# Params \downarrow$	$\mathrm{FID}\downarrow$	#Params↓	$\mathrm{FID}\!\!\downarrow$	
INR-GAN	112.43	-	117.30	16.32	
CIPS	145.12	-	574.32	-	
StyleGAN2	30.28	4.15	30.37	2.84	
PixelFolder	20.84	4.08	21.14	2.98	

Table B: Comparison of high-resolution image generation. The proposed PixelFolder has better performance with fewer parameters than other methods.

Method	Cat	Bedroom	Method	FFHQ	Church
INR-GAN	-	3.40	UDM(RVE)+ST [a]	5.54	-
CIPS StyleGAN2	12.86 9.14	-2.65	UnleashTR [b]	6.11	4.07
PixelFolder	8.41	3.71	StyleGAN2 PixelFolder	3.83 <b>3.77</b>	3.86 <b>2.46</b>

Table C: Comparison of FID scores on Table D: Comparison with diffusionLSUN Cat and LSUN Bedroom.models in terms of FID score.

#### A.2 Comparison of high-resolution image generation

In Tab. B, we compare the performance of high-resolution image generation between the proposed PixelFolder and other SOTA methods, *i.e.*  $512 \times 512$  and  $1024 \times 1024$ . From these results, we observe that PixelFolder still has competitive performance in terms of model efficiency and image quality. Compared to StyleGAN2, PixelFolder obtains competitive FID values with only two-thirds of its parameters, and especially at  $512 \times 512$  resolutions where PixelFolder obtains better performance. Note that the proposed PixelFolder is directly implemented without any tricks like adding additional noises. INR-GAN and CIPS are out of GPU memory to train.

#### A.3 Comparison on more datasets

We conduct additional experiments on two popular datasets as reported in Tab. C, *i.e.* LSUN cat and LSUN bedroom. The proposed PixelFolder still obtains competitive performance, which further demonstrates the generalizability of our PixelFolder.

#### A.4 Comparison with more typical methods

In addition to some GAN-based approaches [1,4], we further compare PixelFolder with the recently popular diffusion models [5,2] in Tab. D, where the merits of our method still can be seen.

## **B** Experiment details

To ablate the pixel folding operations, we compare them with the alternatives of *down-sampling* and *deconvolution* (DeConv) in Tab.3 of the main paper. We can see that both alternative models, *i.e.*, Fold+DeConv and Down-Sampling+DeConv, have more parameters than Fold+Unfold (base). The main reason for the parameter gap is that the pixel unfolding operation reduces the number of channels of the hidden tensors, thus the corresponding convolution kernel has fewer parameters (see blue font in Tab. E).

In order to keep the shape of the hidden tensor with red color consistent, we further modify  $\text{DeConv}_0$  and  $\text{DeConv}_1$ , *i.e.*, the kernels are replaced with  $3 \times 3 \times 512 \times 128$  and  $3 \times 3 \times 128 \times 512$ , respectively. The results of the shape-consistent alternatives are shown in Tab. F. Although Fold+DeConv-sc and Down-Sampling+DeConv-sc guarantee the consistency of the shape of pixel tensors, their performance lags far behind that of Fold+Unfold (base), which further demonstrates the effectiveness of pixel folding operations.

Layers	Fold+UnFold (base)	Fold+DeConv	DownS.+DeConv
Initialization	$16\times16\times512$	$16\times16\times512$	$16\times16\times512$
projection	$16 \times 16 \times 32$	$16 \times 16 \times 32$	-
Folding×2/DownS.	$4 \times 4 \times 512$	$4 \times 4 \times 512$	$4 \times 4 \times 512$
ModConv <sub>0</sub> /DeConv <sub>0</sub>	$4 \times 4 \times 512 (3 \times 3 \times 512 \times 512)$	$\frac{8 \times 8 \times 512}{(3 \times 3 \times 512 \times 512)}$	$\frac{8 \times 8 \times 512}{(3 \times 3 \times 512 \times 512)}$
Unfolding <sub>0</sub>	$8 \times 8 \times 128$	-	-
ModConv <sub>1</sub>	$8 \times 8 \times 512$ $(3 \times 3 \times 128 \times 512)$	$8 \times 8 \times 512$ $(3 \times 3 \times 512 \times 512)$	$\frac{8 \times 8 \times 512}{(3 \times 3 \times 512 \times 512)}$
$ModConv_2/DeConv_1$	$8 \times 8 \times 512$ $(3 \times 3 \times 512 \times 512)$	$\frac{16 \times 16 \times 512}{(3 \times 3 \times 512 \times 512)}$	$\frac{16 \times 16 \times 512}{(3 \times 3 \times 512 \times 512)}$
Unfolding <sub>1</sub>	$16\times16\times128$	-	-
ModConv <sub>3</sub>	$16 \times 16 \times 512$ $(3 \times 3 \times 128 \times 512)$	$ \begin{array}{c} 16 \times 16 \times 512 \\ (3 \times 3 \times 512 \times 512) \end{array} $	$\frac{16 \times 16 \times 512}{(3 \times 3 \times 512 \times 512)}$
ToRGB	$16 \times 16 \times 3$ $(1 \times 1 \times 3 \times 512)$	$ \begin{array}{c} 16 \times 16 \times 3 \\ (1 \times 1 \times 3 \times 512) \end{array} $	$   \begin{array}{c}     16 \times 16 \times 3 \\     (1 \times 1 \times 3 \times 512)   \end{array} $

Table E: The tensor shapes of the output of layers in the first row.  $(k \times k \times c_{in} \times c_{out})$  represents the kernel size of corresponding convolution layer. The red font indicates the tensor shape before feeding into ModConv or DeConv, and the blue font indicates the kernel size of ModConv or DeConv with red tensor as input.

4 J. He et al.

Settings	#Parm (M) $\downarrow$	$\mathrm{GMACs}\downarrow$	$\mathrm{FID}\downarrow$	$\operatorname{Precision} \uparrow$	$\operatorname{Recall} \uparrow$
Fold+Unfold (base)	20.84	23.78	5.49	0.679	0.514
Fold+DeConv Fold+DeConv-sc	$29.41 \\ 12.66$	$86.53 \\ 12.12$	$5.60 \\ 8.41$	$0.667 \\ 0.643$	$0.371 \\ 0.366$
Down-Sampling+DeConv Down-Sampling+DeConv-sc	$29.21 \\ 12.46$	$89.38 \\ 12.97$	$5.53 \\ 8.36$	$0.679 \\ 0.670$	$\begin{array}{c} 0.456 \\ 0.442 \end{array}$

Table F: Additional ablation study on FFHQ. The models of all settings are trained with 200k steps for a quick comparison. "sc" is short for "shape-consistent". The results further prove the advantages of pixel folding operations.

# C Additional samples

In Fig. A, we provide more interpolations on LSUN Church to further demonstrate the generalization capability of PixelFolder. We also provide additional samples on different datasets in Fig. B. Similar to StyleGANv1v2 [3,4] and CIPS [1], our model also has the capability of style mixing controlled by stagewise latent codes as illustrated in Fig. C.



Fig. A: Interpolations by PixelFolder on LSUN Church. The factor  $\alpha$  of the interpolations is set to smaller than that of Fig.4 of the main paper. These interpolations further demonstrate the generalization capability of PixelFolder when the input noise fluctuates slightly.



(c) LSUN Bedroom.

(d) LSUN Cat.

Fig. B: Additional samples of PixelFolder on different datasets,  $\it i.e.,$  FFHQ, LSUN Church, LSUN Bedroom and LSUN Cat.



Fig. C: Stage-wise style mixing of our PixelFolder trained on FFHQ. The left two columns are the images generated by latent code  $z_1$  and  $z_2$ , respectively. The three columns on the right are the images generated by replacing  $z_1$  at stage1, stage2 and stage3 with the corresponding latent code  $z_2$ .

# References

- Anokhin, I., Demochkin, K., Khakhulin, T., Sterkin, G., Lempitsky, V., Korzhenkov, D.: Image generators with conditionally-independent pixel synthesis. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 14278– 14287 (2021)
- Bond-Taylor, S., Hessey, P., Sasaki, H., Breckon, T.P., Willcocks, C.G.: Unleashing transformers: Parallel token prediction with discrete absorbing diffusion for fast high-resolution image generation from vector-quantized codes. arXiv preprint arXiv:2111.12701 (2021)
- 3. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 4401–4410 (2019)
- Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 8110–8119 (2020)
- Kim, H., Choi, Y., Kim, J., Yoo, S., Uh, Y.: Exploiting spatial dimensions of latent in gan for real-time image editing. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 852–861 (2021)
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X.: Improved techniques for training gans. Advances in neural information processing systems 29 (2016)