

# TransGAN: Two Pure Transformers Can Make One Strong GAN, and That Can Scale Up

(NeurIPS 2021)

# Introduction

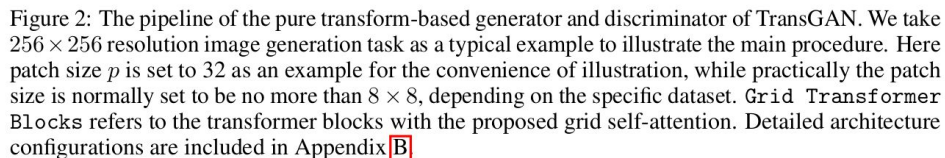
## **TransGAN**

- Generative adversarial networks(GANs) architecture which is transformer based and use convolutions at all.
- Generator is memory friendly and discriminator is multi scale.

## **Contributions :**

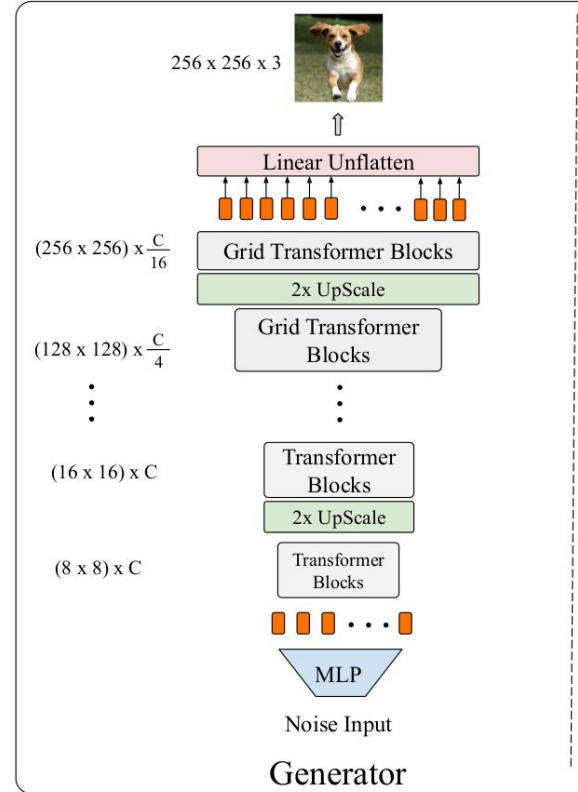
- Introduced new module of grid self-attention.
- Training recipe : Data augmentation, modified normalization, and relative position encoding.
- Competitive performance with s-o-t-a GANs using convolutional backbones.

## Grid Self-Attention



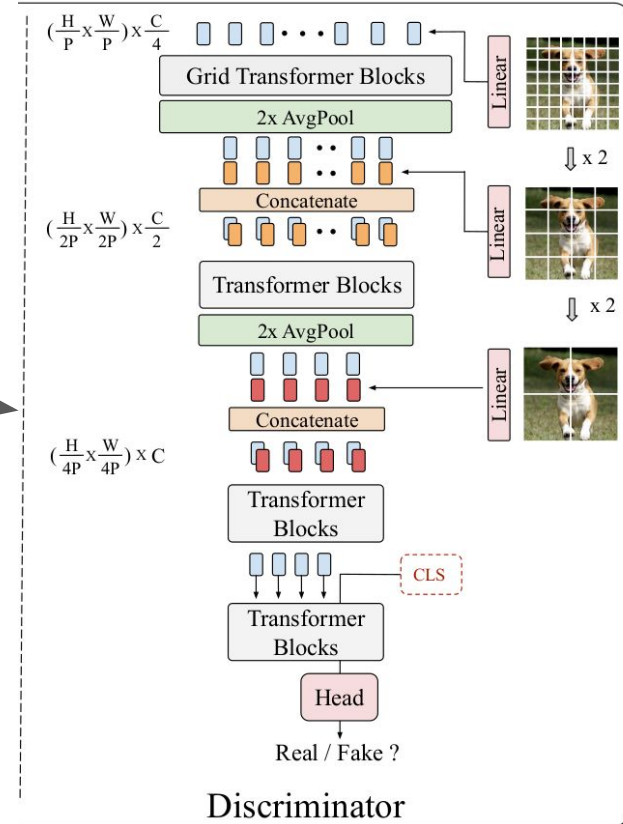
# Architecture : Generator

**Memory-friendly Generator**



# Architecture : Discriminator

## Multi-Scale Discriminator



# Architecture : Grid Self-Attention

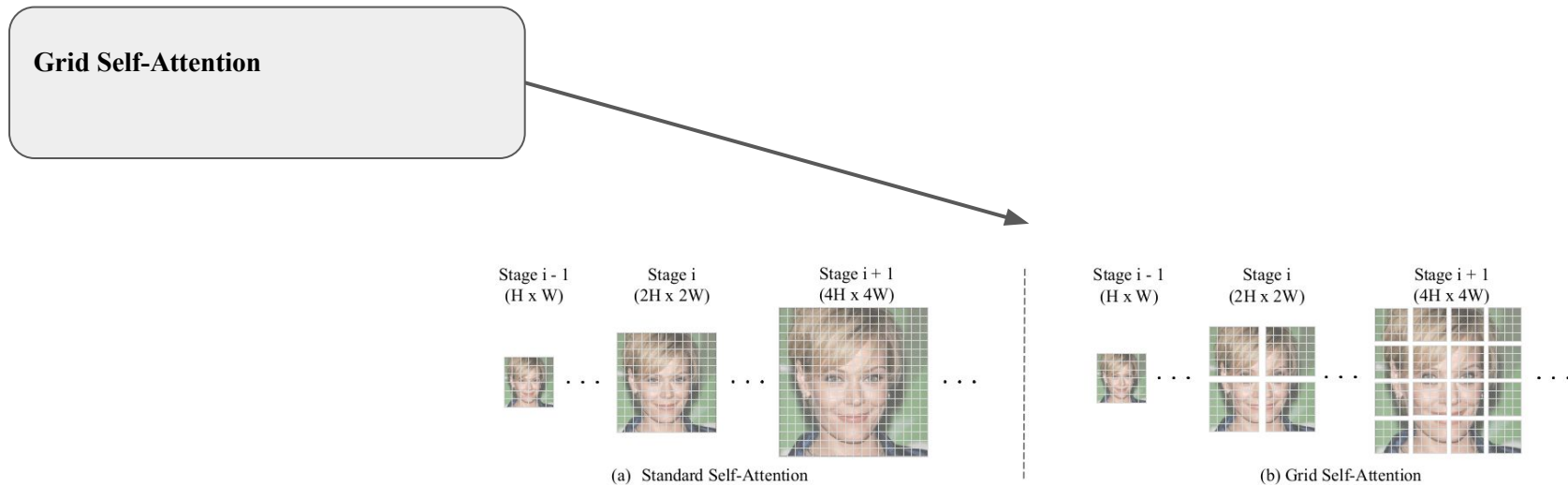


Figure 3: Grid Self-Attention across different transformer stages. We replace Standard Self-Attention with Grid Self-Attention when the resolution is higher than  $32 \times 32$  and the grid size is set to be  $16 \times 16$  by default.

# Training Recipe

## Data Augmentation

Differential augmentation with three basic operators {Translation, Cutout, Color} leads to surprising performance improvement for TransGAN, while CNN-based GANs hardly benefit from it.

## Relative Position Encoding

$$Attention(Q, K, V) = softmax(((\frac{QK^T}{\sqrt{d_k}} + E)V)$$

## Modified Normalization

$$Y = X / \sqrt{\frac{1}{C} \sum_{i=0}^{C-1} (X^i)^2 + \epsilon}, \text{ where } \epsilon = 1e - 8$$

# Results

Table 1: Unconditional image generation results on CIFAR-10, STL-10, and CelebA ( $128 \times 128$ ) dataset. We train the models with their official code if the results are unavailable, denoted as “\*”, others are all reported from references.

Methods	CIFAR-10		STL-10		CelebA
	IS $\uparrow$	FID $\downarrow$	IS $\uparrow$	FID $\downarrow$	FID $\downarrow$
WGAN-GP [1]	6.49 $\pm$ 0.09	39.68	-	-	-
SN-GAN [48]	8.22 $\pm$ 0.05	-	9.16 $\pm$ 0.12	40.1	-
AutoGAN [18]	8.55 $\pm$ 0.10	12.42	9.16 $\pm$ 0.12	31.01	-
AdversarialNAS-GAN [18]	8.74 $\pm$ 0.07	10.87	9.63 $\pm$ 0.19	26.98	-
Progressive-GAN [16]	8.80 $\pm$ 0.05	15.52	-	-	7.30
COCO-GAN [74]	-	-	-	-	5.74
StyleGAN-V2 [69]	9.18	11.07	10.21* $\pm$ 0.14	20.84*	5.59*
StyleGAN-V2 + DiffAug. [69]	<b>9.40</b>	9.89	10.31* $\pm$ 0.12	19.15*	5.40*
<b>TransGAN</b>	9.02 $\pm$ 0.12	<b>9.26</b>	<b>10.43</b> $\pm$ 0.16	<b>18.28</b>	<b>5.28</b>

other “modern” normalization layers [76-78] that need affine parameters for both mean and variances, we find that a simple re-scaling without learnable parameters suffices to stabilize TransGAN training – in fact, it makes TransGAN train better and improves the FID on some common benchmarks, such as CelebA and LSUN-Church.

Table 3: The ablation study of proposed techniques in three common dataset CelebA( $64 \times 64$ ), CelebA( $128 \times 128$ ), and LSUN Church( $256 \times 256$ ). “OOM” represents out-of-memory issue.

Training Configuration	CelebA (64x64)	CelebA (128x128)	LSUN Church (256x256)
(A). Standard Self-Attention	8.92	<b>OOM</b>	<b>OOM</b>
(B). Nyström Self-Attention [64]	13.47	17.42	39.92
(C). Axis Self-Attention [67]	12.39	13.95	29.30
(D). Grid Self-Attention	9.89	10.58	20.39
+ Multi-scale Discriminator	9.28	8.03	15.29
+ Modified Normalization	7.05	7.13	13.27
+ Relative Position Encoding	6.14	6.32	11.93
(E). Converge	<b>5.01</b>	<b>5.28</b>	<b>8.94</b>

larger than CIFAR-10, suggesting that transformer-based architectures benefit much more notably from larger-scale data than CNNs.

Table 2: The effectiveness of Data Augmentation on both CNN-based GANs and TransGAN. We use the full CIFAR-10 training set and DiffAug [69].

Methods	WGAN-GP		AutoGAN		StyleGAN-V2		TransGAN	
	IS $\uparrow$	FID $\downarrow$	IS $\uparrow$	FID $\downarrow$	IS $\uparrow$	FID $\downarrow$	IS $\uparrow$	FID $\downarrow$
Original	<b>6.49</b>	39.68	8.55	<b>12.42</b>	9.18	11.07	8.36	22.53
+ DiffAug [69]	6.29	<b>37.14</b>	<b>8.60</b>	12.72	<b>9.40</b>	<b>9.89</b>	<b>9.02</b>	<b>9.26</b>



# Examples

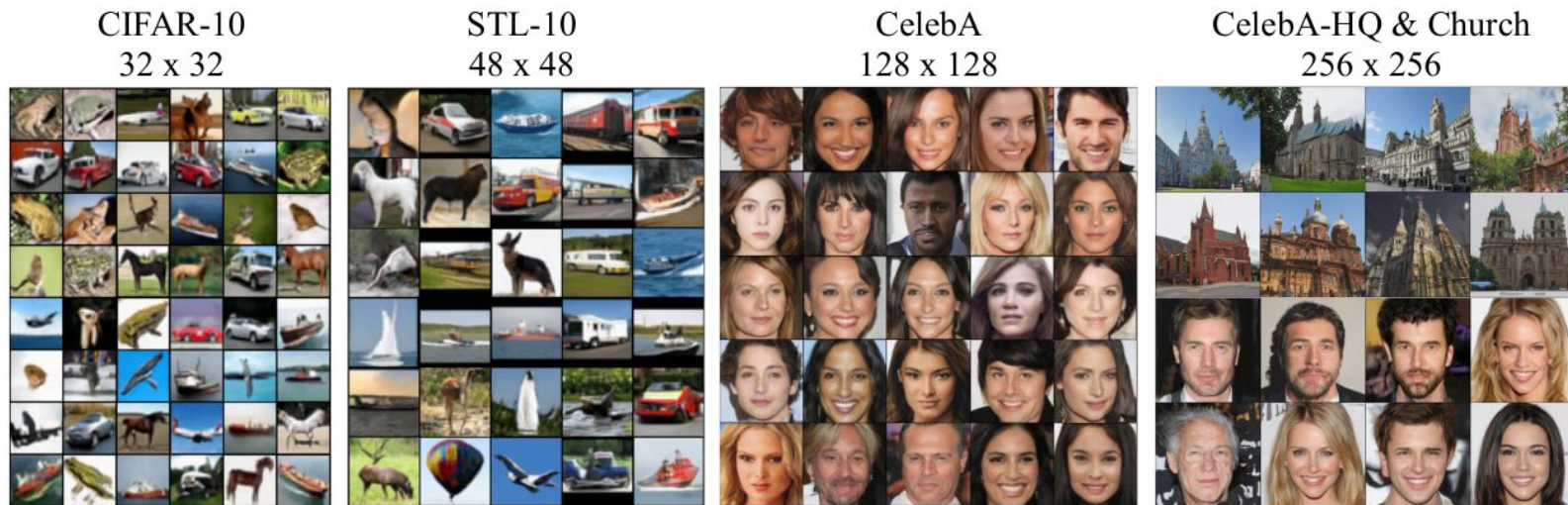


Figure 4: Representative visual results produced by TransGAN on different datasets, as resolution grows from  $32 \times 32$  to  $256 \times 256$ . More visual examples are included in Appendix [F](#)

Thank You