

GLeAD: Improving GANs with A Generator-Leading Task

– Supplementary Material –

Qingyan Bai^{1*} Ceyuan Yang² Yinghao Xu^{3*} Xihui Liu⁴ Yujiu Yang^{1†} Yujun Shen⁵

¹Tsinghua Shenzhen International Graduate School, Tsinghua University ²Shanghai AI Laboratory
³The Chinese University of Hong Kong ⁴The University of Hong Kong ⁵Ant Group

Overview

This supplementary material is organized as follows. Sec. A introduces the detailed discriminator network structure of GLeAD. Sec. B provides comparisons on the computational between our method and the baseline [1].

A. Discriminator Network Structure

Recall that, our D concludes a backbone D_{enc} , a head predicting realness scores, and a decoder h for predicting representative features f and w . Taking images whose resolution are 256×256 as an instance, the backbone D_{enc} is first employed to extract features from the input image. The very last feature map of 4×4 is sent to the scoring head to extract the realness score while the multi-level feature maps are sent to the decoder h to predict the representative features adequate for G to reconstruct the original images. As described in the submission, the representative features consist of latent codes w and the spatial representations f , which concludes a low-level representation and a high-level representation. Recall that, these spatial representations will be sent to the fixed generator to serve as the basis of the reconstruction and will be modulated by latent codes to predict the final results. We illustrate the architectures of the three aforementioned components of D in Tab. S1, Tab. S2, and Tab. S3, respectively.

B. Computational Costs

We first compute the discriminator parameter amounts of the baseline and our method. As in Tab. S4, our method merely brings 7.4% additional parameters over baseline, which is brought by the proposed lightweight design of h composed of 1×1 convolutions. Then we compare the inference time of the discriminators with a single A6000

GPU. At last, we make comparisons on the training time. We separately train the baseline model [1] and our model with 8 A100 GPUs on LSUN Church and record how much time the training costs. From the numbers in Tab. S4, we

Table S1. Network structure of the backbone D_{enc} . The output size is with order $\{C \times H \times W\}$, where C , H , and W respectively denotes the channel dimension, height and weight of the output.

Stage	Block	Output Size
input	-	$3 \times 256 \times 256$
block ₁	$\left[\begin{array}{l} 1 \times 1 \text{ Conv, } 128 \\ 2 \times 3 \times 3 \text{ Conv, } 128 \\ 1 \times 1 \text{ Conv, } 128 \\ \text{Downsample} \\ \text{LeakyReLU, } 0.2 \end{array} \right]$	$128 \times 128 \times 128$
block ₂	$\left[\begin{array}{l} 2 \times 3 \times 3 \text{ Conv, } 256 \\ 1 \times 1 \text{ Conv, } 256 \\ \text{Downsample} \\ \text{LeakyReLU, } 0.2 \end{array} \right]$	$256 \times 64 \times 64$
block ₃	$\left[\begin{array}{l} 2 \times 3 \times 3 \text{ Conv, } 512 \\ 1 \times 1 \text{ Conv, } 512 \\ \text{Downsample} \\ \text{LeakyReLU, } 0.2 \end{array} \right]$	$512 \times 32 \times 32$
block ₄	$\left[\begin{array}{l} 2 \times 3 \times 3 \text{ Conv, } 512 \\ 1 \times 1 \text{ Conv, } 512 \\ \text{Downsample} \\ \text{LeakyReLU, } 0.2 \end{array} \right]$	$512 \times 16 \times 16$
block ₅	$\left[\begin{array}{l} 2 \times 3 \times 3 \text{ Conv, } 512 \\ 1 \times 1 \text{ Conv, } 512 \\ \text{Downsample} \\ \text{LeakyReLU, } 0.2 \end{array} \right]$	$512 \times 8 \times 8$
block ₆	$\left[\begin{array}{l} 2 \times 3 \times 3 \text{ Conv, } 512 \\ 1 \times 1 \text{ Conv, } 512 \\ \text{Downsample} \\ \text{LeakyReLU, } 0.2 \end{array} \right]$	$512 \times 4 \times 4$

* This work was done during an internship at Ant Group.

† Corresponding author.

Table S2. Network structure of the decoder h predicting the low-level spatial representation, the high-level spatial representation and the 512-channel latent codes. Note that h receives multi-level features as inputs due to its feature pyramid architecture [2]. The output size is with order $\{C \times H \times W\}$.

Stage	Block	Output Size
input	—	$512 \times 32 \times 32$
		$512 \times 16 \times 16$
		$512 \times 8 \times 8$
		$512 \times 4 \times 4$
block ₁	$\begin{bmatrix} 1 \times 1 \text{ Conv, } 512 \\ \text{Upsample} \end{bmatrix}$	$512 \times 8 \times 8$
block ₂	$\begin{bmatrix} 1 \times 1 \text{ Conv, } 512 \\ \text{Upsample} \end{bmatrix}$	$512 \times 16 \times 16$
block ₃	$\begin{bmatrix} 1 \times 1 \text{ Conv, } 512 \\ \text{Upsample} \end{bmatrix}$	$512 \times 32 \times 32$
block ₄	$\begin{bmatrix} 1 \times 1 \text{ Conv, } 3 \\ 2 \times 1 \times 1 \text{ Conv, } 512 \\ \text{Downsample} \end{bmatrix}$	$3 \times 32 \times 32$ $512 \times 32 \times 32$ 512

Table S3. Network structure of the head predicting realness scores which are scalars. The output size is with order $\{C \times H \times W\}$.

Stage	Block	Output Size
input	—	$512 \times 4 \times 4$
block ₁	$\begin{bmatrix} \text{Mbstd, } 1 \\ 3 \times 3 \text{ Conv, } 512 \\ \text{LeakyReLU, } 0.2 \\ \text{Downsample} \\ \text{FC, } 512 \\ \text{LeakyReLU, } 0.2 \\ \text{FC, } 1 \end{bmatrix}$	1

Table S4. Computational cost comparisons.

Method	# params	inference time(s)	training time(h)
Baseline	24.00M	0.0184	43.83
GLeaD	25.77M	0.0219	55.78

can conclude that our method improves the synthesis quality without much additional computational burden.

References

- [1] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 1
- [2] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He,

Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017. 2