

---

# Self-Attention Generative Adversarial Networks

---

Han Zhang<sup>1,2</sup> Ian Goodfellow<sup>2</sup> Dimitris Metaxas<sup>1</sup> Augustus Odena<sup>2</sup>

## Abstract

In this paper, we propose the **Self-Attention Generative Adversarial Network (SAGAN)** which allows attention-driven, long-range dependency modeling for image generation tasks. Traditional convolutional GANs generate high-resolution details as a function of only spatially local points in lower-resolution feature maps. In SAGAN, details can be generated using cues from all feature locations. Moreover, the discriminator can check that highly detailed features in distant portions of the image are consistent with each other. Furthermore, recent work has shown that generator conditioning affects GAN performance. Leveraging this insight, we apply spectral normalization to the GAN generator and find that this improves training dynamics. The proposed SAGAN performs better than prior work<sup>1</sup>, boosting the best published Inception score from 36.8 to 52.52 and reducing Fréchet Inception distance from 27.62 to 18.65 on the challenging ImageNet dataset. Visualization of the attention layers shows that the generator leverages neighborhoods that correspond to object shapes rather than local regions of fixed shape.

## 1. Introduction

Image synthesis is an important problem in computer vision. There has been remarkable progress in this direction with the emergence of Generative Adversarial Networks (GANs) (Goodfellow et al., 2014). GANs based on deep convolutional networks (Radford et al., 2016; Karras et al., 2018; Zhang et al.) have been especially successful. However, by carefully examining the generated samples from these models, we can observe that convolutional GANs (Odena et al., 2017; Miyato et al., 2018; Miyato & Koyama, 2018) have much more difficulty in modeling some image classes than others when trained on multi-class datasets (e.g., ImageNet (Russakovsky et al., 2015)). For example, while the state-of-the-art ImageNet GAN model (Miyato & Koyama, 2018) excels at synthesizing image classes with few structural constraints (e.g., ocean, sky and landscape classes, which are distinguished more by texture than by geometry), it fails to capture geometric or structural patterns that occur consistently in some classes (for example, dogs are often drawn with realistic fur texture but without clearly defined separate feet). One possible explanation for this is that previous models rely heavily on convolution to model the dependencies across different image regions. Since the convolution operator has a local receptive field, long range dependencies can only be processed after passing through several convolutional layers. This could prevent learning about long-term dependencies for a variety of reasons: a small model may not be able to represent them, optimization algorithms may have trouble discovering parameter values that carefully coordinate multiple layers to capture these dependencies, and these parameterizations may be statistically brittle and prone to failure when applied to previously unseen inputs. Increasing the size of the convolution kernels can increase the representational capacity of the network but doing so also loses the computational and statistical efficiency obtained by using local convolutional structure. **Self-attention** (Cheng et al., 2016; Parikh et al., 2016; Vaswani et al., 2017), on the other hand, exhibits a better balance between the ability to model long-range dependencies and the computational and statistical efficiency. The self-attention module calculates response at a position as a weighted sum of the features at all positions, where the weights – or attention vectors – are calculated with only a small computational cost.

<sup>1</sup>Department of Computer Science, Rutgers University <sup>2</sup>Google Research, Brain Team. Correspondence to: Han Zhang <zhanghan@google.com>.

*Proceedings of the 36<sup>th</sup> International Conference on Machine Learning*, Long Beach, California, PMLR 97, 2019. Copyright 2019 by the author(s).

<sup>1</sup>Brock et al. (2018), which builds heavily on this work, has since improved those results substantially.

In this work, we propose Self-Attention Generative Adversarial Networks (SAGANs), which introduce a self-attention mechanism into convolutional GANs. The self-attention module is complementary to convolutions and helps with modeling long range, multi-level dependencies across image regions. Armed with self-attention, the generator can draw images in which fine details at every location are carefully coordinated with fine details in distant portions of the image. Moreover, the discriminator can also more accurately enforce complicated geometric constraints on the global image



**Figure 1.** The proposed SAGAN generates images by leveraging complementary features in distant portions of the image rather than local regions of fixed shape to generate consistent objects/scenarios. In each row, the first image shows five representative query locations with color coded dots. The other five images are attention maps for those query locations, with corresponding color coded arrows summarizing the most-attended regions.

structure.

In addition to self-attention, we also incorporate recent insights relating network conditioning to GAN performance. The work by (Odena et al., 2018) showed that well-conditioned generators tend to perform better. We propose enforcing good conditioning of GAN generators using the spectral normalization technique that has previously been applied only to the discriminator (Miyato et al., 2018).

We have conducted extensive experiments on the ImageNet dataset to validate the effectiveness of the proposed self-attention mechanism and stabilization techniques. SAGAN significantly outperforms prior work in image synthesis by boosting the best reported Inception score from **36.8 to 52.52** and reducing Fréchet Inception distance from **27.62 to 18.65**. Visualization of the attention layers shows that the generator leverages neighborhoods that correspond to object shapes rather than local regions of fixed shape. Our code is available at <https://github.com/brain-research/self-attention-gan>.

## 2. Related Work

**Generative Adversarial Networks.** GANs have achieved great success in various image generation tasks, including image-to-image translation (Isola et al., 2017; Zhu et al., 2017; Taigman et al., 2017; Liu & Tuzel, 2016; Xue et al., 2018; Park et al., 2019), image super-resolution (Ledig et al., 2017; Snderby et al., 2017) and text-to-image synthesis (Reed et al., 2016b;a; Zhang et al., 2017; Hong et al., 2018). Despite this success, the training of GANs is known to be unstable and sensitive to the choices of hyperparameters. Several works have attempted to stabilize the GAN training dynamics and improve the sample diversity by designing new network architectures (Radford et al., 2016;

Zhang et al., 2017; Karras et al., 2018; 2019), modifying the learning objectives and dynamics (Arjovsky et al., 2017; Salimans et al., 2018; Metz et al., 2017; Che et al., 2017; Zhao et al., 2017; Jolicoeur-Martineau, 2019), adding regularization methods (Gulrajani et al., 2017; Miyato et al., 2018) and introducing heuristic tricks (Salimans et al., 2016; Odena et al., 2017). Recently, Miyato et al. (Miyato et al., 2018) proposed limiting the spectral norm of the weight matrices in the discriminator in order to constrain the Lipschitz constant of the discriminator function. Combined with the projection-based discriminator (Miyato & Koyama, 2018), the spectrally normalized model greatly improves class-conditional image generation on ImageNet.

**Attention Models.** Recently, attention mechanisms have become an integral part of models that must capture global dependencies (Bahdanau et al., 2014; Xu et al., 2015; Yang et al., 2016; Gregor et al., 2015; Chen et al., 2018). In particular, self-attention (Cheng et al., 2016; Parikh et al., 2016), also called intra-attention, calculates the response at a position in a sequence by attending to all positions within the same sequence. Vaswani et al. (Vaswani et al., 2017) demonstrated that machine translation models could achieve state-of-the-art results by solely using a self-attention model. Parmar et al. (Parmar et al., 2018) proposed an Image Transformer model to add self-attention into an autoregressive model for image generation. Wang et al. (Wang et al., 2018) formalized self-attention as a non-local operation to model the spatial-temporal dependencies in video sequences. In spite of this progress, self-attention has not yet been explored in the context of GANs. (AttnGAN (Xu et al., 2018) uses attention over word embeddings within an *input* sequence, but not self-attention over *internal model states*). SAGAN learns to efficiently find global, long-range dependencies within internal representations of images.

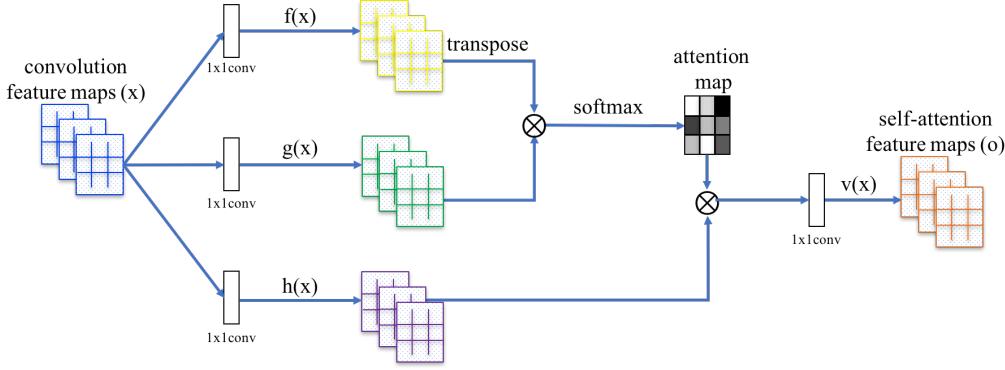


Figure 2. The proposed self-attention module for the SAGAN. The  $\otimes$  denotes matrix multiplication. The softmax operation is performed on each row.

### 3. Self-Attention Generative Adversarial Networks

Most GAN-based models (Radford et al., 2016; Salimans et al., 2016; Karras et al., 2018) for image generation are built using convolutional layers. Convolution processes the information in a local neighborhood, thus using convolutional layers alone is computationally inefficient for modeling long-range dependencies in images. In this section, we adapt the non-local model of (Wang et al., 2018) to introduce self-attention to the GAN framework, enabling both the generator and the discriminator to efficiently model relationships between widely separated spatial regions. We call the proposed method Self-Attention Generative Adversarial Networks (SAGAN) because of its self-attention module (see Figure 2).

The image features from the previous hidden layer  $x \in \mathbb{R}^{C \times N}$  are first transformed into two feature spaces  $f, g$  to calculate the attention, where  $f(x) = W_f x$ ,  $g(x) = W_g x$

$$\beta_{j,i} = \frac{\exp(s_{ij})}{\sum_{i=1}^N \exp(s_{ij})}, \text{ where } s_{ij} = \mathbf{f}(x_i)^T \mathbf{g}(x_j), \quad (1)$$

and  $\beta_{j,i}$  indicates the extent to which the model attends to the  $i^{th}$  location when synthesizing the  $j^{th}$  region. Here,  $C$  is the number of channels and  $N$  is the number of feature locations of features from the previous hidden layer. The output of the attention layer is  $\mathbf{o} = (\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_j, \dots, \mathbf{o}_N) \in \mathbb{R}^{C \times N}$ , where,

$$\mathbf{o}_j = \mathbf{v} \left( \sum_{i=1}^N \beta_{j,i} \mathbf{h}(x_i) \right), \quad \mathbf{h}(x_i) = W_h x_i, \quad \mathbf{v}(x_i) = W_v x_i \quad (2)$$

In the above formulation,  $W_g \in \mathbb{R}^{\bar{C} \times C}$ ,  $W_f \in \mathbb{R}^{\bar{C} \times C}$ ,  $W_h \in \mathbb{R}^{\bar{C} \times C}$ , and  $W_v \in \mathbb{R}^{C \times \bar{C}}$  are the learned weight matrices, which are implemented as  $1 \times 1$  convolutions. Since

We did not notice any significant performance decrease when reducing the channel number of  $C$  to be  $C/k$ , where  $k = 1, 2, 4, 8$  after few training epochs on ImageNet. For memory efficiency, we choose  $k = 8$  (i.e.,  $\bar{C} = C/8$ ) in all our experiments.

In addition, we further multiply the output of the attention layer by a scale parameter and add back the input feature map. Therefore, the final output is given by,

$$y_i = \gamma o_i + x_i, \quad (3)$$

where  $\gamma$  is a learnable scalar and it is initialized as 0. Introducing the learnable  $\gamma$  allows the network to first rely on the cues in the local neighborhood – since this is easier – and then gradually learn to assign more weight to the non-local evidence. The intuition for why we do this is straightforward: we want to learn the easy task first and then progressively increase the complexity of the task. In the SAGAN, the proposed attention module has been applied to both the generator and the discriminator, which are trained in an alternating fashion by minimizing the hinge version of the adversarial loss (Lim & Ye, 2017; Tran et al., 2017; Miyato et al., 2018),

$$\begin{aligned} L_D = & -\mathbb{E}_{(x,y) \sim p_{data}} [\min(0, -1 + D(x, y))] \\ & -\mathbb{E}_{z \sim p_z, y \sim p_{data}} [\min(0, -1 - D(G(z), y))], \quad (4) \\ L_G = & -\mathbb{E}_{z \sim p_z, y \sim p_{data}} D(G(z), y), \end{aligned}$$

### 4. Techniques to Stabilize the Training of GANs

We also investigate two techniques to stabilize the training of GANs on challenging datasets. First, we use spectral normalization (Miyato et al., 2018) in the generator as well as in the discriminator. Second, we confirm that the two-timescale update rule (TTUR) (Heusel et al., 2017) is effective, and we advocate using it specifically to address slow learning in regularized discriminators.

#### 4.1. Spectral normalization for both generator and discriminator

Miyato *et al.* (Miyato et al., 2018) originally proposed stabilizing the training of GANs by applying spectral normalization to the discriminator network. Doing so constrains the Lipschitz constant of the discriminator by restricting the spectral norm of each layer. Compared to other normalization techniques, spectral normalization does not require extra hyper-parameter tuning (setting the spectral norm of all weight layers to 1 consistently performs well in practice). Moreover, the computational cost is also relatively small.

We argue that the generator can also benefit from spectral normalization, based on recent evidence that the conditioning of the generator is an important causal factor in GANs’ performance (Odena et al., 2018). Spectral normalization in the generator can prevent the escalation of parameter magnitudes and avoid unusual gradients. We find empirically that spectral normalization of both generator and discriminator makes it possible to use fewer discriminator updates per generator update, thus significantly reducing the computational cost of training. The approach also shows more stable training behavior.

#### 4.2. Imbalanced learning rate for generator and discriminator updates

In previous work, regularization of the discriminator (Miyato et al., 2018; Gulrajani et al., 2017) often slows down the GANs’ learning process. In practice, methods using regularized discriminators typically require multiple (*e.g.*, 5) discriminator update steps per generator update step during training. Independently, Heusel *et al.* (Heusel et al., 2017) have advocated using separate learning rates (TTUR) for the generator and the discriminator. We propose using TTUR specifically to compensate for the problem of slow learning in a regularized discriminator, making it possible to use fewer discriminator steps per generator step. Using this approach, we are able to produce better results given the same wall-clock time.

### 5. Experiments

To evaluate the proposed methods, we conducted extensive experiments on the LSVRC2012 (ImageNet) dataset (Russakovsky et al., 2015). First, in Section 5.1, we present experiments designed to evaluate the effectiveness of the two proposed techniques for stabilizing GANs’ training. Next, the proposed self-attention mechanism is investigated in Section 5.2. Finally, our SAGAN is compared with state-of-the-art methods (Odena et al., 2017; Miyato & Koyama, 2018) on the image generation task in Section 5.3.

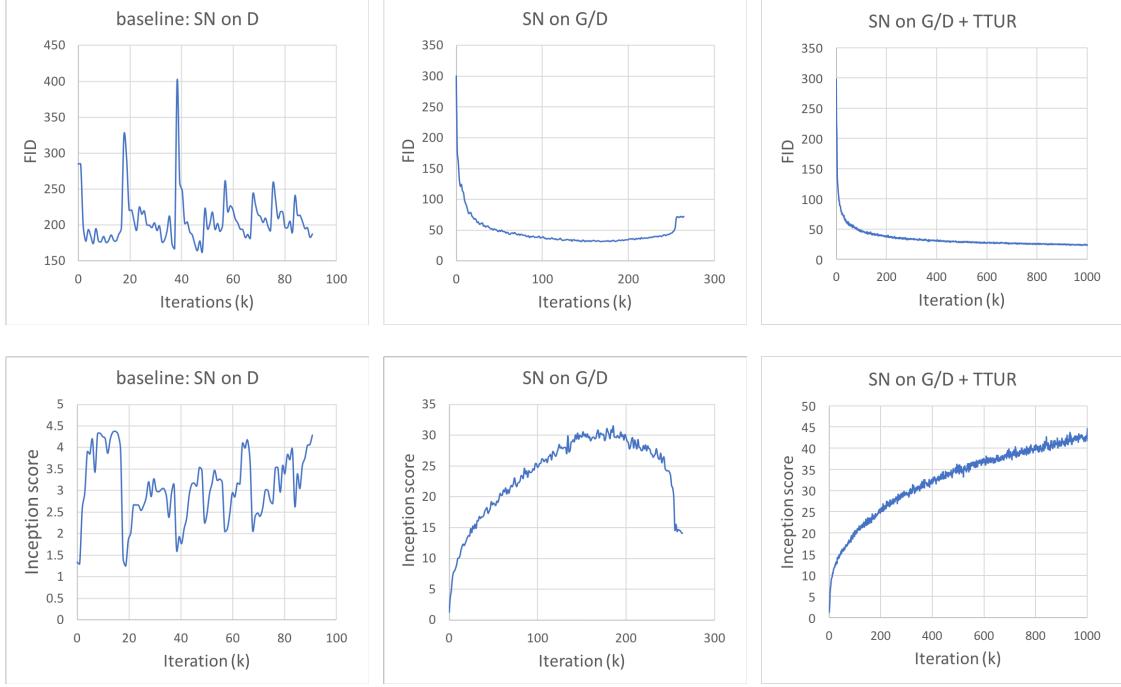
**Evaluation metrics.** We choose the Inception score (IS) (Salimans et al., 2016) and the Fréchet Inception dis-

tance (FID) (Heusel et al., 2017) for quantitative evaluation. The Inception score (Salimans et al., 2016) computes the KL divergence between the conditional class distribution and the marginal class distribution. Higher Inception score indicates better image quality. We include the Inception score because it is widely used and thus makes it possible to compare our results to previous work. However, it is important to understand that Inception score has serious limitations—it is intended primarily to ensure that the model generates samples that can be confidently recognized as belonging to a specific class, and that the model generates samples from many classes, not necessarily to assess realism of details or intra-class diversity. FID is a more principled and comprehensive metric, and has been shown to be more consistent with human evaluation in assessing the realism and variation of the generated samples (Heusel et al., 2017). FID calculates the Wasserstein-2 distance between the generated images and the real images in the feature space of an Inception-v3 network. Besides the FID calculated over the whole data distribution (*i.e.*, all 1000 classes of images in ImageNet), we also compute FID between the generated images and dataset images within each class (called intra FID (Miyato & Koyama, 2018)). Lower FID and intra FID values mean closer distances between synthetic and real data distributions. In all our experiments, 50k samples are randomly generated for each model to compute the Inception score, FID and intra FID.

**Network structures and implementation details.** All the SAGAN models we train are designed to generate  $128 \times 128$  images. By default, spectral normalization (Miyato et al., 2018) is used for the layers in both the generator and the discriminator. Similar to (Miyato & Koyama, 2018), SAGAN uses conditional batch normalization in the generator and projection in the discriminator. For all models, we use the Adam optimizer (Kingma & Ba, 2015) with  $\beta_1 = 0$  and  $\beta_2 = 0.9$  for training. By default, the learning rate for the discriminator is 0.0004 and the learning rate for the generator is 0.0001.

#### 5.1. Evaluating the proposed stabilization techniques

In this section, experiments are conducted to evaluate the effectiveness of the proposed stabilization techniques, *i.e.*, applying spectral normalization (SN) to the generator and utilizing imbalanced learning rates (TTUR). In Figure 3, our models “SN on G/D” and “SN on G/D+TTUR” are compared with a baseline model, which is implemented based on the state-of-the-art image generation method (Miyato et al., 2018). In this baseline model, SN is only utilized in the discriminator. When we train it with 1:1 balanced updates for the discriminator ( $D$ ) and the generator ( $G$ ), the training becomes very unstable, as shown in the leftmost sub-figures of Figure 3. It exhibits mode collapse very early in training. For example, the top-left sub-figure of Figure 4



**Figure 3.** Training curves for the baseline model and our models with the proposed stabilization techniques, “SN on G/D” and two-timescale learning rates (TTUR). All models are trained with 1:1 balanced updates for  $G$  and  $D$ .

illustrates some images randomly generated by the baseline model at the 10k-th iteration. Although in the the original paper (Miyato et al., 2018) this unstable training behavior is greatly mitigated by using 5:1 imbalanced updates for  $D$  and  $G$ , the ability to be stably trained with 1:1 balanced updates is desirable for improving the convergence speed of the model. Thus, using our proposed techniques means that the model can produce better results given the same wall-clock time. Given this, there is no need to search for a suitable update ratio for the generator and discriminator. As shown in the middle sub-figures of Figure 3, adding SN to both the generator and the discriminator greatly stabilized our model “SN on G/D”, even when it was trained with 1:1 balanced updates. However, the quality of samples does not improve monotonically during training. For example, the image quality as measured by FID and IS is starting to drop at the 260k-th iteration. Example images randomly generated by this model at different iterations can be found in Figure 4. When we also apply the imbalanced learning rates to train the discriminator and the generator, the quality of images generated by our model “SN on G/D+TTUR” improves monotonically during the whole training process. As shown in Figure 3 and Figure 4, we do not observe any significant decrease in sample quality or in the FID or the Inception score during one million training iterations. Thus, both quantitative results and qualitative results demonstrate the effectiveness of the proposed stabilization techniques

for GANs’ training. They also demonstrate that the effect of the two techniques is at least partly additive. In the rest of experiments, all models use spectral normalization for both the generator and discriminator and use the imbalanced learning rates to train the generator and the discriminator with 1:1 updates.

## 5.2. Self-attention mechanism.

To explore the effect of the proposed self-attention mechanism, we build several SAGAN models by adding the self-attention mechanism to different stages of the generator and the discriminator. As shown in Table 1, the SAGAN models with the self-attention mechanism at the middle-to-high level feature maps (e.g.,  $feat_{32}$  and  $feat_{64}$ ) achieve better performance than the models with the self-attention mechanism at the low level feature maps (e.g.,  $feat_8$  and  $feat_{16}$ ). For example, the FID of the model “SAGAN,  $feat_8$ ” is improved from 22.98 to 18.28 by “SAGAN,  $feat_{32}$ ”. The reason is that self-attention receives more evidence and enjoys more freedom to choose conditions with larger feature maps (*i.e.*, it is complementary to convolution for large feature maps), however, it plays a similar role as the local convolution when modeling dependencies for small (*e.g.*,  $8 \times 8$ ) feature maps. It demonstrates that the attention mechanism gives more power to both the generator and the discriminator to directly model the long-range dependencies in the feature maps. In addition, the comparison of our SAGAN



Figure 4. 128×128 examples randomly generated by the baseline model and our models “SN on  $G/D$ ” and “SN on  $G/D+TTUR$ ”.

Model	no attention	SAGAN				Residual			
		$feat_8$	$feat_{16}$	$feat_{32}$	$feat_{64}$	$feat_8$	$feat_{16}$	$feat_{32}$	$feat_{64}$
FID	22.96	22.98	22.14	<b>18.28</b>	18.65	42.13	22.40	27.33	28.82
IS	42.87	43.15	45.94	51.43	<b>52.52</b>	23.17	44.49	38.50	38.96

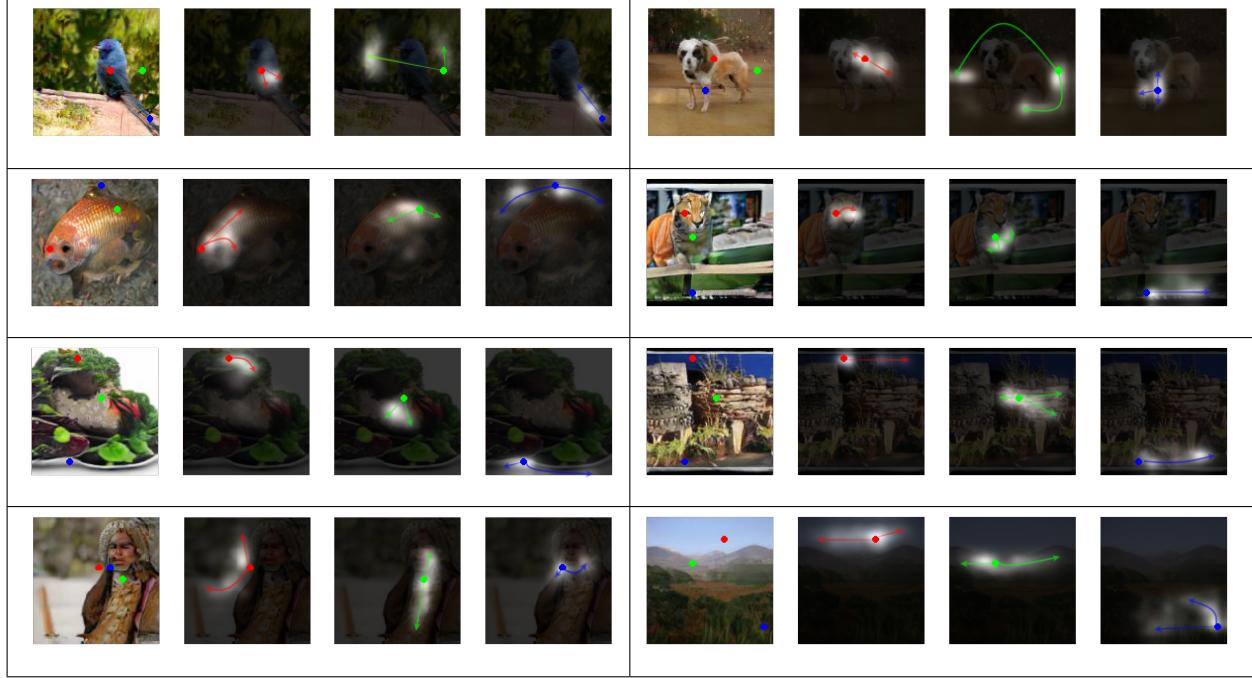
Table 1. Comparison of Self-Attention and Residual block on GANs. These blocks are added into different layers of the network. All models have been trained for one million iterations, and the best Inception scores (IS) and Fréchet Inception distance (FID) are reported.  $feat_k$  means adding self-attention to the  $k \times k$  feature maps.

and the baseline model without attention (2nd column of Table 1) further shows the effectiveness of the proposed self-attention mechanism.

Compared with residual blocks with the same number of parameters, the self-attention blocks also achieve better results. For example, the training is not stable when we replace the self-attention block with the residual block in  $8 \times 8$  feature maps, which leads to a significant decrease in performance (*e.g.*, FID increases from 22.98 to 42.13). Even for the cases when the training goes smoothly, replacing the self-attention block with the residual block still leads to worse results in terms of FID and Inception score. (*e.g.*, FID 18.28 vs 27.33 in feature map  $32 \times 32$ ). This comparison demonstrates that the performance improvement given by using SAGAN is not simply due to an increase in model depth and capacity.

To better understand what has been learned during the generation process, we visualize the attention weights of the generator in SAGAN for different images. Some sample images with attention are shown in Figure 5 and Figure 1.

We observe that the network learns to allocate attention according to similarity of color and texture, rather than just spatial adjacency. For example, in the top-left cell of Figure 1, the red point attends mostly to the body of the bird around it, however, the green point learns to attend to other side of the image. In this way, the image has a consistent background (*i.e.*, trees from the left to the right though they are separated by the bird). Similarly, the blue point allocates the attention to the whole tail of the bird to make the generated part coherent. Those long-range dependencies could not be captured by convolutions with local receptive fields. We also find that although some query points are quite close in spatial location, their attention maps can be very different, as shown in the bottom-left cell. The red point attends mostly to the background regions, whereas the blue point, though adjacent to red point, puts most of the attention on the foreground object. This also reduces the chance for the local errors to propagate, since the adjacent position has the freedom to choose to attend to other distant locations. These observations further demonstrate that self-attention



**Figure 5.** Visualization of attention maps. These images were generated by SAGAN. We visualize the attention maps of the last generator layer that used attention, since this layer is the closest to the output pixels and is the most straightforward to project into pixel space and interpret. In each cell, the first image shows three representative query locations with color coded dots. The other three images are attention maps for those query locations, with corresponding color coded arrows summarizing the most-attended regions. We observe that the network learns to allocate attention according to similarity of color and texture, rather than just spatial adjacency (see the top-left cell). We also find that although some query points are quite close in spatial location, their attention maps can be very different, as shown in the bottom-left cell. As shown in the top-right cell, SAGAN is able to draw dogs with clearly separated legs. The blue query point shows that attention helps to get the structure of the joint area correct. See the text for more discussion about the properties of learned attention maps.

is complementary to convolutions for image generation in GANs. As shown in the top-right cell, SAGAN is able to draw dogs with clearly separated legs. The blue query point shows that attention helps to get the structure of the joint area correct.

### 5.3. Comparison with the state-of-the-art

Our SAGAN is also compared with the state-of-the-art GAN models (Odena et al., 2017; Miyato & Koyama, 2018) for class conditional image generation on ImageNet. As shown in Table 2, our proposed SAGAN achieves the best Inception score, intra FID and FID. The proposed SAGAN significantly improves the best published Inception score from 36.8 to 52.52. The lower FID (18.65) and intra FID (83.7) achieved by the SAGAN also indicates that the SAGAN can better approximate the original image distribution by using the self-attention module to model the long-range dependencies between image regions.

Figure 6 shows some comparison results and generated images for representative classes of ImageNet. We observe that our SAGAN achieves much better performance (*i.e.*, lower intra FID) than the state-of-the-art GAN model (Miy-

ato & Koyama, 2018) for synthesizing image classes with complex geometric or structural patterns, such as goldfish and Saint Bernard. For classes with few structural constraints (*e.g.*, valley, stone wall and coral fungus, which are distinguished more by texture than by geometry), our SAGAN shows less superiority compared with the baseline model (Miyato & Koyama, 2018). Again, the reason is that the self-attention in SAGAN is complementary to the convolution for capturing long-range, global-level dependencies occurring consistently in geometric or structural patterns, but plays a similar role as the local convolution when modeling dependencies for simple texture.

## 6. Conclusion

In this paper, we proposed Self-Attention Generative Adversarial Networks (SAGANs), which incorporate a self-attention mechanism into the GAN framework. The self-attention module is effective in modeling long-range dependencies. In addition, we show that spectral normalization applied to the generator stabilizes GAN training and that TTUR speeds up training of regularized discriminators. SAGAN achieves the state-of-the-art performance on class-conditional image generation on ImageNet.

Model	Inception Score	Intra FID	FID
AC-GAN (Odena et al., 2017)	28.5	260.0	/
SNGAN-projection (Miyato & Koyama, 2018)	36.8	92.4	27.62*
<b>SAGAN</b>	<b>52.52</b>	<b>83.7</b>	<b>18.65</b>

Table 2. Comparison of the proposed SAGAN with state-of-the-art GAN models (Odena et al., 2017; Miyato & Koyama, 2018) for class conditional image generation on ImageNet. FID of SNGAN-projection is calculated from officially released weights.



Figure 6. 128x128 example images generated by SAGAN for different classes. Each row shows examples from one class. In the leftmost column, the intra FID of our SAGAN (left) and the state-of-the-art method (Miyato & Koyama, 2018) (right) are listed.

## Acknowledgments

We thank Surya Bhupatiraju for feedback on drafts of this article. We also thank David Berthelot and Tom B. Brown for help with implementation details. Finally, we thank Jakob Uszkoreit, Tao Xu, and Ashish Vaswani for helpful discussions.

## References

- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein GAN. *arXiv:1701.07875*, 2017.
- Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate. *arXiv:1409.0473*, 2014.
- Brock, A., Donahue, J., and Simonyan, K. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- Che, T., Li, Y., Jacob, A. P., Bengio, Y., and Li, W. Mode regularized generative adversarial networks. In *ICLR*, 2017.
- Chen, X., Mishra, N., Rohaninejad, M., and Abbeel, P. Pixelsnail: An improved autoregressive generative model. In *ICML*, 2018.
- Cheng, J., Dong, L., and Lapata, M. Long short-term memory-networks for machine reading. In *EMNLP*, 2016.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. C., and Bengio, Y. Generative adversarial nets. In *NIPS*, 2014.
- Gregor, K., Danihelka, I., Graves, A., Rezende, D. J., and Wierstra, D. DRAW: A recurrent neural network for image generation. In *ICML*, 2015.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. C. Improved training of wasserstein GANs. In *NIPS*, 2017.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *NIPS*, pp. 6629–6640, 2017.
- Hong, S., Yang, D., Choi, J., and Lee, H. Inferring semantic layout for hierarchical text-to-image synthesis. In *CVPR*, 2018.
- Isola, P., Zhu, J.-Y., Zhou, T., and Efros, A. A. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017.
- Jolicoeur-Martineau, A. The relativistic discriminator: a key element missing from standard GAN. In *ICLR*, 2019.
- Karras, T., Aila, T., Laine, S., and Lehtinen, J. Progressive growing of GANs for improved quality, stability, and variation. In *ICLR*, 2018.
- Karras, T., Laine, S., and Aila, T. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- Ledig, C., Theis, L., Huszar, F., Caballero, J., Aitken, A., Tejani, A., Totz, J., Wang, Z., and Shi, W. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017.
- Lim, J. H. and Ye, J. C. Geometric GAN. *arXiv:1705.02894*, 2017.
- Liu, M. and Tuzel, O. Coupled generative adversarial networks. In *NIPS*, 2016.
- Metz, L., Poole, B., Pfau, D., and Sohl-Dickstein, J. Unrolled generative adversarial networks. In *ICLR*, 2017.
- Miyato, T. and Koyama, M. cGANs with projection discriminator. In *ICLR*, 2018.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. In *ICLR*, 2018.
- Odena, A., Olah, C., and Shlens, J. Conditional image synthesis with auxiliary classifier GANs. In *ICML*, 2017.
- Odena, A., Buckman, J., Olsson, C., Brown, T. B., Olah, C., Raffel, C., and Goodfellow, I. Is generator conditioning causally related to GAN performance? In *ICML*, 2018.
- Parikh, A. P., Täckström, O., Das, D., and Uszkoreit, J. A decomposable attention model for natural language inference. In *EMNLP*, 2016.
- Park, T., Liu, M., Wang, T., and Zhu, J. Semantic image synthesis with spatially-adaptive normalization. In *CVPR*, 2019.

- Parmar, N., Vaswani, A., Uszkoreit, J., ukasz Kaiser, Shazeer, N., and Ku, A. Image transformer. *arXiv:1802.05751*, 2018.
- Radford, A., Metz, L., and Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016.
- Reed, S., Akata, Z., Mohan, S., Tenka, S., Schiele, B., and Lee, H. Learning what and where to draw. In *NIPS*, 2016a.
- Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., and Lee, H. Generative adversarial text-to-image synthesis. In *ICML*, 2016b.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. ImageNet large scale visual recognition challenge. *IJCV*, 2015.
- Salimans, T., Goodfellow, I. J., Zaremba, W., Cheung, V., Radford, A., and Chen, X. Improved techniques for training GANs. In *NIPS*, 2016.
- Salimans, T., Zhang, H., Radford, A., and Metaxas, D. N. Improving GANs using optimal transport. In *ICLR*, 2018.
- Snderby, C. K., Caballero, J., Theis, L., Shi, W., and Huszar, F. Amortised map inference for image super-resolution. In *ICLR*, 2017.
- Taigman, Y., Polyak, A., and Wolf, L. Unsupervised cross-domain image generation. In *ICLR*, 2017.
- Tran, D., Ranganath, R., and Blei, D. M. Deep and hierarchical implicit models. *arXiv:1702.08896*, 2017.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *arXiv:1706.03762*, 2017.
- Wang, X., Girshick, R., Gupta, A., and He, K. Non-local neural networks. In *CVPR*, 2018.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A. C., Salakhutdinov, R., Zemel, R. S., and Bengio, Y. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.
- Xu, T., Zhang, P., Huang, Q., Zhang, H., Gan, Z., Huang, X., and He, X. AttnGAN: Fine-grained text to image generation with attentional generative adversarial networks. In *CVPR*, 2018.
- Xue, Y., Xu, T., Zhang, H., Long, L. R., and Huang, X. SegAN: Adversarial network with multi-scale L1 loss for medical image segmentation. *Neuroinformatics*, pp. 1–10, 2018.
- Yang, Z., He, X., Gao, J., Deng, L., and Smola, A. J. Stacked attention networks for image question answering. In *CVPR*, 2016.
- Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., and Metaxas, D. N. StackGAN++: Realistic image synthesis with stacked generative adversarial networks. *TPAMI*.
- Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., and Metaxas, D. StackGAN: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV*, 2017.
- Zhao, J., Mathieu, M., and LeCun, Y. Energy-based generative adversarial network. In *ICLR*, 2017.
- Zhu, J.-Y., Park, T., Isola, P., and Efros, A. A. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017.