

Spectral Normalization on GANs

Bingbing Hu

July 17, 2019

SIST, ShanghaiTech

Outline

1. Motivation

2. Method

3. Experiments

4. Conclusion

Motivation

Motivation

GAN training is unstable!

Motivation

GAN training is unstable!

- Discriminator can be crucial to the GAN training
- There exists a discriminator that can perfectly distinguish the model distribution from target
- Such discriminator can cause zero gradient, hence training collapse

Motivation

GAN training is unstable!

- Discriminator can be crucial to the GAN training
- There exists a discriminator that can perfectly distinguish the model distribution from target
- Such discriminator can cause zero gradient, hence training collapse

All of these motivate us to bring some constraints to discriminator.

Method

GAN Structure

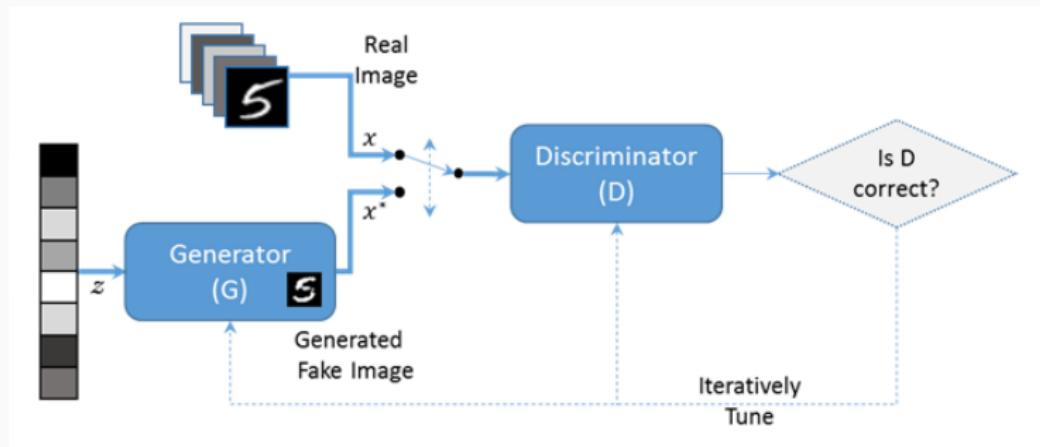


Figure 1: GAN Structure

GAN Structure

The standard formulation of GAN is given by

$$\min_G \max_D V(G, D)$$

The conventional form of $V(G, D)$ is given by

$$\mathbb{E}_{\mathbf{x} \sim q_{\text{data}}} [\log D(\mathbf{x})] + \mathbb{E}_{\hat{\mathbf{x}} \sim p_G} [\log(1 - D(G(\hat{\mathbf{x}})))] \quad (1)$$

Why Unstable?

For the conventional form of $V(G, D)$ (1) and a fixed generator G , it is known that [1], the optimal discriminator is given by

$$D_G^*(\mathbf{x}) = \frac{q_{\text{data}}(\mathbf{x})}{q_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})}.$$

This gives

$$D_G^*(\mathbf{x}) = \sigma(f^*(\mathbf{x})),$$

where $f^*(\mathbf{x}) = \log q_{\text{data}}(\mathbf{x}) - \log p_G(\mathbf{x})$

$$\implies \nabla_{\mathbf{x}} f^*(\mathbf{x}) = \frac{1}{q_{\text{data}}(\mathbf{x})} \nabla_{\mathbf{x}} q_{\text{data}}(\mathbf{x}) - \frac{1}{p_G(\mathbf{x})} \nabla_{\mathbf{x}} p_G(\mathbf{x}).$$

The derivate of $f^*(\mathbf{x})$ can be unbounded or even incomputable.

Method

Formulate the discriminator as follows:

$$f(\mathbf{x}, \theta) = W^{L+1} a_L(W^L(a_{L-1}(W^{L-1}(\dots a_1(W^1 \mathbf{x}) \dots))))$$

- \mathbf{x} : input of the network
- $\theta := \{W^1, \dots, W^L, W^{L+1}\}$: learning parameters set
- a_l : element-wise non-linear function

Method

Formulate the discriminator as follows:

$$f(\mathbf{x}, \theta) = W^{L+1} a_L(W^L(a_{L-1}(W^{L-1}(\dots a_1(W^1 \mathbf{x}) \dots))))$$

- \mathbf{x} : input of the network
- $\theta := \{W^1, \dots, W^L, W^{L+1}\}$: learning parameters set
- a_l : element-wise non-linear function

The final output of the discriminator is given by

$$D(\mathbf{x}, \theta) = \mathcal{A}(f(\mathbf{x}, \theta))$$

Spectral Normalization

The spectral normalization controls the Lipschitz constant of the discriminator function f by iteratively constraining the spectral norm of each layer $g : \mathbf{h}_{in} \mapsto \mathbf{h}_{out}$.

Definition 1

The Lipschitz norm $\|f\|_{\text{Lip}}$ of function f is smallest value M such that

$$\frac{\|f(\mathbf{x}) - f(\mathbf{x}')\|_2}{\|\mathbf{x} - \mathbf{x}'\|_2} \leq M$$

for any $\mathbf{x} \neq \mathbf{x}'$.

Spectral Normalization

Definition 2

The spectral norm of a matrix \mathbf{A} is defined by

$$\sigma(\mathbf{A}) := \sup_{\mathbf{h}: \mathbf{h} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{h}\|_2}{\|\mathbf{h}\|_2} = \sup_{\|\mathbf{h}\|_2 \leq 1} \|\mathbf{A}\mathbf{h}\|_2$$

- spectral norm is actually matrix 2-norm
- spectral norm is the max singular value of \mathbf{A}

Spectral Normalization

Therefore, for a linear layer $g(\mathbf{h}) = W\mathbf{h}$,

$$\|g\|_{\text{Lip}} = \sup_{\mathbf{h}} \sigma(\nabla g(\mathbf{h})) = \sup_{\mathbf{h}} \sigma(W) = \sigma(W)$$

If $\|a_l\|_{\text{Lip}} = 1$ ¹, we can use the inequality

$$\|g_1 \circ g_2\|_{\text{Lip}} \leq \|g_1\|_{\text{Lip}} \cdot \|g_2\|_{\text{Lip}}$$

to get

$$\begin{aligned} \|f\|_{\text{Lip}} &\leq \|(\mathbf{h}_L \mapsto W^{L+1}\mathbf{h}_L)\|_{\text{Lip}} \cdot \|a_L\|_{\text{Lip}} \cdot \\ &\quad \|(\mathbf{h}_{L-1} \mapsto W^L\mathbf{h}_{L-1})\|_{\text{Lip}} \cdots \|a_1\|_{\text{Lip}} \cdot \|(\mathbf{h}_0 \mapsto W^1\mathbf{h}_0)\|_{\text{Lip}} \\ &= \prod_{l=1}^{L+1} \|(\mathbf{h}_{l-1} \mapsto W^l\mathbf{h}_{l-1})\|_{\text{Lip}} = \prod_{l=1}^{L+1} \sigma(W^l). \end{aligned}$$

¹e.g., ReLU and leaky ReLU

Spectral Normalization

Normalize the spectral norm of the weight matrix W such that $\sigma(W) = 1$:

$$\bar{W}_{\text{SN}} := W/\sigma(W). \quad (2)$$

If we normalize each W^l using (2), together with the fact $\sigma(\bar{W}_{\text{SN}}) = 1$, we see that $\|f\|_{\text{Lip}}$ is bounded above by 1.

Power Iteration

How to compute the max singular value $\sigma(W)$?

- SVD? – computation inefficient
- Power iteration – ok

Power Iteration

Brief view of power method. Suppose $A \in \mathbf{R}^{n \times n}$ is diagonalizable, with eigenvalues ordered that

$$|\lambda_1| \geq |\lambda_2| \geq \cdots \geq |\lambda_n|.$$

Then we have $A = V\Lambda V^{-1}$ where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$, now

$$\begin{aligned} A^k &= V\Lambda^k V^{-1} \\ \implies A^k V &= V\Lambda^k. \end{aligned}$$

Randomly chose $\mathbf{x} = V\tilde{\mathbf{x}} \in \mathbf{R}^n$ and consider

$$\begin{aligned} A^k \mathbf{x} &= A^k V \tilde{\mathbf{x}} = V \Lambda^k \tilde{\mathbf{x}} = \sum_{j=1}^n \mathbf{v}_j \lambda_j^k \tilde{x}_j \\ &= \lambda_1^k \sum_{j=1}^n \left(\frac{\lambda_j}{\lambda_1} \right)^k \tilde{x}_j \mathbf{v}_j \rightarrow \lambda_1^k \tilde{x}_1 \mathbf{v}_1. \end{aligned}$$

Power Iteration

If $|\lambda_1| > |\lambda_2|$, $A^k \mathbf{x} \rightarrow c\mathbf{v}_1$ for sufficient large k . This is the idea behind power iteration:

$$\mathbf{x}^{(k+1)} = \frac{A\mathbf{x}^{(k)}}{\|A\mathbf{x}^{(k)}\|} = \frac{A^k \mathbf{x}^{(0)}}{\|A^k \mathbf{x}^{(0)}\|}.$$

As \mathbf{v}_1 is the eigenvector corresponding to the max eigenvalue, noting that V is orthonormal, we have $\mathbf{x}^{(k)} \rightarrow \mathbf{v}_1$. So the max eigen value of A is given by

$$\lambda_1 = \mathbf{v}_1^T A \mathbf{v}_1 = \left(\mathbf{x}^{(k)}\right)^T \mathbf{x}^{(k+1)}.$$

Gradient Analysis

Recall that

$$\bar{W}_{\text{SN}} := W/\sigma(W).$$

The gradient of the normalized weight matrix w.r.t. W_{ij} is

$$\frac{\partial \bar{W}_{\text{SN}}}{\partial W_{ij}} = \frac{E_{ij}}{\sigma(W)} - \frac{W}{\sigma^2(W)} \frac{\partial \sigma(W)}{\partial W_{ij}} \quad (3)$$

$$= \frac{E_{ij}}{\sigma(W)} - \frac{[\mathbf{u}_1 \mathbf{v}_1^T]_{ij}}{\sigma^2(W)} W \quad (4)$$

$$= \frac{1}{\sigma(W)} (E_{ij} - [\mathbf{u}_1 \mathbf{v}_1^T]_{ij} \bar{W}_{\text{SN}}), \quad (5)$$

where E_{ij} is the matrix whose (i, j) -th entry is 1 and zero otherwise, and \mathbf{u}_1 and \mathbf{v}_1 are respectively the first left and right singular vectors of W .

Gradient Analysis

If \mathbf{h} is the hidden layer to be transformed by \bar{W}_{SN} , the derivative of the $V(G, D)$ calculated over the mini-batch w.r.t. W of the discriminator D is given by

$$\frac{\partial V(G, D)}{\partial W} = \frac{1}{\sigma(W)} (\hat{E}[\boldsymbol{\delta} \mathbf{h}^T] - (\hat{E}[\boldsymbol{\delta}^T \bar{W}_{\text{SN}} \mathbf{h}]) \mathbf{u}_1 \mathbf{v}_1^T) \quad (6)$$

$$= \frac{1}{\sigma(W)} (\hat{E}[\boldsymbol{\delta} \mathbf{h}^T] - \lambda \mathbf{u}_1 \mathbf{v}_1^T) \quad (7)$$

- $\boldsymbol{\delta} := (\partial V(G, D) / \partial (\bar{W}_{\text{SN}} \mathbf{h}))^T$
- $\lambda := \hat{E}[\boldsymbol{\delta}^T (\bar{W}_{\text{SN}} \mathbf{h})]$
- $\hat{E}[\cdot]$: empirical expectation over the mini-batch

Gradient Analysis

$$\frac{\partial V(G, D)}{\partial W} = \frac{1}{\sigma(W)} (\hat{E}[\delta \mathbf{h}^T] - (\hat{E}[\delta^T \bar{W}_{\text{SN}} \mathbf{h}]) \mathbf{u}_1 \mathbf{v}_1^T) \quad (6)$$

$$= \frac{1}{\sigma(W)} (\hat{E}[\delta \mathbf{h}^T] - \lambda \mathbf{u}_1 \mathbf{v}_1^T) \quad (7)$$

- The first term $\hat{E}[\delta \mathbf{h}^T]$ of (7) is same as the derivative w/o spectral normalization
- The second term can be seen as a regularizer with adaptive coefficient λ
- $\lambda > 0$ if δ and $\bar{W}_{\text{SN}} \mathbf{h}$ point to same direction, thus gives a penalty

Gradient Analysis

$$\frac{\partial V(G, D)}{\partial W} = \frac{1}{\sigma(W)} (\hat{E}[\boldsymbol{\delta} \mathbf{h}^T] - (\hat{E}[\boldsymbol{\delta}^T \bar{W}_{\text{SN}} \mathbf{h}]) \mathbf{u}_1 \mathbf{v}_1^T) \quad (6)$$

$$= \frac{1}{\sigma(W)} (\hat{E}[\boldsymbol{\delta} \mathbf{h}^T] - \lambda \mathbf{u}_1 \mathbf{v}_1^T) \quad (7)$$

Spectral normalization prevents the transformation of each layer from being too sensitive in one direction.

Experiments

Metrics

Assessment metrics:

- Inception Score (IS) [2]: KL divergence between conditional and marginal class distribution, higher is better
- Fréchet Inception Distance (FID) [3]: Wassertein-2 distance between generated and real images, lower is better

Experiments

The spectral norm of each layer in the discriminator doesn't change much in the course of training, they tends to be stabilized.

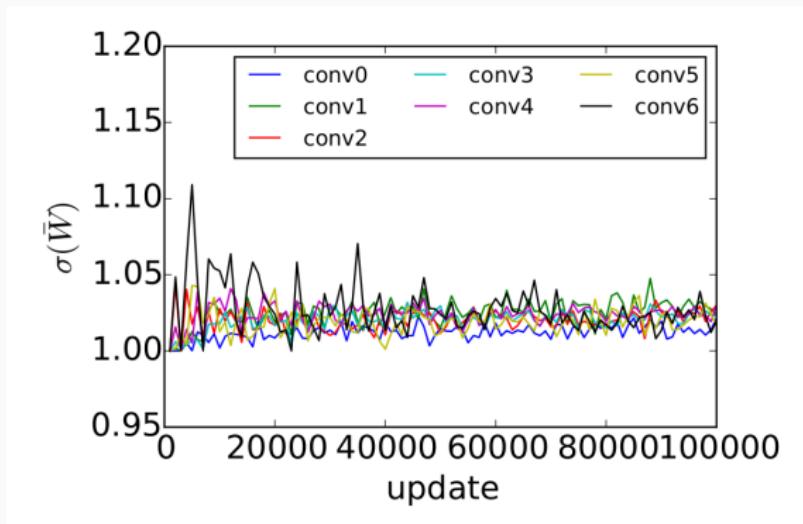


Figure 2: Spectral norm changes during training

Experiments

Spectral normalization is more robust than others w.r.t. different settings of hyperparameters.

Setting	α	β_1	β_2	n_{dis}
A [†]	0.0001	0.5	0.9	5
B [‡]	0.0001	0.5	0.999	1
C [*]	0.0002	0.5	0.999	1
D	0.001	0.5	0.9	5
E	0.001	0.5	0.999	5
F	0.001	0.9	0.999	5

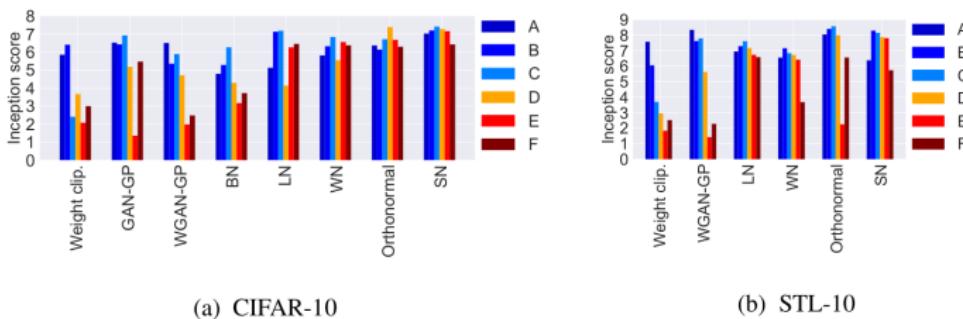


Figure 3: Inception scores on CIFAR-10 and STL-10

Spectral normalized GAN (SN-GAN) generates better images.
(Hope you can see it.)

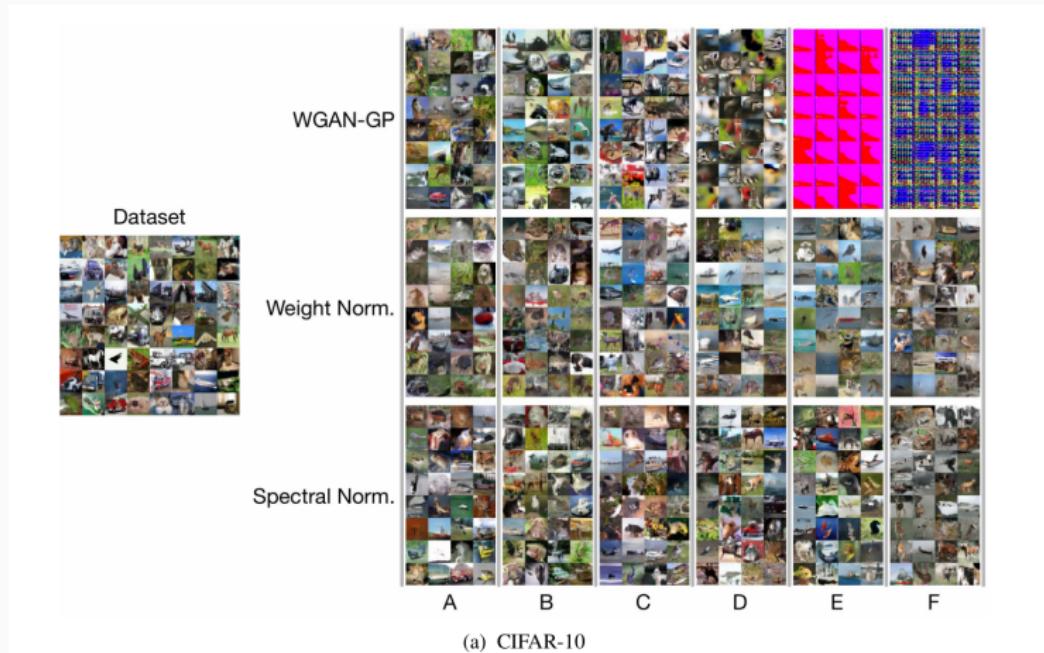


Figure 4: Generated images on different methods on CIFAR-10

ILSVRC2012 dataset: large high dimensional.

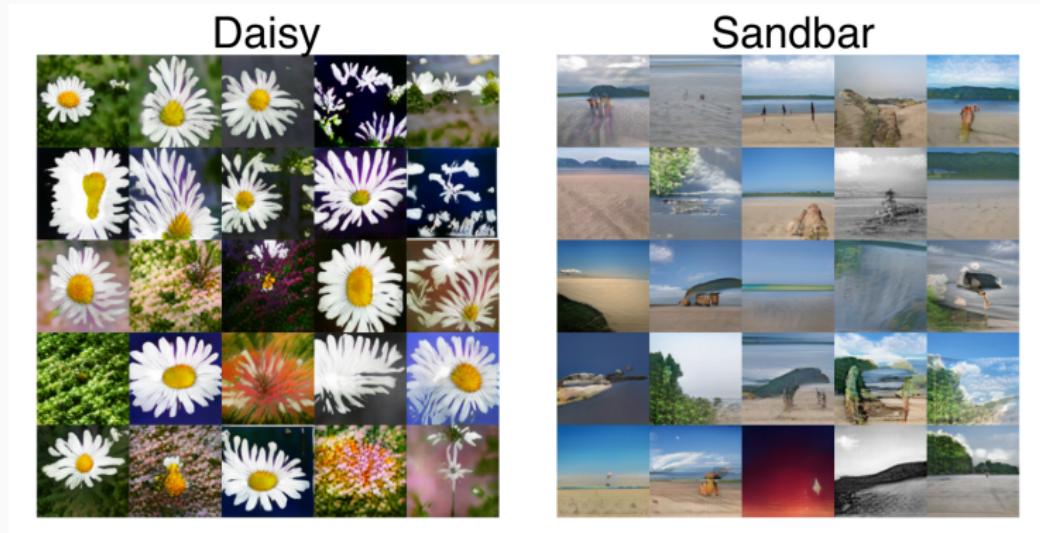


Figure 5: 128x128 pixel images generated by SN-GANs trained on ILSVRC2012 dataset

Singular Value Visualization

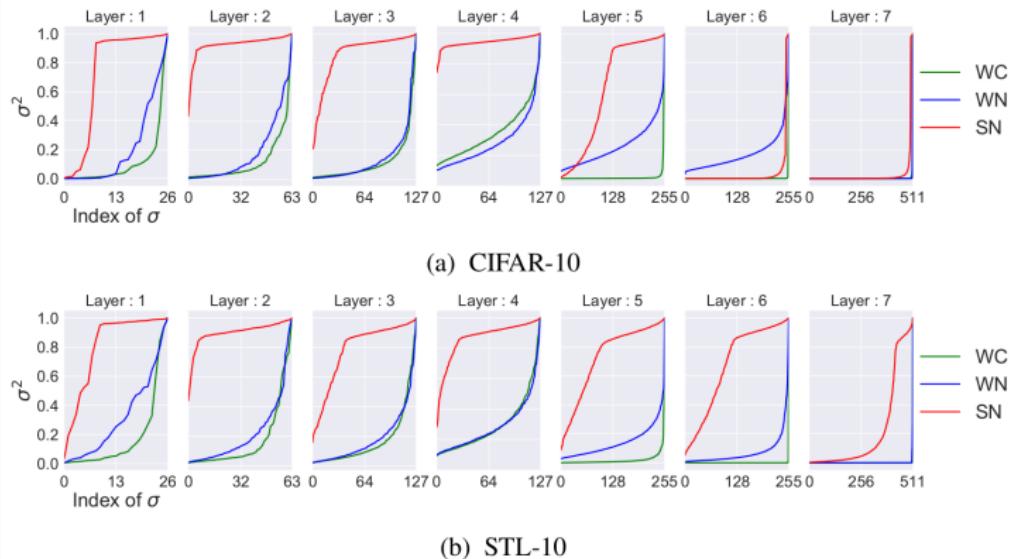


Figure 6: Squared singular value of weight matrix trained with different methods: Weight Clipping (WC), Weight Normalization (WN), Spectral Normalization (SN)

Singular Value Visualization

- recall that SN can prevent the column space of weight matrix from being concentrating on one direction
- thus the singular values is more broadly distributed
- thus SN-GAN can generate more diversity, relatively

Training Time

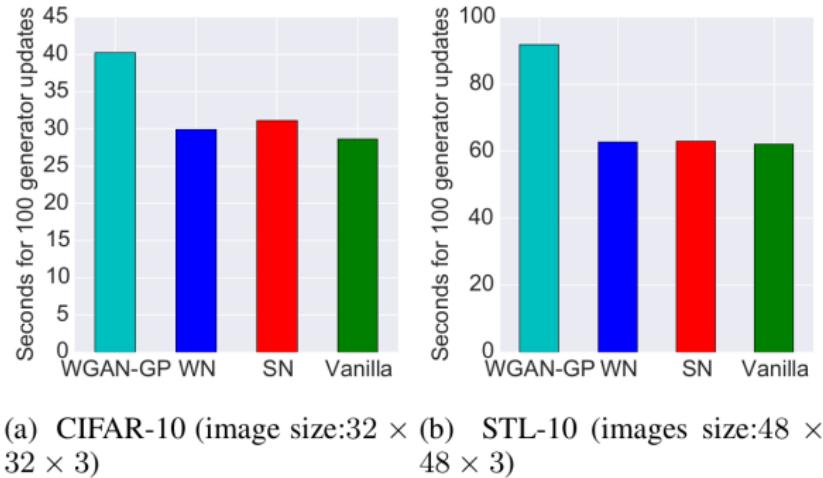


Figure 6: Computational time for 100 updates

Conclusion

Conclusion

- Proposed spectral normalization as a stabilizer
- SN-GANs generate more diversity and have better performance
- Spectral normalization is computational light

References

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- [2] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, “Improved techniques for training gans,” in *Advances in neural information processing systems*, pp. 2234–2242, 2016.
- [3] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, G. Klambauer, and S. Hochreiter, “Gans trained by a two time-scale update rule converge to a nash equilibrium,” *arXiv preprint arXiv:1706.08500*, vol. 12, no. 1, 2017.

Thank you!