



On Computation and Generalization of Generative Adversarial Networks under Spectrum Control

Haoming Jiang*, Zhehui Chen*, Minshuo Chen*, Feng Liu[△], Dingding Wang[△], Tuo Zhao*

*Georgia Tech [△]Florida Atlantic University



Background

GANs solve:

$$\min_{\theta} \max_{\mathcal{W}} f(\theta, \mathcal{W}) := \frac{1}{n} \sum_{i=1}^n \phi(\mathcal{A}(D_{\mathcal{W}}(x_i))) + \mathbb{E}_{x \sim \mathcal{D}_G} [\phi(1 - \mathcal{A}(D_{\mathcal{W}}(x)))]$$

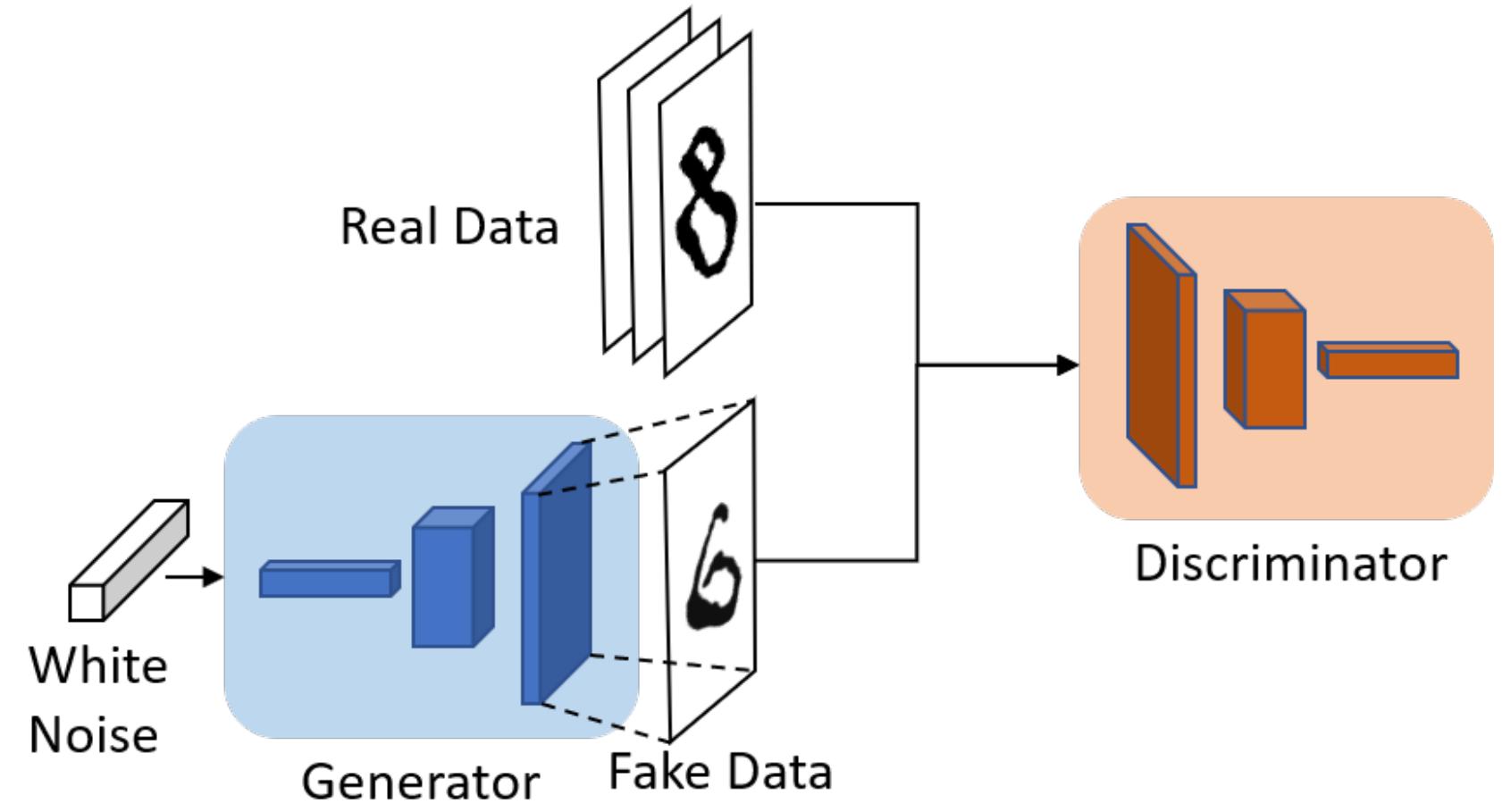
- G_{θ} denotes the **generator**.

- $D_{\mathcal{W}}$ denotes the **discriminator**.

An L -layer discriminator can be formulated as follows:

$$D_{\mathcal{W}}(x) = W_L \sigma_{L-1}(W_{L-1} \cdots \sigma_1(W_1 x) \cdots)$$

- For DC-GAN, $\phi(x) = \log(x)$, $\mathcal{A}(x) = \frac{1}{1+\exp(-x)}$.



SN-GAN [1] stabilizes training by spectral normalization on W_i 's: $\tilde{W}_i = \frac{W_i}{\|W_i\|_2}$. $\|W_i\|_2$ is obtained by power iteration.

Questions:

- Why does spectral normalization help?
- Can we make any further improvement?

Generalization Bound

Theorem 1. Under some technical assumptions and assume

- $\|W_i\|_2 \leq B_{W_i}$ for $i \in [L]$ and $\|x_k\|_2 \leq B_x$ for $k \in [n]$.

- Generator and discriminator are well trained, i.e.,

$$d_{\mathcal{F},\phi}(\hat{\mu}_n, \nu_n) - \inf_{\nu \in \mathcal{D}_G} d_{\mathcal{F},\phi}(\hat{\mu}_n, \nu) \leq \epsilon,$$

where $d_{\mathcal{F},\phi}(\cdot, \cdot)$ is the **neural distance** (similar to \mathcal{F} -divergence with \mathcal{F} being the class of discriminators).

\Rightarrow with probability at least $1 - \delta$, we have

$$d_{\mathcal{F},\phi}(\mu, \nu_n) - \inf_{\nu \in \mathcal{D}_G} d_{\mathcal{F},\phi}(\mu, \nu) \leq \tilde{O}\left(\frac{B_x \prod_{i=1}^L B_{W_i} \sqrt{d^2 L}}{\sqrt{n}}\right).$$

Spectral normalization \Rightarrow polynomial bound $\tilde{O}\left(\sqrt{\frac{d^2 L}{n}}\right)$.

Spectral normalization controls the product of spectral norms to benefit GANs in generalization.

SVD Reparameterization

$$\min_{\theta} \max_{\mathcal{W}} f(\theta, \mathcal{W})$$

\Downarrow SVD reparameterization: $W_i = U_i E_i V_i^T$, where $E_i = \text{diag}(e_1^i, \dots, e_{r_i}^i)$ with $e_1^i \geq \dots \geq e_{r_i}^i \geq 0$. Denote $\mathcal{E} = \{E_i\}_{i=1}^L, \mathcal{U} = \{U_i\}_{i=1}^L, \mathcal{V} = \{V_i\}_{i=1}^L$.

$$\min_{\theta} \max_{\mathcal{E}, \mathcal{U}, \mathcal{V}} f(\theta, \mathcal{E}, \mathcal{U}, \mathcal{V}),$$

subject to $U_i^T U_i = I_i, V_i^T V_i = I_i$.

\Downarrow Control spectrum via regularizer $\mathcal{R}(\mathcal{E})$, and constraint $\mathcal{E} \in \Omega$.

$$\min_{\theta} \max_{\mathcal{E}, \mathcal{U}, \mathcal{V}} f(\theta, \mathcal{E}, \mathcal{U}, \mathcal{V}) - \gamma \mathcal{R}(\mathcal{E}),$$

subject to $\mathcal{E} \in \Omega, U_i^T U_i = I_i, V_i^T V_i = I_i$.

\Downarrow Relax orthogonal constraint for computational efficiency.

$$\min_{\theta} \max_{\mathcal{E}, \mathcal{U}, \mathcal{V}} f(\theta, \mathcal{E}, \mathcal{U}, \mathcal{V}) - \lambda \mathcal{L}_o(\mathcal{U}, \mathcal{V}) - \gamma \mathcal{R}(\mathcal{E}), \text{ s.t. } \mathcal{E} \in \Omega,$$

$$\text{where } \mathcal{L}_o(\mathcal{U}, \mathcal{V}) = \sum_{i=1}^L (\|U_i^T U_i - I_i\|_F^2 + \|V_i^T V_i - I_i\|_F^2).$$

Spectrum Control

Layer-wise Control:

$$\Omega = \{\mathcal{E} : 0 \leq e_1^i \leq 1 \forall i \in [L]\}.$$

- Spectral Constraint: $\tilde{e}_j^i = \min(1, e_j^i)$.
- Spectrum Normalization: $\tilde{e}_j^i = e_j^i / \max_j e_j^i$.

Overall Control: Alternatively, we can control the overall Lipschitz constant, $\prod_{i=1}^L e_1^i$, by Lipschitz regularizer:

$$\mathcal{R}(\mathcal{E}) := \max \left(\log \prod_{i=1}^L e_1^i, 0 \right) = \max \left(\sum_{i=1}^L \log e_1^i, 0 \right).$$

D-Optimal Regularizer: Penalty focusing on small singular values:

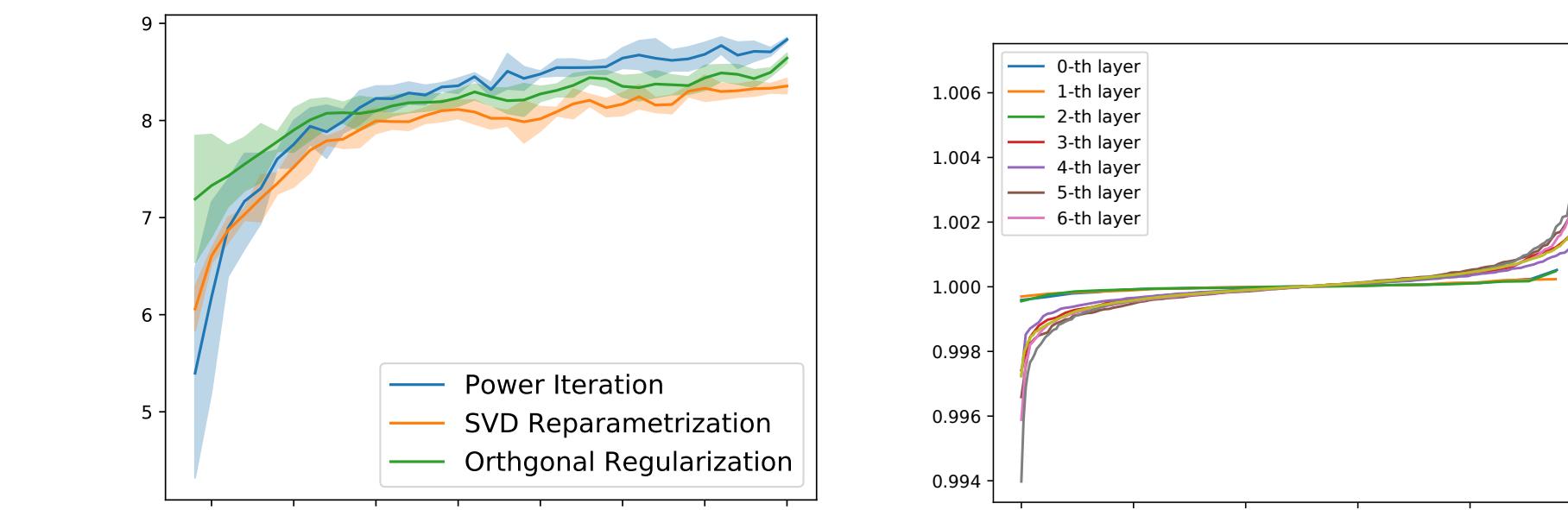
$$\mathcal{R}(\mathcal{E}) = \frac{1}{2} \sum_{i=1}^{L-1} \log \left(|(E_i^T E_i)^{-1}| \right) = - \sum_{i=1}^{L-1} \log \left(\prod_{k=1}^{r_i} e_k^i \right).$$

Divergence Regularizer: Empirical KL divergence to a half-normal distribution p :

$$\mathcal{R}(\mathcal{E}) = \sum_{i=1}^{L-1} \frac{1}{r_i - 1} \sum_{k=1}^{r_i-1} \log \left[\frac{(r_i - 1)^{-1}}{(e_{k+1}^i - e_k^i)p(e_k^i)} \right].$$

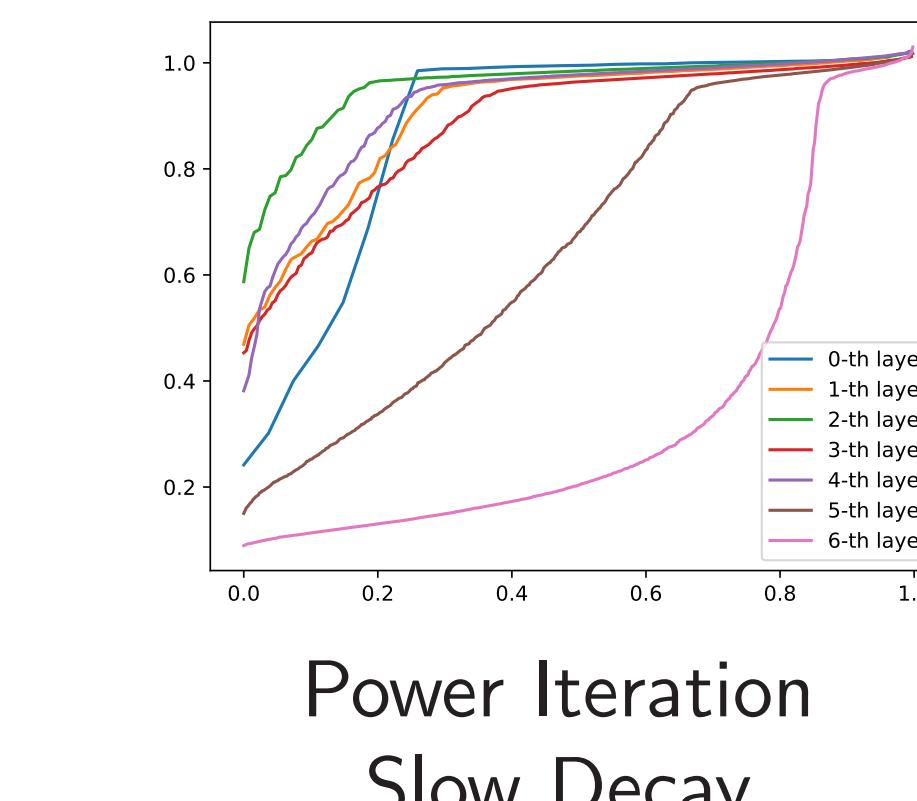
Motivation – Slow Singular Value Decay

Comparison between spectral normalization (Power iteration, SVD reparameterization) and orthogonal regularizer:

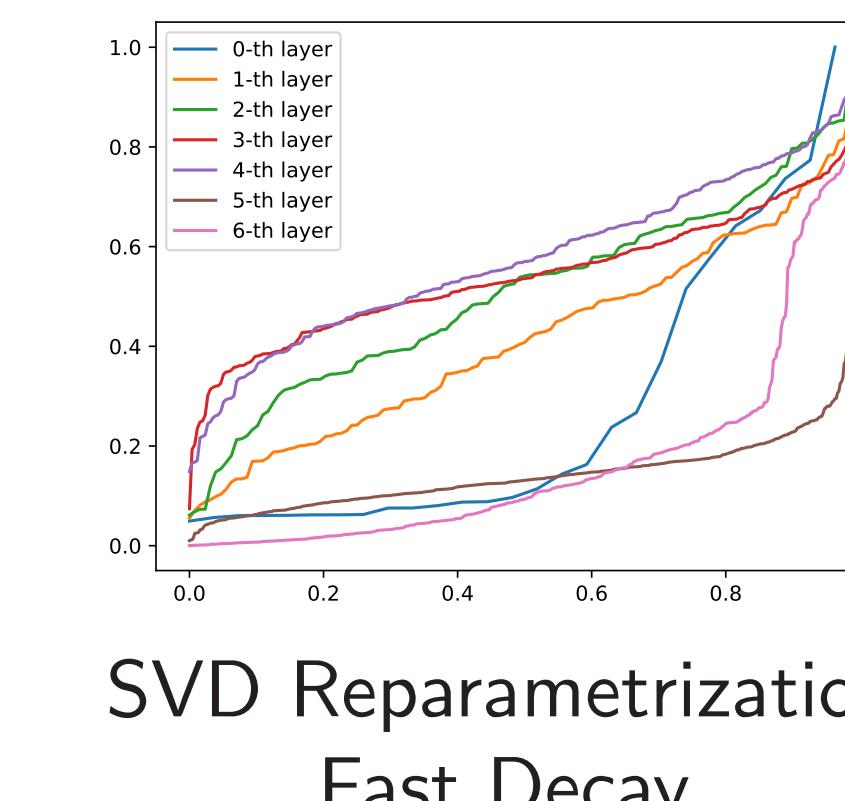


Performance

Orthogonal Regularizer
No Decay



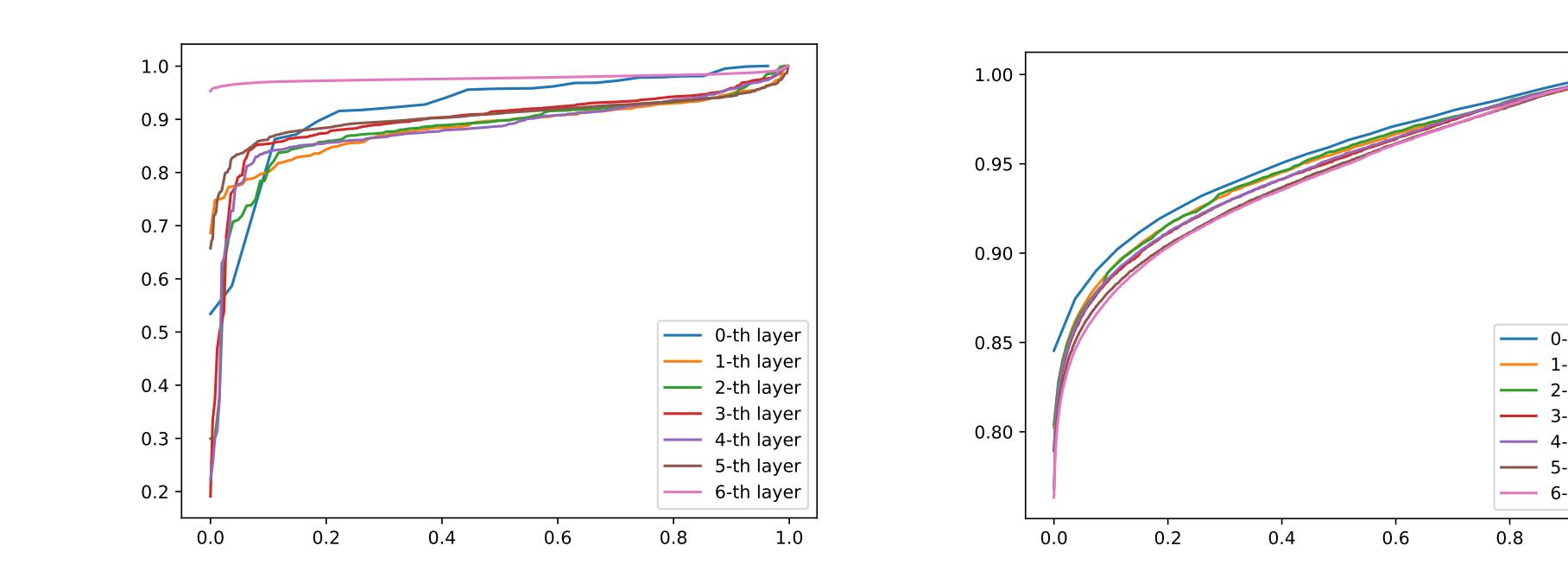
Power Iteration
Slow Decay



SVD Reparameterization
Fast Decay

Observation: Slow singular value decay is better than both no decay and fast decay.

Both D-Optimal and divergence regularizers yield slow singular value decay.

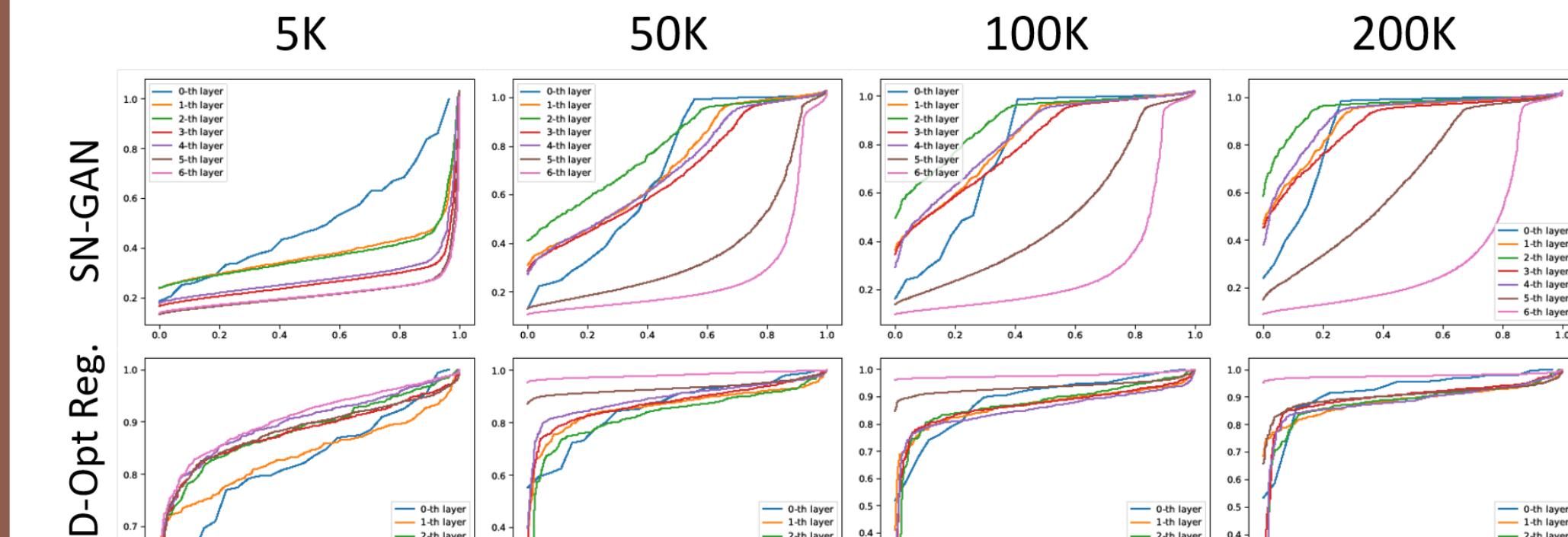


SN+D-Optimal Regularizer

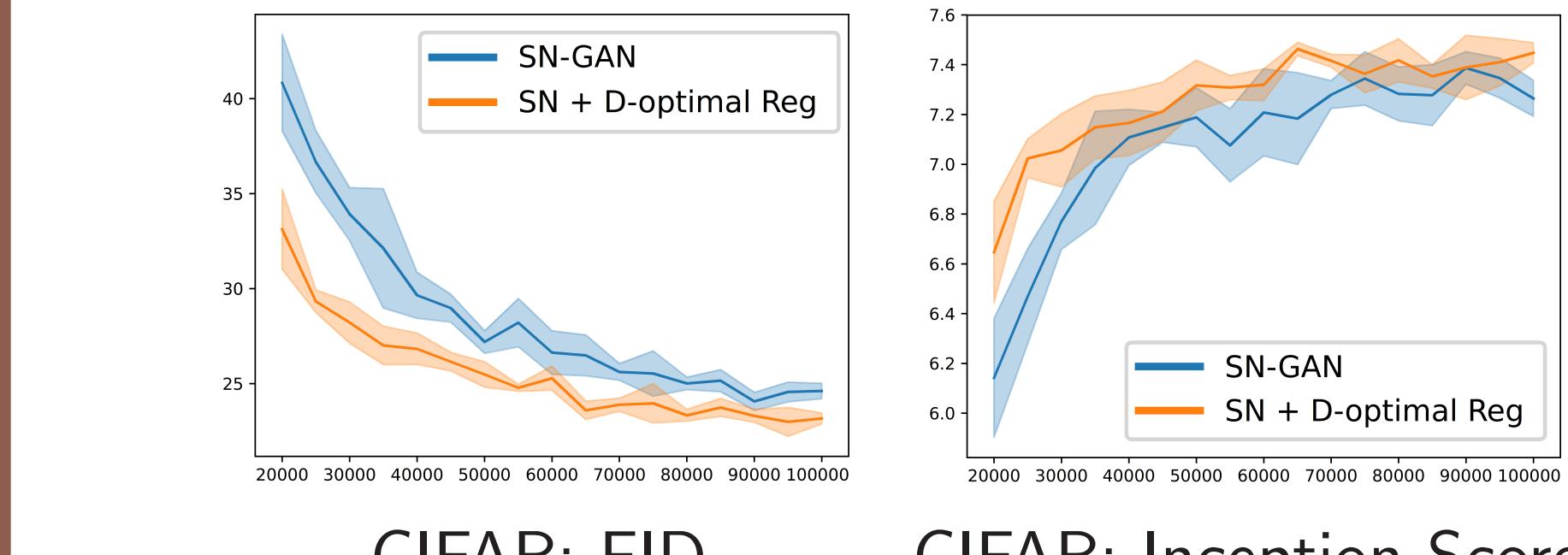
Divergence Regularizer

Image Generation

Singular Value Decay:



DC-GAN on CIFAR-10 and STL-10:



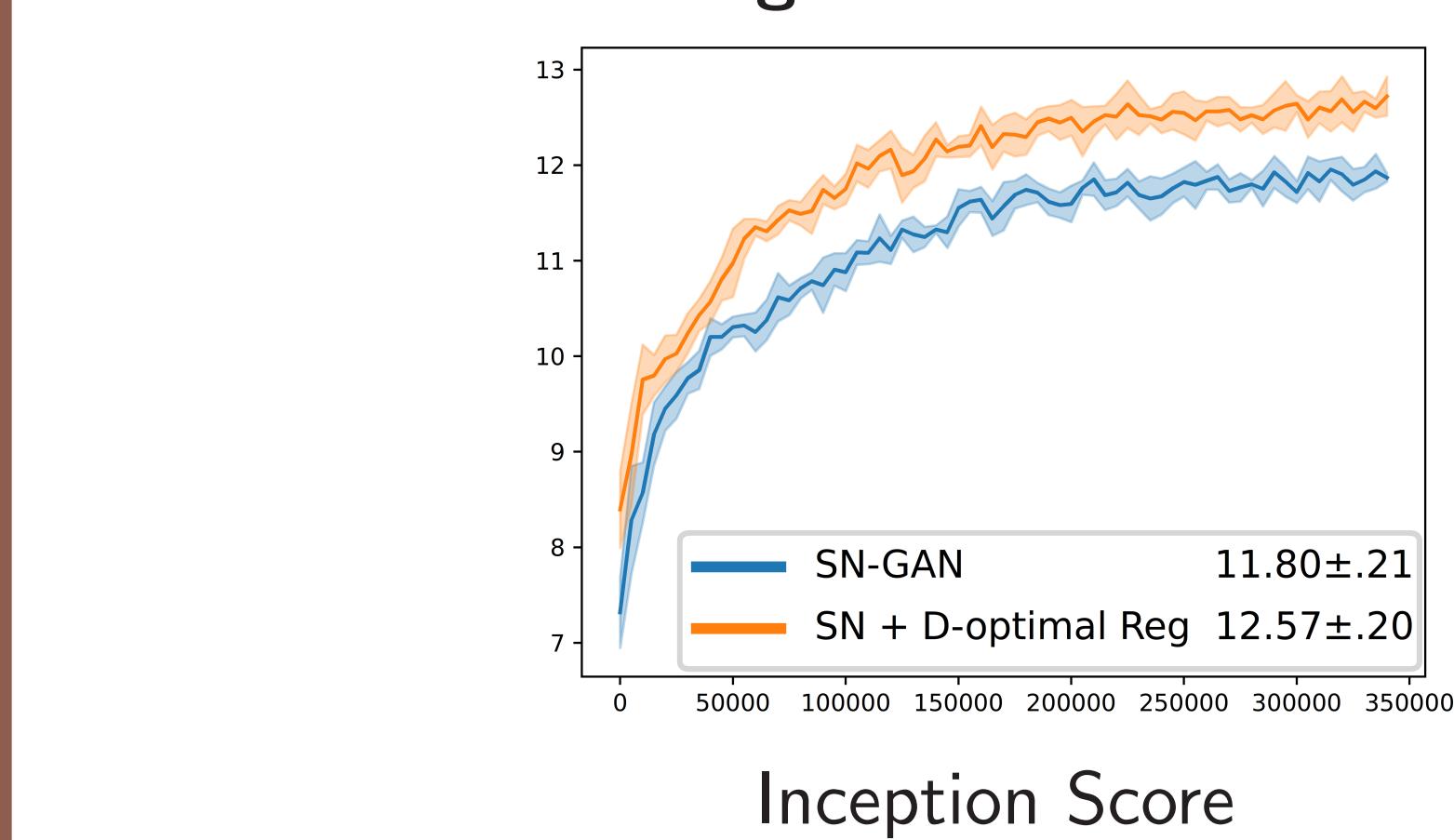
CIFAR: FID

CIFAR: Inception Score

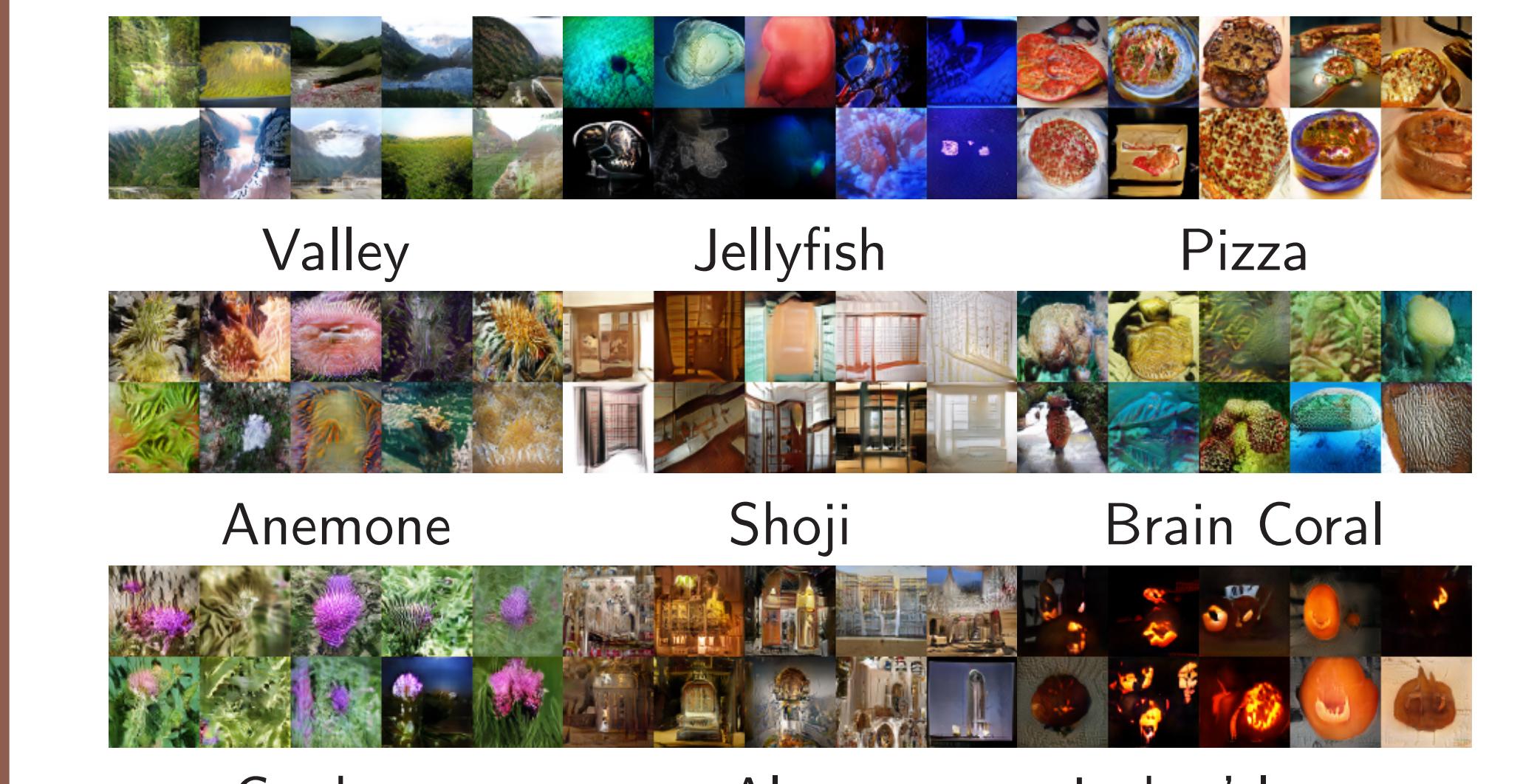
STL: FID

STL: Inception Score

ResNet-GAN on ImageNet:



Inception Score



Reference

- [1] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.