



Australian  
National  
University



# CHAIN: Enhancing Generalization in Data-Efficient GANs via *lipsCHitz* continuity constrAIned Normalization

Yao Ni<sup>†</sup>    Piotr Koniusz<sup>§,†</sup>

<sup>†</sup>The Australian National University

<sup>§</sup>Data61♥CSIRO

CVPR 2024

# Background: Challenges in Data-Efficient GANs and BN

## Challenges in Data-Efficient GANs:

- Discriminator overfitting.
- Training instability.

# Background: Challenges in Data-Efficient GANs and BN

## Challenges in Data-Efficient GANs:

- Discriminator overfitting.
- Training instability.

## Advantages of Batch Normalization (BN):

- **aids generalization** by aligning training and test distributions and reducing sharpness of loss landscape.
- **stabilizes training** process by mitigating internal covariate shift.

# Background: Challenges in Data-Efficient GANs and BN

## Challenges in Data-Efficient GANs:

- Discriminator overfitting.
- Training instability.

## Advantages of Batch Normalization (BN):

- **aids generalization** by aligning training and test distributions and reducing sharpness of loss landscape.
- **stabilizes training** process by mitigating internal covariate shift.

BN appears to benefit the discriminator, **but is rarely used as it impairs performance.**

# Background: Challenges in Data-Efficient GANs and BN

## Challenges in Data-Efficient GANs:

- Discriminator overfitting.
- Training instability.

## Advantages of Batch Normalization (BN):

- **aids generalization** by aligning training and test distributions and reducing sharpness of loss landscape.
- **stabilizes training** process by mitigating internal covariate shift.

BN appears to benefit the discriminator, **but is rarely used as it impairs performance.**

**Goal: Integrate BN into discriminator for improved generalization.**

# Methods: Generalization error of GANs

Lemma 3.1 (GAN generalization error on function set):

$$\epsilon_{\text{gan}} \leq 2d_{\mathcal{H}}(\mu, \hat{\mu}_n) + 2d_{\mathcal{H}}(\nu_n^*, \hat{\nu}_n)$$

$d_{\mathcal{H}}$ : discrepancy over discriminator set.  $\mu, \hat{\mu}_n$ : unseen/seen real data.  $\nu_n^*, \hat{\nu}_n$ : ideal/seen fake.

# Methods: Generalization error of GANs

Lemma 3.1 (GAN generalization error on function set):

$$\epsilon_{\text{gan}} \leq 2d_{\mathcal{H}}(\mu, \hat{\mu}_n) + 2d_{\mathcal{H}}(\nu_n^*, \hat{\nu}_n)$$

$d_{\mathcal{H}}$ : discrepancy over discriminator set.  $\mu, \hat{\mu}_n$ : unseen/seen real data.  $\nu_n^*, \hat{\nu}_n$ : ideal/seen fake.

$\nu_n^* \approx \hat{\mu}_n \implies$  Lowering real/fake discrepancy aids generalization

# Methods: Generalization error of GANs

Lemma 3.1 (GAN generalization error on function set):

$$\epsilon_{\text{gan}} \leq 2d_{\mathcal{H}}(\mu, \hat{\mu}_n) + 2d_{\mathcal{H}}(\nu_n^*, \hat{\nu}_n)$$

$d_{\mathcal{H}}$ : discrepancy over discriminator set.  $\mu, \hat{\mu}_n$ : unseen/seen real data.  $\nu_n^*, \hat{\nu}_n$ : ideal/seen fake.

$\nu_n^* \approx \hat{\mu}_n \implies$  Lowering real/fake discrepancy aids generalization

$\mu$  is inaccessible. We need further analyze  $d_{\mathcal{H}}(\mu, \hat{\mu}_n)$ .



# Methods: Generalization error of GANs

Lemma 3.2 (GAN generalization error on neural network):

$$\epsilon_{\text{gan}}^{\text{nn}} \leq 2\omega \left( \|\nabla_{\theta_d}\|_2 + \|\tilde{\nabla}_{\theta_d}\|_2 \right) + 4R \left( \frac{\|\theta_d\|_2^2}{\omega^2}, \frac{1}{n} \right) + \omega^2 (|\lambda_{\max}^{\mathbf{H}}| + |\lambda_{\max}^{\tilde{\mathbf{H}}}|)$$

$\epsilon_{\text{gan}}^{\text{nn}}$ :  $\epsilon_{\text{gan}}$  on neural networks,  $\omega > 0$ .  $\theta_d$ :  $D$ 's weights.  $\nabla_{\theta_d}$ ,  $\lambda_{\max}^{\mathbf{H}}$ : real gradient, top Hessian eigenvalue.  $\tilde{\nabla}_{\theta_d}$ ,  $\lambda_{\max}^{\tilde{\mathbf{H}}}$ : fake versions.  $R$ : related to  $\|\theta_d\|_2^2$  and data size  $n$ .

# Methods: Generalization error of GANs

Lemma 3.2 (GAN generalization error on neural network):

$$\epsilon_{\text{gan}}^{\text{nn}} \leq 2\omega \left( \|\nabla_{\theta_d}\|_2 + \|\tilde{\nabla}_{\theta_d}\|_2 \right) + 4R \left( \frac{\|\theta_d\|_2^2}{\omega^2}, \frac{1}{n} \right) + \omega^2 (|\lambda_{\text{max}}^H| + |\lambda_{\text{max}}^{\tilde{H}}|)$$

$\epsilon_{\text{gan}}^{\text{nn}}$ :  $\epsilon_{\text{gan}}$  on neural networks,  $\omega > 0$ .  $\theta_d$ :  $D$ 's weights.  $\nabla_{\theta_d}, \lambda_{\text{max}}^H$ : real gradient, top Hessian eigenvalue.  $\tilde{\nabla}_{\theta_d}, \lambda_{\text{max}}^{\tilde{H}}$ : fake versions.  $R$ : related to  $\|\theta_d\|_2^2$  and data size  $n$ .

Reducing weight gradient norms of discriminator aids generalization.

# Methods: Motivation for BN

Better generalization on GANs requires:

- Lowering real-fake discrepancy
- Reducing weight gradient norms of discriminator

# Methods: Motivation for BN

Better generalization on GANs requires:

- Lowering real-fake discrepancy
- Reducing weight gradient norms of discriminator

Motivation of BN



# Methods: Motivation for BN

Better generalization on GANs requires:

- Lowering real-fake discrepancy
- Reducing weight gradient norms of discriminator

Motivation of BN



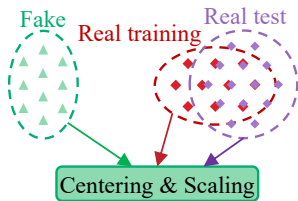
Applying BN **separately** on real/fake batches

# Methods: Motivation for BN

Better generalization on GANs requires:

- Lowering real-fake discrepancy
- Reducing weight gradient norms of discriminator

## Motivation of BN



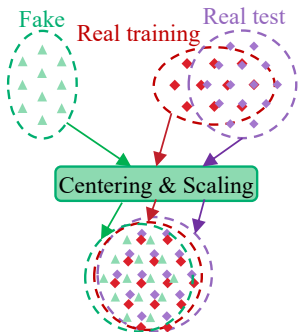
Applying BN **separately** on real/fake batches

# Methods: Motivation for BN

Better generalization on GANs requires:

- Lowering real-fake discrepancy
- Reducing weight gradient norms of discriminator

## Motivation of BN



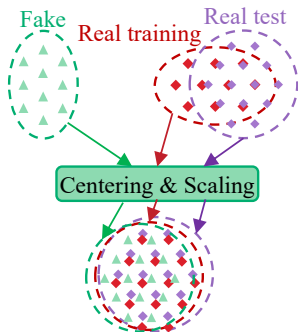
Applying BN **separately** on real/fake batches reduces the fake-real discrepancy via standardization.

# Methods: Motivation for BN

Better generalization on GANs requires:

- Lowering real-fake discrepancy
- Reducing weight gradient norms of discriminator

## Motivation of BN



Applying BN **separately** on real/fake batches reduces the fake-real discrepancy via standardization.

**But incorporating BN risks gradient explosion issues.**



# Methods: Gradient Issues of BN

Standardization in BN:

Linear transformation:  $\mathbf{Y} = \mathbf{A}\mathbf{W}$

Centering:  $\mathbf{\hat{Y}}^c = \mathbf{Y} - \mu$

Scaling:  $\mathbf{\hat{Y}}^s = \mathbf{\hat{Y}}^c / \sigma$ .

# Methods: Gradient Issues of BN

Linear transformation:  $\mathbf{Y} = \mathbf{A}\mathbf{W}$

Standardization in BN:

Centering:  $\mathbf{\hat{Y}}^c = \mathbf{Y} - \mu$

Scaling:  $\mathbf{\hat{Y}}^s = \mathbf{\hat{Y}}^c / \sigma$ .

Theorem 3.1 (Issue in centering, **similarity dropping causes feature divergence**):

$$\mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [\cos(\mathbf{y}_1, \mathbf{y}_2)] \geq \mathbb{E}_{\mathbf{\hat{y}}_1^c, \mathbf{\hat{y}}_2^c} [\cos(\mathbf{\hat{y}}_1^c, \mathbf{\hat{y}}_2^c)] = 0$$

$\mathbf{y}_1, \mathbf{\hat{y}}_1^c$ : pre- & post-centering features. Features similar in early layers diverge in later layers.

# Methods: Gradient Issues of BN

Linear transformation:  $\mathbf{Y} = \mathbf{A}\mathbf{W}$

Standardization in BN:

Centering:  $\mathbf{\hat{Y}}^c = \mathbf{Y} - \boldsymbol{\mu}$

Scaling:  $\mathbf{\hat{Y}}^s = \mathbf{\hat{Y}}^c / \boldsymbol{\sigma}$ .

Theorem 3.1 (Issue in centering, **similarity dropping causes feature divergence**):

$$\mathbb{E}_{\mathbf{y}_1, \mathbf{y}_2} [\cos(\mathbf{y}_1, \mathbf{y}_2)] \geq \mathbb{E}_{\mathbf{\hat{y}}_1^c, \mathbf{\hat{y}}_2^c} [\cos(\mathbf{\hat{y}}_1^c, \mathbf{\hat{y}}_2^c)] = 0$$

$\mathbf{y}_1, \mathbf{\hat{y}}_1^c$ : pre- & post-centering features. Features similar in early layers diverge in later layers.

Theorem 3.2 (Issue in scaling, **unbounded Lipschitz causes gradient explosion**):

$$\left\| \text{diag}\left(\frac{1}{\boldsymbol{\sigma}}\right) \right\|_{\text{lc}} = \frac{1}{\sigma_{\min}}$$

Lipschitz constant (lc) is large when  $\sigma_{\min} = \min_c \sigma_c$ , is small.

# Method: CHAIN replaces centering/scaling with OMR/ARMS

mean  $\mu$  and root mean square  $\psi$ :

$$\mu_c = \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}$$

$$\psi_c = \sqrt{\left( \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}^2 \right) + \epsilon}$$

# Method: CHAIN replaces centering/scaling with OMR/ARMS

mean  $\mu$  and root mean square  $\psi$ :

$$\mu_c = \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}$$

$$\psi_c = \sqrt{\left( \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}^2 \right) + \epsilon}$$

0-mean regularization:

$$\ell^{\text{OMR}}(\mathbf{Y}) = \lambda \cdot p \cdot \|\mu\|_2^2$$

$\psi_{\min} = \min_c \psi_c$ .  $\epsilon = 10^{-5}$ .  $\lambda$ : a hyperparameter.  $p \in [0, 1]$  controls  $\ell^{\text{OMR}}$  and Bernoulli mask  $\mathbf{M} \sim \mathcal{B}(p)$

---

Pytorch-style pseudo code for CHAIN<sub>batch</sub>

---

```
# Y: B x d x H x W size; lbd: hyperparameter  $\lambda$   
def CHAIN_batch(Y, p, lbd, eps=1e-5):  
    reg=Y.mean([0,2,3]).square().sum()*(p*lbd)
```

---

# Method: CHAIN replaces centering/scaling with OMR/ARMS

mean  $\mu$  and root mean square  $\psi$ :

$$\mu_c = \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}$$

$$\psi_c = \sqrt{\left( \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}^2 \right) + \epsilon}$$

0-mean regularization:

$$\ell^{\text{OMR}}(\mathbf{Y}) = \lambda \cdot p \cdot \|\mu\|_2^2$$

Adaptive RMS normalization:

$$\text{ARMS}(\mathbf{Y}) = (1 - \mathbf{M}) \odot \mathbf{Y} + \mathbf{M} \odot \frac{\mathbf{Y}}{\psi} \cdot \psi_{\min}$$

$\psi_{\min} = \min_c \psi_c$ .  $\epsilon = 10^{-5}$ .  $\lambda$ : a hyperparameter.  $p \in [0, 1]$  controls  $\ell^{\text{OMR}}$  and Bernoulli mask  $\mathbf{M} \sim \mathcal{B}(p)$

---

Pytorch-style pseudo code for CHAIN<sub>batch</sub>

---

```
# Y: B x d x H x W size; lbd: hyperparameter  $\lambda$ 
def CHAIN_batch(Y, p, lbd, eps=1e-5):
    reg=Y.mean([0,2,3]).square().sum()*(p*lbd)
    M= (torch.rand(*Y.shape[:2], 1, 1)<p)*1.0
    psi_s=Y.square().mean([0,2,3],keepdim=True)
    psi = (psi_s + eps).sqrt()
    psi_min = psi.min().detach()
    Y_arms = (1 - M) * Y + M * (Y/psi*psi_min)
    return Y_arms, reg
```

---

# Method: CHAIN reduces feature and weight gradients in $D$

$$\begin{aligned}\|\Delta \mathbf{y}_c\|_2^2 &\leq \|\Delta \dot{\mathbf{y}}_c\|_2^2 \left( \frac{(1-p)\psi_c + p\psi_{\min}}{\psi_c} \right)^2 - \frac{2(1-p)p\psi_{\min}}{B\psi_c} (\Delta \dot{\mathbf{y}}_c^T \tilde{\mathbf{y}}_c)^2 \\ \|\Delta \mathbf{w}_c\|_2^2 &\leq \lambda_{\max}^2 \|\Delta \mathbf{y}_c\|_2^2\end{aligned}$$

$\Delta \mathbf{y}_c, \Delta \dot{\mathbf{y}}_c$ :  $c$ -th column of gradient for CHAIN input/output  $\mathbf{Y}, \dot{\mathbf{Y}}$ .  $\lambda_{\max}$ : top eigenvalue of  $\mathbf{A}$ .  $\tilde{\mathbf{y}}_c$ :  $c$ -th column of  $\tilde{\mathbf{Y}} = \mathbf{Y}/\psi$ .  $\Delta \mathbf{w}_c$ :  $c$ -th column of grad for  $\mathbf{W}$ .  $p \in [0, 1]$

# Method: CHAIN reduces feature and weight gradients in $D$

$$\begin{aligned}\|\Delta \mathbf{y}_c\|_2^2 &\leq \|\Delta \dot{\mathbf{y}}_c\|_2^2 \left( \frac{(1-p)\psi_c + p\psi_{\min}}{\psi_c} \right)^2 - \frac{2(1-p)p\psi_{\min}}{B\psi_c} (\Delta \dot{\mathbf{y}}_c^T \tilde{\mathbf{y}}_c)^2 \\ \|\Delta \mathbf{w}_c\|_2^2 &\leq \lambda_{\max}^2 \|\Delta \mathbf{y}_c\|_2^2\end{aligned}$$

$\Delta \mathbf{y}_c, \Delta \dot{\mathbf{y}}_c$ :  $c$ -th column of gradient for CHAIN input/output  $\mathbf{Y}, \dot{\mathbf{Y}}$ .  $\lambda_{\max}$ : top eigenvalue of  $\mathbf{A}$ .  $\tilde{\mathbf{y}}_c$ :  $c$ -th column of  $\tilde{\mathbf{Y}} = \mathbf{Y}/\psi$ .  $\Delta \mathbf{w}_c$ :  $c$ -th column of grad for  $\mathbf{W}$ .  $p \in [0, 1]$

$$\frac{(1-p)\psi_c + p\psi_{\min}}{\psi_c} \leq 1$$



# Method: CHAIN reduces feature and weight gradients in $D$

$$\begin{aligned}\|\Delta \mathbf{y}_c\|_2^2 &\leq \|\Delta \dot{\mathbf{y}}_c\|_2^2 \left( \frac{(1-p)\psi_c + p\psi_{\min}}{\psi_c} \right)^2 - \frac{2(1-p)p\psi_{\min}}{B\psi_c} (\Delta \dot{\mathbf{y}}_c^T \tilde{\mathbf{y}}_c)^2 \\ \|\Delta \mathbf{w}_c\|_2^2 &\leq \lambda_{\max}^2 \|\Delta \mathbf{y}_c\|_2^2\end{aligned}$$

$\Delta \mathbf{y}_c, \Delta \dot{\mathbf{y}}_c$ :  $c$ -th column of gradient for CHAIN input/output  $\mathbf{Y}, \dot{\mathbf{Y}}$ .  $\lambda_{\max}$ : top eigenvalue of  $\mathbf{A}$ .  $\tilde{\mathbf{y}}_c$ :  $c$ -th column of  $\tilde{\mathbf{Y}} = \mathbf{Y}/\psi$ .  $\Delta \mathbf{w}_c$ :  $c$ -th column of grad for  $\mathbf{W}$ .  $p \in [0, 1]$

$$\frac{(1-p)\psi_c + p\psi_{\min}}{\psi_c} \leq 1 \quad (\Delta \dot{\mathbf{y}}_c^T \tilde{\mathbf{y}}_c)^2 \geq 0$$


# Method: CHAIN reduces feature and weight gradients in $D$

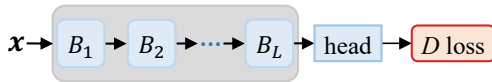
$$\begin{aligned}\|\Delta \mathbf{y}_c\|_2^2 &\leq \|\Delta \dot{\mathbf{y}}_c\|_2^2 \left( \frac{(1-p)\psi_c + p\psi_{\min}}{\psi_c} \right)^2 - \frac{2(1-p)p\psi_{\min}}{B\psi_c} (\Delta \dot{\mathbf{y}}_c^T \tilde{\mathbf{y}}_c)^2 \\ \|\Delta \mathbf{w}_c\|_2^2 &\leq \lambda_{\max}^2 \|\Delta \mathbf{y}_c\|_2^2\end{aligned}$$

$\Delta \mathbf{y}_c, \Delta \dot{\mathbf{y}}_c$ :  $c$ -th column of gradient for CHAIN input/output  $\mathbf{Y}, \dot{\mathbf{Y}}$ .  $\lambda_{\max}$ : top eigenvalue of  $\mathbf{A}$ .  $\tilde{\mathbf{y}}_c$ :  $c$ -th column of  $\tilde{\mathbf{Y}} = \mathbf{Y}/\psi$ .  $\Delta \mathbf{w}_c$ :  $c$ -th column of grad for  $\mathbf{W}$ .  $p \in [0, 1]$

$$\frac{(1-p)\psi_c + p\psi_{\min}}{\psi_c} \leq 1 \quad (\Delta \dot{\mathbf{y}}_c^T \tilde{\mathbf{y}}_c)^2 \geq 0 \implies \text{CHAIN reduces } \|\Delta \mathbf{y}_c\|_2 \text{ and } \|\Delta \mathbf{w}_c\|_2.$$


# Pipeline

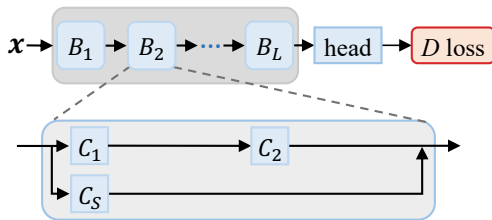
Discriminator with CHAIN 



$x$ : Real image.  $B_l$ :  $l$ -th block.  $D$ : Discriminator.

# Pipeline

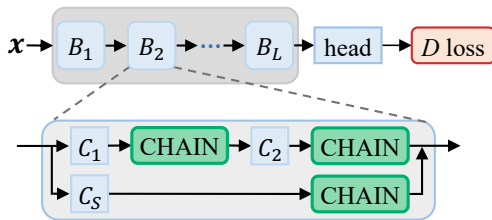
Discriminator with CHAIN 



$x$ : Real image.  $B_l$ :  $l$ -th block.  $D$ : Discriminator.  $C_S$ : Convolution in skip branch.

# Pipeline

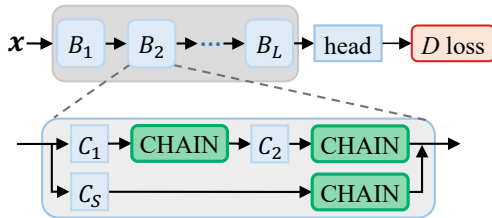
## Discriminator with CHAIN



$x$ : Real image.  $B_l$ :  $l$ -th block.  $D$ : Discriminator.  $C_S$ : Convolution in skip branch.

# Pipeline

Discriminator with CHAIN



CHAIN: 0MR & ARMS

0-Mean Regularization loss:

$$\ell^{0MR}(Y) = \lambda \cdot p \cdot \|\mu\|_2^2$$

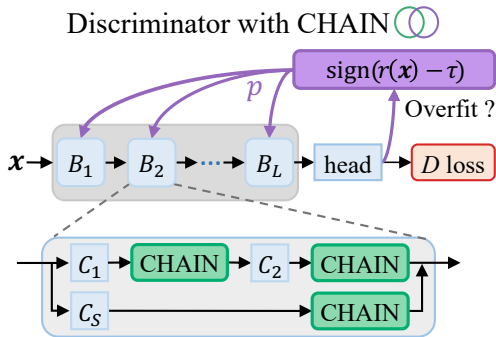
Adaptive Root Mean Square normalization:

$$ARMS(Y) = (1 - M) \odot Y + M \odot \frac{Y}{\psi} \cdot \psi_{\min}$$

$M \sim \mathcal{B}(p)$ .  $\mu, \psi$  are mean and root mean square of feature map  $Y$ .

$x$ : Real image.  $B_l$ :  $l$ -th block.  $D$ : Discriminator.  $C_S$ : Convolution in skip branch.  $\mathcal{B}$ : Bernoulli noise.  $p$ : Bernoulli probability and  $\ell^{0MR}$  strength.  $\lambda$ : A hyperparameter.

# Pipeline



## CHAIN: 0MR & ARMS

0-Mean Regularization loss:

$$\ell^{\text{0MR}}(\mathbf{Y}) = \lambda \cdot p \cdot \|\boldsymbol{\mu}\|_2^2$$

Adaptive Root Mean Square normalization:

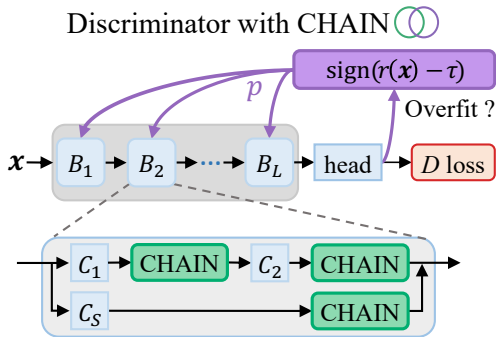
$$\text{ARMS}(\mathbf{Y}) = (1 - \mathbf{M}) \odot \mathbf{Y} + \mathbf{M} \odot \frac{\mathbf{Y}}{\boldsymbol{\psi}} \cdot \psi_{\min}$$

$\mathbf{M} \sim \mathcal{B}(p)$ .  $\boldsymbol{\mu}, \boldsymbol{\psi}$  are mean and root mean square of feature map  $\mathbf{Y}$ .

$\mathbf{x}$ : Real image.  $B_l$ :  $l$ -th block.  $D$ : Discriminator.  $C_S$ : Convolution in skip branch.  $\mathcal{B}$ : Bernoulli noise.  $p$ : Bernoulli probability and  $\ell^{\text{0MR}}$  strength.  $\lambda$ : A hyperparameter.  $\tau$ : A predefined threshold.  $\Delta_p$ : A small value.

Control  $p$ :  $r(\mathbf{x}) = \mathbb{E}[\text{sign}(D(\mathbf{x}))]$ ,  $p_{t+1} = p_t + \Delta_p \cdot \text{sign}(r(\mathbf{x}) - \tau)$ .

# Pipeline



## CHAIN: 0MR & ARMS

0-Mean Regularization loss:

$$\ell^{\text{0MR}}(\mathbf{Y}) = \lambda \cdot p \cdot \|\boldsymbol{\mu}\|_2^2$$

Adaptive Root Mean Square normalization:

$$\text{ARMS}(\mathbf{Y}) = (1 - \mathbf{M}) \odot \mathbf{Y} + \mathbf{M} \odot \frac{\mathbf{Y}}{\boldsymbol{\psi}} \cdot \psi_{\min}$$

$\mathbf{M} \sim \mathcal{B}(p)$ .  $\boldsymbol{\mu}, \boldsymbol{\psi}$  are mean and root mean square of feature map  $\mathbf{Y}$ .

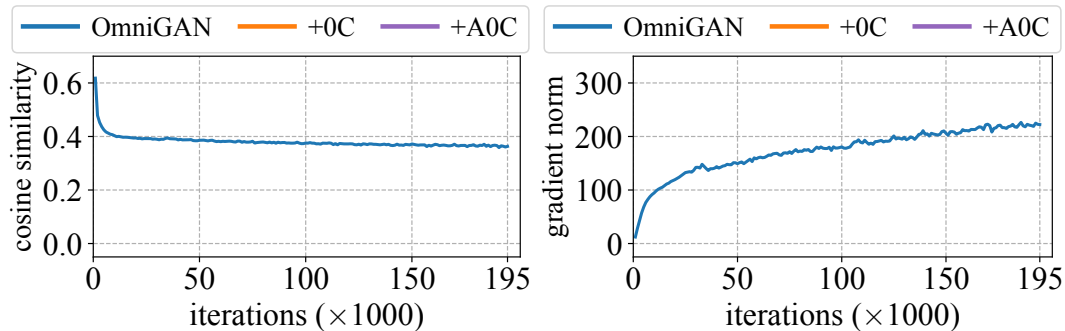
$\mathbf{x}$ : Real image.  $B_l$ :  $l$ -th block.  $D$ : Discriminator.  $C_S$ : Convolution in skip branch.  $\mathcal{B}$ : Bernoulli noise.  $p$ : Bernoulli probability and  $\ell^{\text{0MR}}$  strength.  $\lambda$ : A hyperparameter.  $\tau$ : A predefined threshold.  $\Delta_p$ : A small value.

Control  $p$ :  $r(\mathbf{x}) = \mathbb{E}[\text{sign}(D(\mathbf{x}))]$ ,  $p_{t+1} = p_t + \Delta_p \cdot \text{sign}(r(\mathbf{x}) - \tau)$ .

CHAIN is applied **separately** to real and fake data batches.

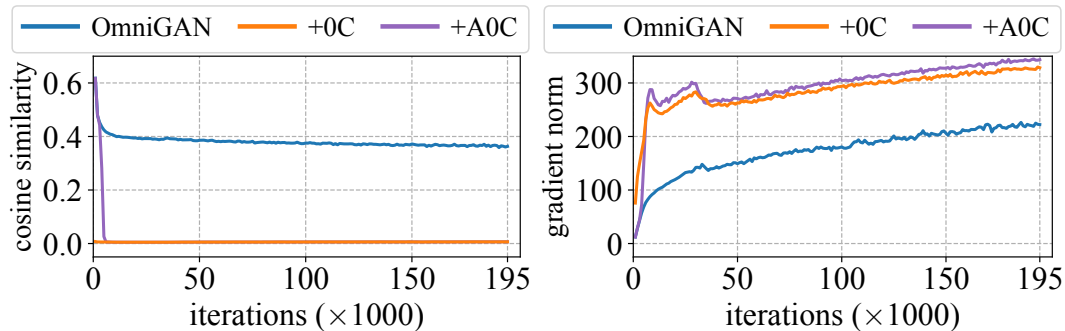


# Experiments: Analysis of gradient issue in centering



0C: centering. A0C: adaptive centering. On 10% CIFAR-10.

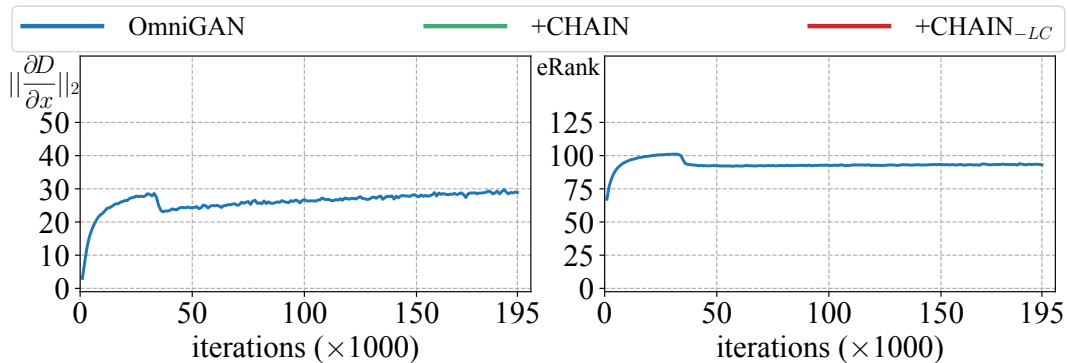
# Experiments: Analysis of gradient issue in centering



0C: centering. A0C: adaptive centering. On 10% CIFAR-10.

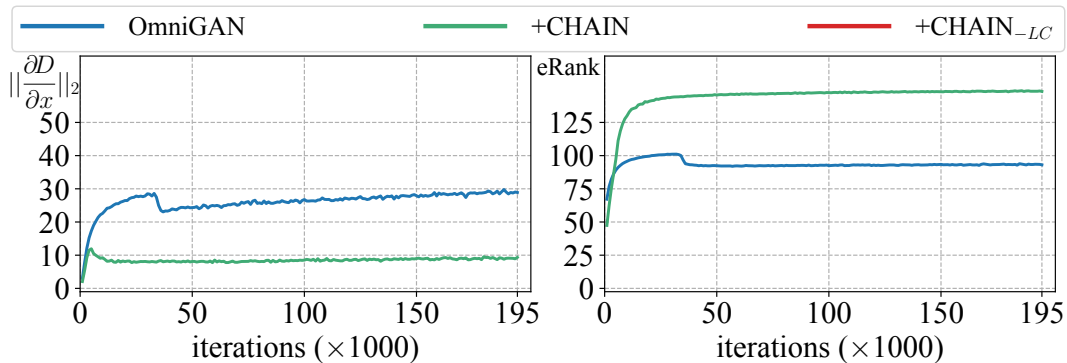
**Centering reduces similarity and raises gradient.**

# Experiments: Analysis of gradient issue in scaling



—LC: without Lipschitz constraint. On 10% CIFAR-10.

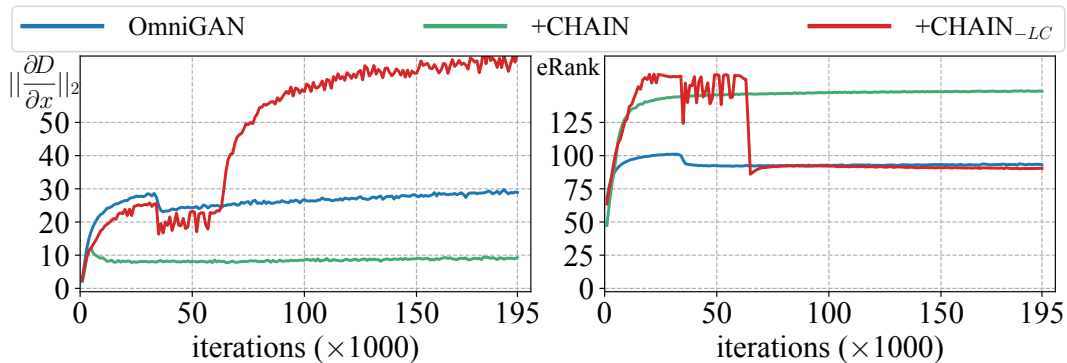
# Experiments: Analysis of gradient issue in scaling



—LC: without Lipschitz constraint. On 10% CIFAR-10.

**CHAIN reduces latent gradients**

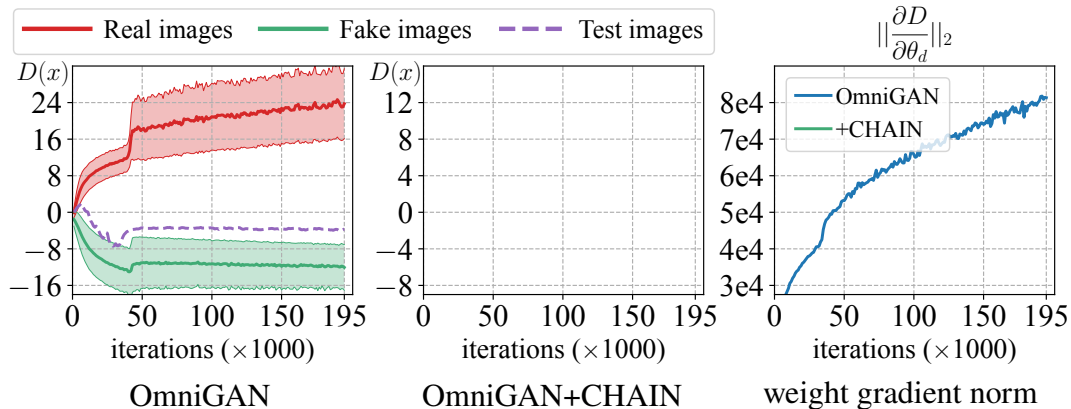
# Experiments: Analysis of gradient issue in scaling



—LC: without Lipschitz constraint. On 10% CIFAR-10.

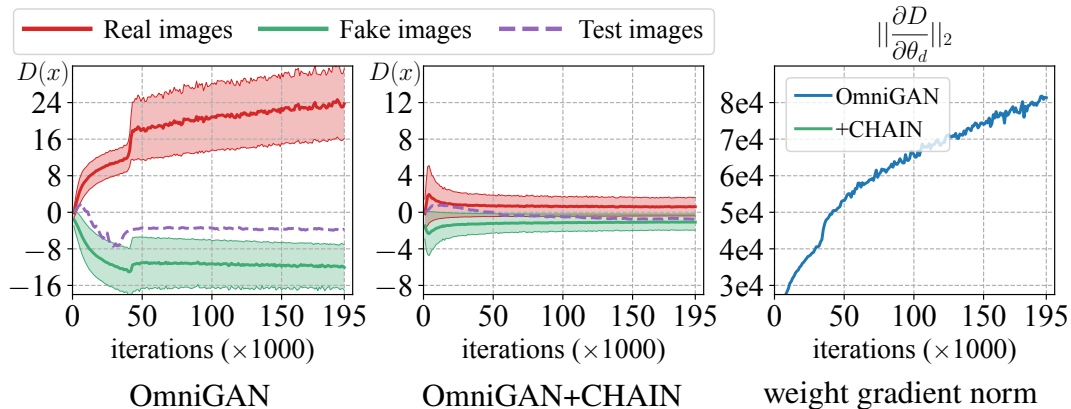
**CHAIN reduces latent gradients while removing LC raises gradient and impairs feature eRank.**

# Experiments: Analysis of generalization of CHAIN



$D(x)$ : discriminator output. On 10% CIFAR-10.

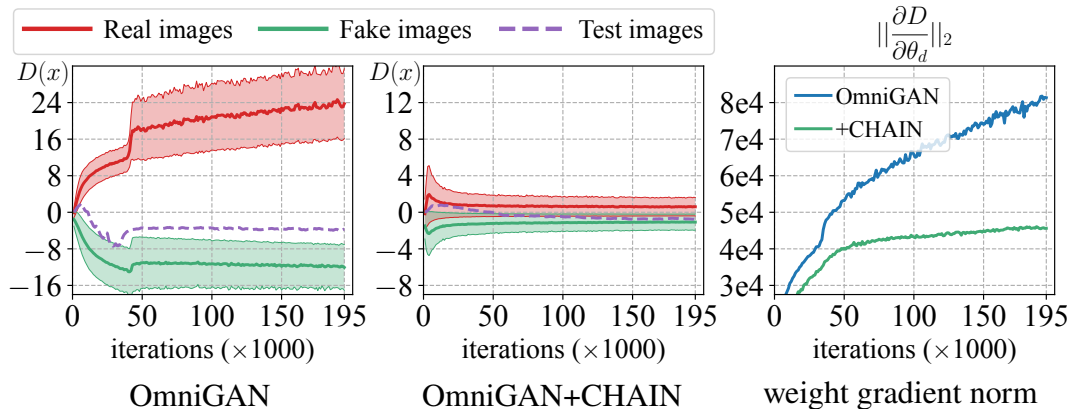
# Experiments: Analysis of generalization of CHAIN



$D(x)$ : discriminator output. On 10% CIFAR-10.

**CHAIN reduces discrepancies among real/fake/testing data**

# Experiments: Analysis of generalization of CHAIN



$D(x)$ : discriminator output. On 10% CIFAR-10.

**CHAIN reduces discrepancies among real/fake/testing data and  $D$ 's weight gradients.**



# Experiments: Comparison with state of the arts

Method	CIFAR-10						CIFAR-100					
	10% data		20% data		100% data		10% data		20% data		100% data	
	IS↑	tFID↓	IS↑	tFID↓	IS↑	tFID↓	IS↑	tFID↓	IS↑	tFID↓	IS↑	tFID↓
BigGAN	8.24	31.45	8.74	16.20	9.21	5.48	7.58	50.79	9.94	25.83	11.02	7.86
+CHAIN	<b>8.63</b>	<b>12.02</b>	<b>8.98</b>	<b>8.15</b>	<b>9.49</b>	<b>4.18</b>	<b>10.04</b>	<b>13.13</b>	<b>10.15</b>	<b>11.58</b>	<b>11.16</b>	<b>6.04</b>
LeCam+DA	8.81	12.64	9.01	8.53	9.45	4.32	9.17	22.75	10.12	15.96	11.25	6.45
+CHAIN	<b>8.96</b>	<b>8.54</b>	<b>9.27</b>	<b>5.92</b>	<b>9.52</b>	<b>3.51</b>	<b>10.11</b>	<b>12.69</b>	<b>10.62</b>	<b>9.02</b>	<b>11.37</b>	<b>5.26</b>
OmniGAN+ADA	7.86	40.05	9.41	27.04	10.24	4.95	8.95	44.65	12.07	13.54	13.07	6.12
+CHAIN	<b>10.10</b>	<b>6.22</b>	<b>10.26</b>	<b>3.98</b>	<b>10.31</b>	<b>2.22</b>	<b>12.70</b>	<b>9.49</b>	<b>12.98</b>	<b>7.02</b>	<b>13.98</b>	<b>4.02</b>

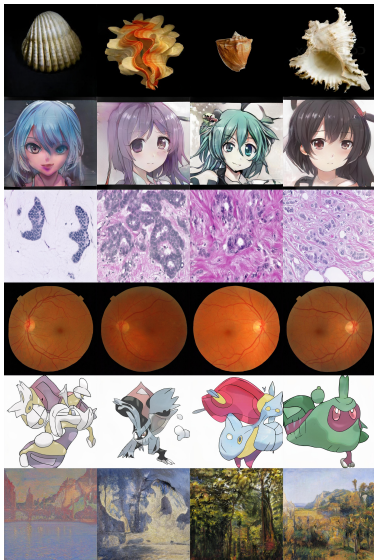
Method (FID↓)	Shells	Skulls	AnimeFace	BreCaHAD	MessidorSet1	Pokemon	ArtPainting
	64 imgs	97 imgs	120 imgs	162 imgs	400 imgs	833 imgs	1000 imgs
FastGAN	138.50	97.87	54.05	63.83	38.33	45.70	43.21
FreGAN	123.75	84.58	49.09	<b>57.87</b>	34.61	39.09	43.14
FastGAN- $D_{\text{big}}$	171.35	165.64	76.02	68.63	37.38	53.48	43.04
+CHAIN	<b>78.62</b>	<b>82.47</b>	<b>46.27</b>	58.98	<b>28.76</b>	<b>31.94</b>	<b>38.83</b>

# Experiments: Comparison with state of the arts

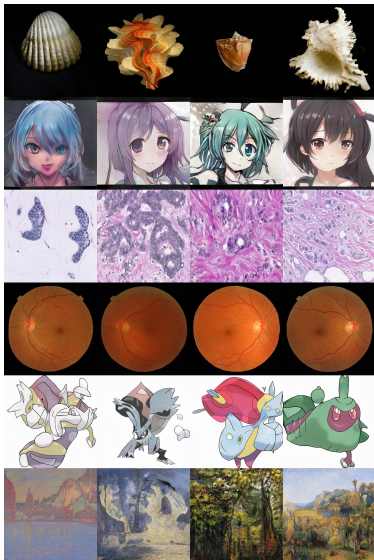
Method	2.5% ImageNet			5% ImageNet			10% ImageNet		
	IS↑	tFID↓	vFID↓	IS↑	tFID↓	vFID↓	IS↑	tFID↓	vFID↓
BigGAN	8.61	101.62	100.09	6.27	90.32	88.01	12.44	50.75	49.84
+CHAIN	<b>14.68</b>	<b>30.66</b>	<b>29.32</b>	<b>17.34</b>	<b>21.13</b>	<b>19.95</b>	<b>20.45</b>	<b>14.70</b>	<b>13.84</b>
ADA	7.93	67.84	66.55	11.56	47.56	46.25	14.82	31.75	30.68
+CHAIN	<b>16.57</b>	<b>23.01</b>	<b>21.90</b>	<b>19.15</b>	<b>16.14</b>	<b>15.17</b>	<b>22.04</b>	<b>12.91</b>	<b>12.17</b>

Method (FID↓)	100-shot			Animal Face	
	Obama	GrumpyCat	Panda	Cat	Dog
StyleGAN2	80.20	48.90	34.27	71.71	131.90
+CHAIN	<b>28.72</b>	<b>27.21</b>	<b>9.51</b>	<b>38.93</b>	<b>53.27</b>
AdvAug	52.86	31.02	14.75	47.40	68.28
DA	46.87	27.08	12.06	42.44	58.85
InsGen	32.42	22.01	9.85	33.01	44.93
FakeCLR	26.95	19.56	8.42	26.34	42.02
KDDLGAN	29.38	19.65	8.41	31.89	50.22
AugSelfGAN	26.00	19.81	8.36	30.53	48.19
DA+CHAIN	<b>22.87</b>	<b>17.57</b>	<b>6.93</b>	<b>19.58</b>	<b>30.88</b>

# Generated Images and Conclusions



# Generated Images and Conclusions



## Conclusions:

- CHAIN reduces real-fake discrepancy and discriminator weight gradients, improving generalization.
- CHAIN lowers latent feature gradients in discriminator, enhancing stability.