



Australian
National
University



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

Yao Ni[†] Piotr Koniusz^{§,†}

[†]The Australian National University

[§]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):

$$\epsilon_{\text{gan}} \propto 2 d_{\mathcal{H}}(p_{\cdot, \hat{n}} \parallel q_{\cdot, \hat{n}}) - 2 d_{\mathcal{H}}(p_{\cdot, \hat{n}^*} \parallel q_{\cdot, \hat{n}^*})$$

$d_{\mathcal{H}}$: discrepancy over discriminator set. $p_{\cdot, \hat{n}}$: unseen/seen real data. $q_{\cdot, \hat{n}}$: ideal/seen fake.

Methods: Generalization error of GANs

Lemma 3.1 (GAN generalization error on function set):

$$\epsilon_{\text{gan}} \propto 2 d_{\mathcal{H}}(p_{\cdot, \hat{n}}^q, 2 d_{\mathcal{H}}(p_{\hat{n}^*}^*, \hat{n}^q))$$

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

$\hat{n}^* \hat{n} \rightarrow \tilde{n}$ Lowering real/fake discrepancy aids generalization

Methods: Generalization error of GANs

Lemma 3.1 (GAN generalization error on function set):

$$\epsilon_{\text{gan}} \propto 2 d_{\mathcal{H}}(p, \hat{q}_n) - 2 d_{\mathcal{H}}(p, \hat{q}_n^*)$$

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

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

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

Methods: Generalization error of GANs

Lemma 3.2 (GAN generalization error on neural network):

$$\epsilon_{\text{gan}}^{\text{nn}} \propto 2! \|r_{d'}\|_2 \|F_{d'}\|_2 \left(4R \frac{\|d\|_2^2}{\epsilon^2} + \frac{1}{n} \right) \left(\epsilon^2 \|\mathbf{H}_{\max}\| + \|\hat{\mathbf{H}}_{\max}\| \right)$$

$\epsilon_{\text{gan}}^{\text{nn}}$: GAN generalization error on neural networks, $\epsilon > 0$. d : D 's weights. $r_{d'}$: real gradient, top Hessian eigenvalue. $F_{d'}$: fake versions. R : related to $\|d\|_2^2$ and data size n .

Methods: Generalization error of GANs

Lemma 3.2 (GAN generalization error on neural network):

$$\epsilon_{\text{gan}}^{\text{nn}} \propto 2! \left(k r_{d'} k_2 + k \mathcal{F}_{d'} k_2 \right) 4R \frac{k_{d'} k_2^2}{\lambda^2}, \frac{1}{n} \left(\lambda^2 \left| \mathbf{H}_{\max} \right| + \left| \hat{\mathbf{H}}_{\max} \right| \right)$$

$\epsilon_{\text{gan}}^{\text{nn}}$: gan on neural networks, $\lambda > 0$. d' : D 's weights. $r_{d'}$: real gradient, top Hessian eigenvalue. $\mathcal{F}_{d'}$: fake versions. R : related to $k_{d'} k_2^2$ 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

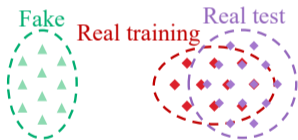


Methods: Motivation for BN

Better generalization on GANs requires:

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

Motivation of BN



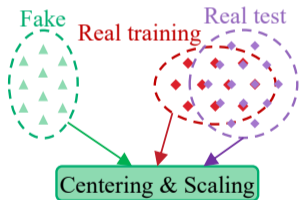
Applying BN **separately** on real/fake batches

Methods: Motivation for BN

Better generalization on GANs requires:

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

Motivation of BN



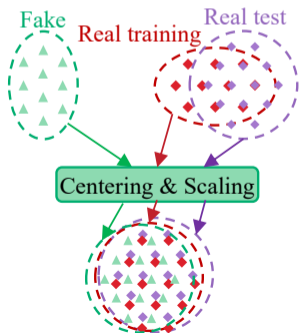
Applying BN **separately** on real/fake batches

Methods: Motivation for BN

Better generalization on GANs requires:

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

Motivation of BN



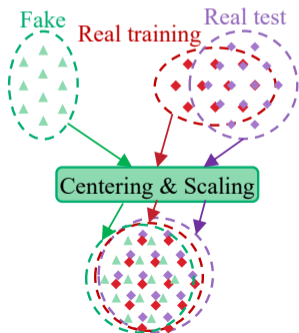
Applying BN **separately** on real/fake batches reduces the fake-real discrepancy via standardization.

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:	Y	AW
Centering:	$\overset{c}{Y}$	Y
Scaling:	$\overset{s}{Y}$	$\overset{c}{Y} \{ \cdot$

Methods: Gradient Issues of BN

Standardization in BN:

Linear transformation:	Y	AW
Centering:	$\overset{c}{Y}$	Y
Scaling:	$\overset{s}{Y}$	$\overset{c}{Y} \{ \cdot \}$

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

$$E_{\mathbf{y}_1, \mathbf{y}_2} \cos \langle \mathbf{y}_1, \mathbf{y}_2 \rangle \neq E_{\overset{c}{\mathbf{y}}_1, \overset{c}{\mathbf{y}}_2} \cos \langle \overset{c}{\mathbf{y}}_1, \overset{c}{\mathbf{y}}_2 \rangle \quad 0$$

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

Methods: Gradient Issues of BN

Standardization in BN:

Linear transformation: $Y \rightarrow AW$

Centering: $\hat{Y}^c \rightarrow Y$

Scaling: $\hat{Y}^s \rightarrow \hat{Y}^c \{ \cdot \}$

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

$$E_{\mathbf{y}_1, \mathbf{y}_2} \cos \langle \mathbf{y}_1, \mathbf{y}_2 \rangle \neq E_{\hat{\mathbf{y}}_1^c, \hat{\mathbf{y}}_2^c} \cos \langle \hat{\mathbf{y}}_1^c, \hat{\mathbf{y}}_2^c \rangle \quad 0$$

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

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

$$\text{diag} \left[\frac{1}{\min_c} \right]_{lc} \quad \frac{1}{\min_c}$$

Lipschitz constant (lc) is large when $\min_c \min_c c$, is small.

Method: CHAIN replaces centering/scaling with OMR/ARMS

mean and root mean square :

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

$$c \quad \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}^2$$

Method: CHAIN replaces centering/scaling with OMR/ARMS

mean and root mean square :

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

$$c = \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}^2$$

0-mean regularization:

$$\lambda^{0MR} p Y q \quad p \quad k \quad k_2^2$$

$\min_{\rho} \min_c c \cdot 10^{-5}$: a hyperparameter.
 $\rho \in [0, 1]$ controls λ^{0MR} and Bernoulli mask $M \in B_{p,q}$

Pytorch-style pseudo code for CHAIN_{batch}

```
# Y: B d H W size; lbd: hyperparameter
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 and root mean square :

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

$$c \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}^2$$

0-mean regularization:

$$\|Y - \text{OMR}(Y)\|_2^2$$

Adaptive RMS normalization:

$$\text{ARMS}(Y) = \frac{Y}{\sqrt{M}} \quad \min_{p, \rho} \mathbb{E} \|Y - \text{OMR}(Y)\|_2^2$$

$\rho \in [0, 1]$ controls $\|Y - \text{OMR}(Y)\|_2^2$ and Bernoulli mask $M \in \{0, 1\}^{B \times p \times q}$

Pytorch-style pseudo code for CHAIN_{batch}

```
# Y: B d H W size; lbd: hyperparameter
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

$$k \mathbf{y}_c k_2^2 \propto k \mathbf{y}_c k_2^2 \frac{(1-\rho) \min_c \{c+\rho\}}{c} \frac{2\rho}{B_c} \mathbf{y}_c^T \mathbf{g}_c \rho^2$$

$$k \mathbf{w}_c k_2^2 \propto \max_c k \mathbf{y}_c k_2^2$$

$\mathbf{y}_c, \mathbf{y}_c$: c -th column of gradient for CHAIN input/output Y, Ψ . λ_{\max} : top eigenvalue of A . \mathbf{g}_c : c -th column of $\mathbf{Y}^T \mathbf{Y}$. \mathbf{w}_c : c -th column of grad for W . $\rho \in [0, 1]$

Method: CHAIN reduces feature and weight gradients in D

$$k \mathbf{y}_c k_2^2 \propto k \mathbf{y}_c k_2^2 \frac{(1-\rho) \frac{c+\rho}{c} \min}{c}^2 \frac{2\rho^2}{B_c} \rho \mathbf{y}_c^T \mathbf{g}_c \rho^2$$

$$k \mathbf{w}_c k_2^2 \propto \max k \mathbf{y}_c k_2^2$$

$\mathbf{y}_c, \mathbf{y}_c$: c -th column of gradient for CHAIN input/output Y, Ψ . λ_{\max} : top eigenvalue of A . \mathbf{g}_c : c -th column of $\mathbf{Y}^T \mathbf{Y}$. \mathbf{w}_c : c -th column of grad for W . $\rho \in [0, 1]$

$$\frac{(1-\rho) \frac{c+\rho}{c} \min}{c} \propto 1$$

Method: CHAIN reduces feature and weight gradients in D

$$k \mathbf{y}_c k_2^2 \propto k \mathbf{y}_c k_2^2 \frac{(1-\rho) \frac{c+\rho}{c} \min}{c}^2 \frac{2\rho^2}{B_c} \rho \mathbf{y}_c^T \mathbf{g}_c \rho^2$$

$$k \mathbf{w}_c k_2^2 \propto \max k \mathbf{y}_c k_2^2$$

$\mathbf{y}_c, \mathbf{y}_c$: c -th column of gradient for CHAIN input/output Y, Ψ . λ_{\max} : top eigenvalue of A . \mathbf{g}_c : c -th column of $\mathbf{Y}^T \mathbf{Y}$. \mathbf{w}_c : c -th column of grad for W . $\rho \in [0, 1]$

$$\frac{(1-\rho) \frac{c+\rho}{c} \min}{c} \propto 1 \quad \rho \mathbf{y}_c^T \mathbf{g}_c \rho^2 \neq 0$$

Method: CHAIN reduces feature and weight gradients in D


$$k \mathbf{y}_c k_2^2 \propto k \mathbf{y}_c k_2^2 \frac{(1-\rho) \frac{c+\rho}{c} \min}{c}^2 \frac{2\rho^2}{B_c} \rho \mathbf{y}_c^T \mathbf{g}_c \rho^2$$

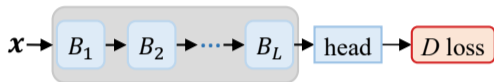
$$k \mathbf{w}_c k_2^2 \propto \max_c k \mathbf{y}_c k_2^2$$

$\mathbf{y}_c, \mathbf{y}_c$: c -th column of gradient for CHAIN input/output Y, Ψ . λ_{\max} : top eigenvalue of A . \mathbf{g}_c : c -th column of $\mathbf{H} = Y \{ \dots \}$. \mathbf{w}_c : c -th column of grad for W . $\rho \in [0, 1]$

$\frac{(1-\rho) \frac{c+\rho}{c} \min}{c} \propto 1 - \rho \mathbf{y}_c^T \mathbf{g}_c \rho^2 \neq 0 \iff$ CHAIN reduces $k \mathbf{y}_c k_2$ and $k \mathbf{w}_c k_2$.


Pipeline

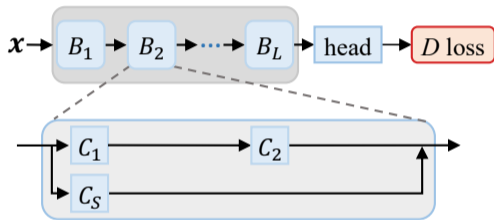
Discriminator with CHAIN 



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

Pipeline

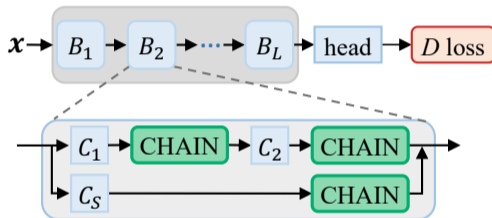
Discriminator with CHAIN 



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


Pipeline

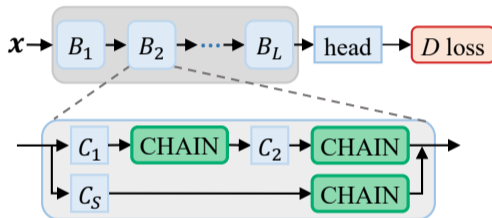
Discriminator with CHAIN



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:

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

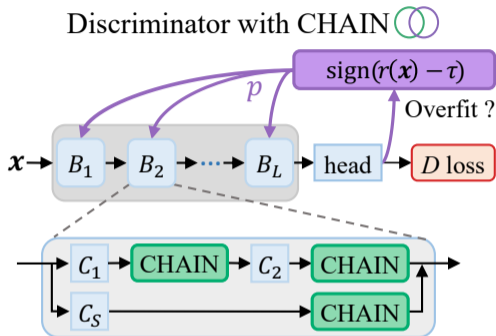
Adaptive Root Mean Square normalization:

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

$\mathbf{M} \sim \mathcal{B}(p)$. $\boldsymbol{\mu}$, $\boldsymbol{\psi}$ 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. B : Bernoulli noise. p : Bernoulli probability and 0MR strength. λ : A hyperparameter.

Pipeline



CHAIN: 0MR & ARMS

0-Mean Regularization loss:

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

Adaptive Root Mean Square normalization:

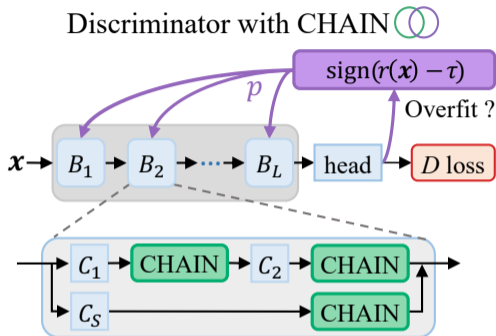
$$\text{ARMS}(Y) = (1 - M) \odot Y + M \odot \frac{Y}{\psi} \cdot \psi_{\min}$$

$M \sim \mathcal{B}(p)$. μ, ψ 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. B : Bernoulli noise. p : Bernoulli probability and 0MR strength. λ : A hyperparameter. τ : A predefined threshold. ρ : A small value.

Control p : $r(p, x) = \text{sign}(p - D(p, x)) \cdot \rho$, $p_{t+1} = p_t + r(p_t, x) \cdot \rho$

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:

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

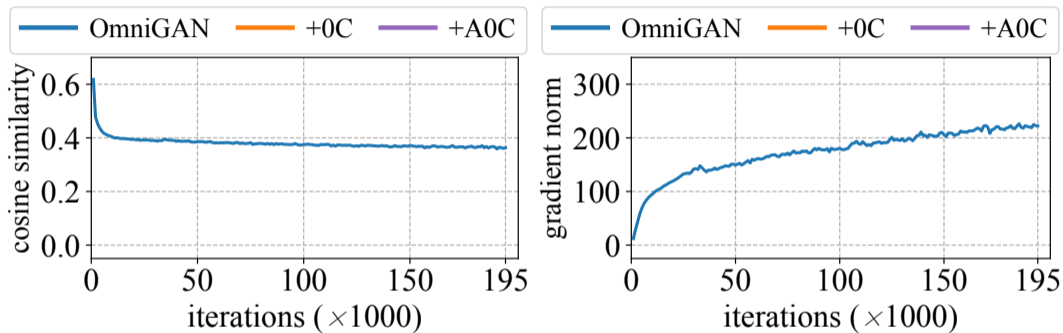
$\mathbf{M} \sim \mathcal{B}(p)$. $\boldsymbol{\mu}, \boldsymbol{\psi}$ are mean and root mean square of feature map \mathbf{Y} .

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

Control p : $r \times x \times q$ $\rightarrow \text{sign}(p) \times D \times p \times x \times q$, $p_{t+1} = \rho_t \times p + \rho \times \text{sign}(p) \times r \times x \times q$

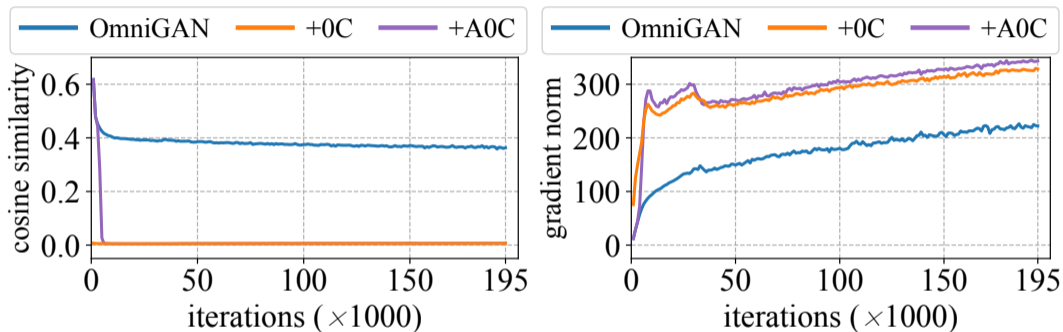
CHAIN is applied **separately** to real and fake data batches.

Experiments: Analysis of gradient issue in centering



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

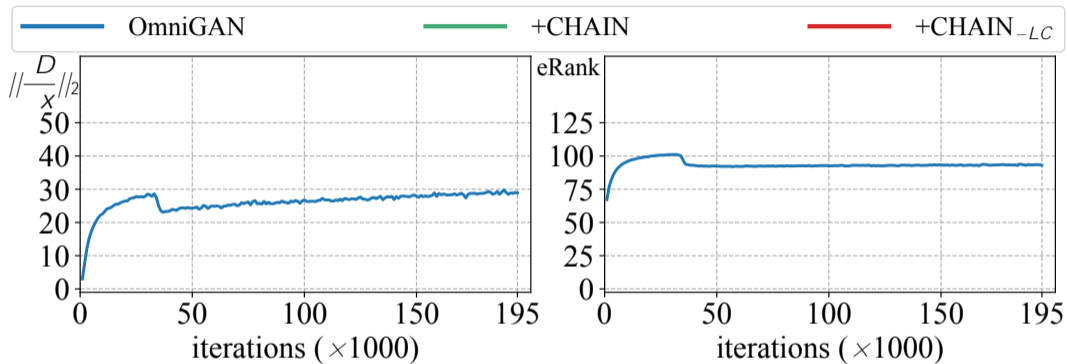
Experiments: Analysis of gradient issue in centering



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

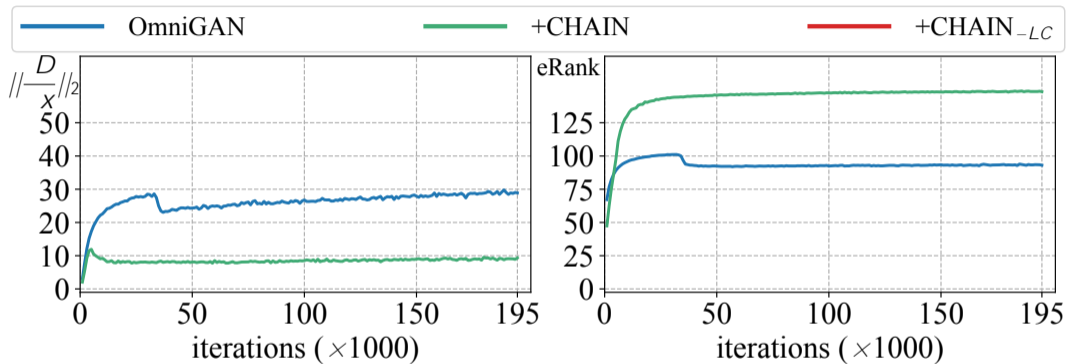
Centering reduces similarity and raises gradient.

Experiments: Analysis of gradient issue in scaling



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

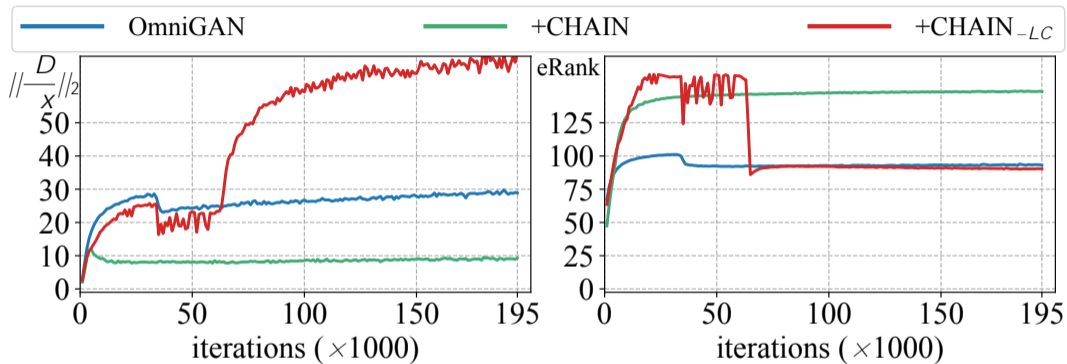
Experiments: Analysis of gradient issue in scaling



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

CHAIN reduces latent gradients

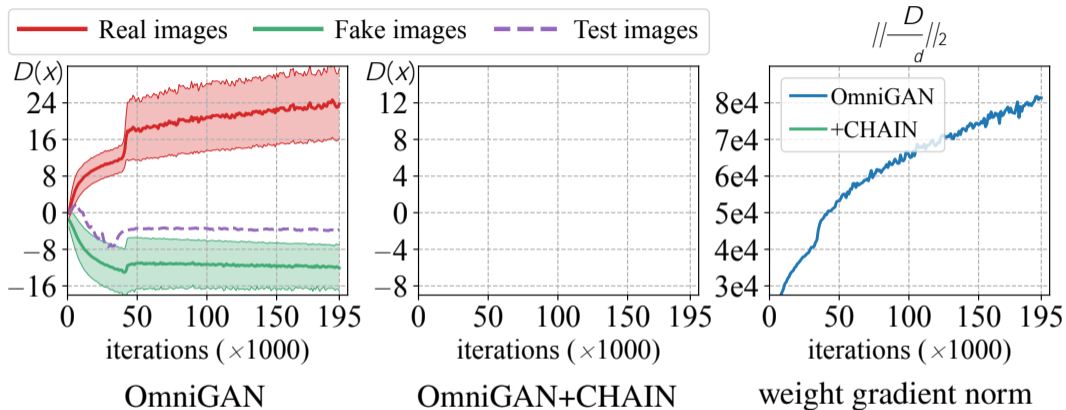
Experiments: Analysis of gradient issue in scaling



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

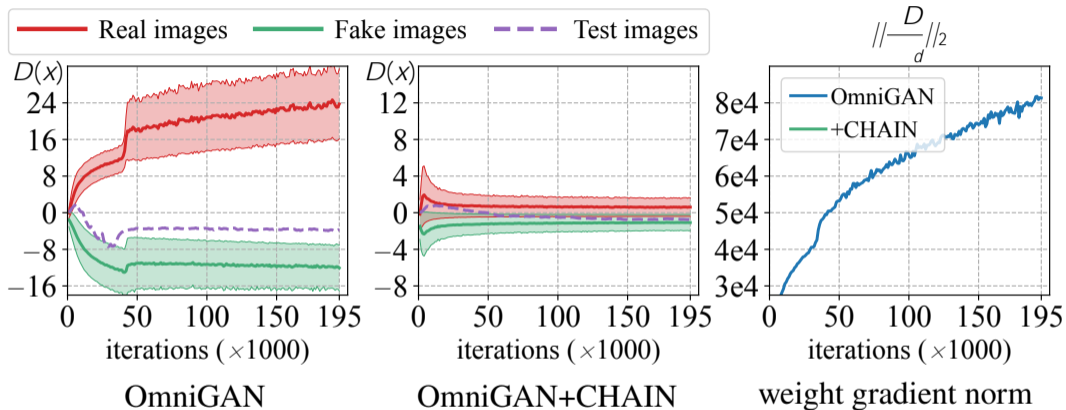
CHAIN reduces latent gradients while removing LC raises gradient and impairs feature eRank.

Experiments: Analysis of generalization of CHAIN



$D_p x q$: discriminator output. On 10% CIFAR-10.

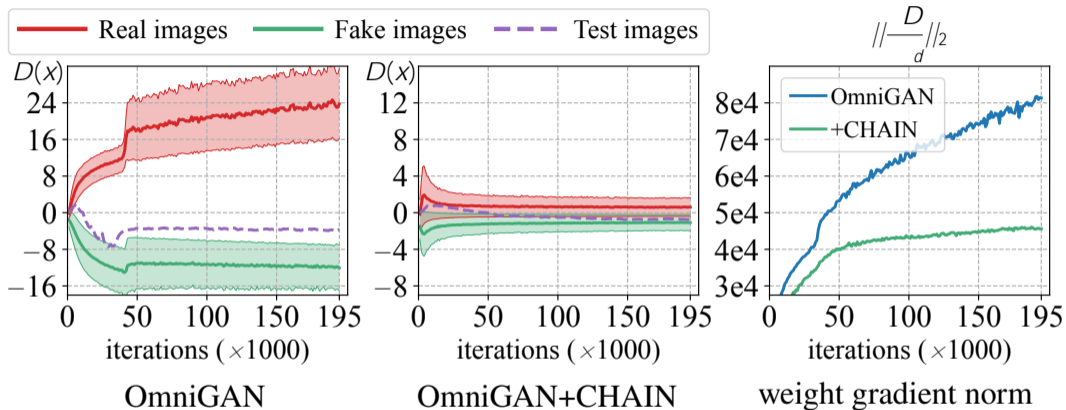
Experiments: Analysis of generalization of CHAIN



$D_p x q$: discriminator output. On 10% CIFAR-10.

CHAIN reduces discrepancies among real/fake/testing data

Experiments: Analysis of generalization of CHAIN



$D_p x q$: 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 $\bar{\sigma}$	tFID $\bar{\sigma}$	IS $\bar{\sigma}$	tFID $\bar{\sigma}$	IS $\bar{\sigma}$	tFID $\bar{\sigma}$	IS $\bar{\sigma}$	tFID $\bar{\sigma}$	IS $\bar{\sigma}$	tFID $\bar{\sigma}$	IS $\bar{\sigma}$	tFID $\bar{\sigma}$
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	8.63	12.02	8.98	8.15	9.49	4.18	10.04	13.13	10.15	11.58	11.16	6.04
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	8.96	8.54	9.27	5.92	9.52	3.51	10.11	12.69	10.62	9.02	11.37	5.26
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	10.10	6.22	10.26	3.98	10.31	2.22	12.70	9.49	12.98	7.02	13.98	4.02

Method (FID $\bar{\sigma}$)	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	57.87	34.61	39.09	43.14
FastGAN- D_{big}	171.35	165.64	76.02	68.63	37.38	53.48	43.04
+CHAIN	78.62	82.47	46.27	58.98	28.76	31.94	38.83

Experiments: Comparison with state of the arts

Method	2.5% ImageNet			5% ImageNet			10% ImageNet		
	IS0	tFID0	vFID0	IS0	tFID0	vFID0	IS0	tFID0	vFID0
BigGAN	8.61	101.62	100.09	6.27	90.32	88.01	12.44	50.75	49.84
+CHAIN	14.68	30.66	29.32	17.34	21.13	19.95	20.45	14.70	13.84
ADA	7.93	67.84	66.55	11.56	47.56	46.25	14.82	31.75	30.68
+CHAIN	16.57	23.01	21.90	19.15	16.14	15.17	22.04	12.91	12.17

Method (FID0)	100-shot			Animal Face	
	Obama	GrumpyCat	Panda	Cat	Dog
StyleGAN2	80.20	48.90	34.27	71.71	131.90
+CHAIN	28.72	27.21	9.51	38.93	53.27
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	22.87	17.57	6.93	19.58	30.88

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.