

Summary

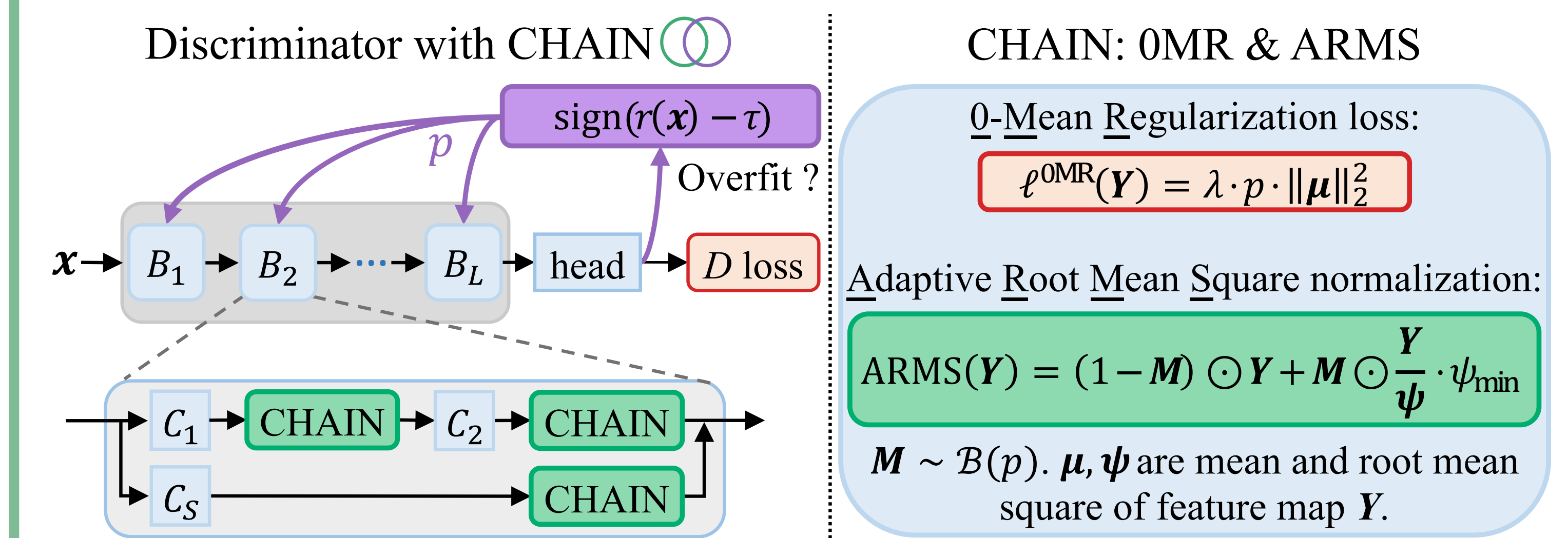
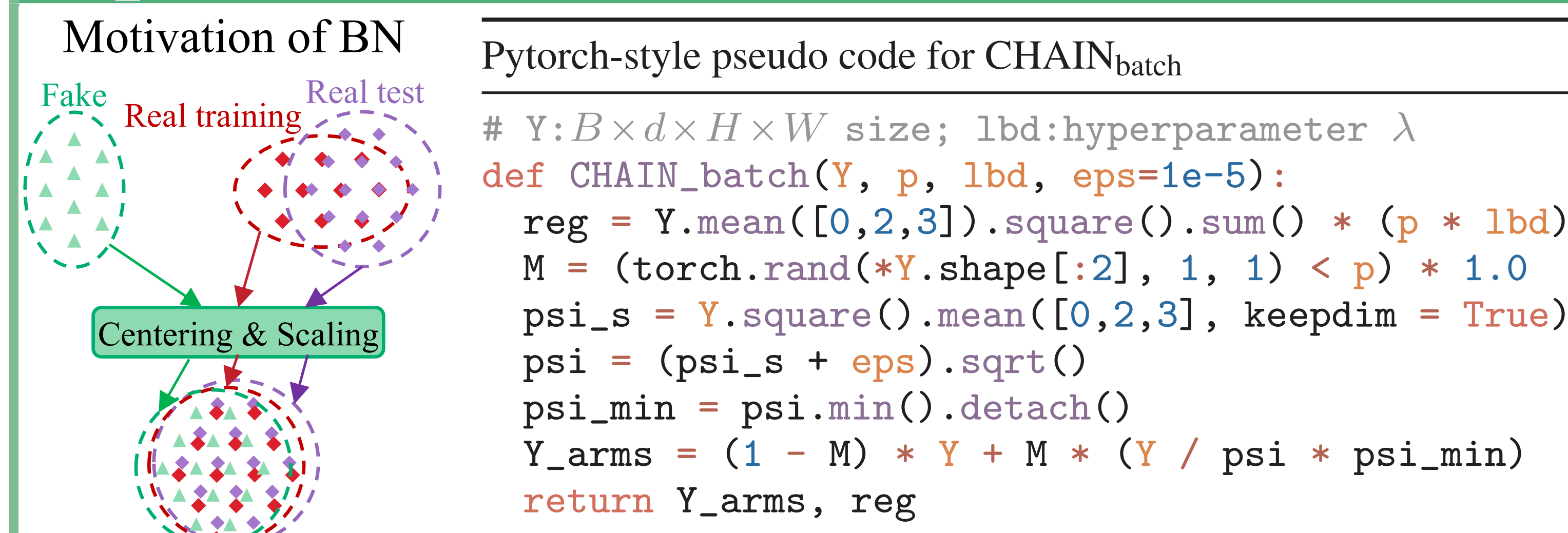
Background: Limited data in GANs causes discriminator overfitting and unstable training. Batch Normalization (BN) boosts generalization and training stability but is rarely used in data-efficient GAN discriminators as it impairs performance.

Goal: Integrate BN into data-efficient GAN discriminators.

Contributions:

- We derive a new error bound emphasizing reduced discriminator weight gradients to enhanced GAN generalization.
- We find that BN in the discriminator tends to cause gradient explosion due to its centering and scaling steps.
- We design CHAIN by replacing BN’s centering and scaling steps with 0-mean regularization and adaptive RMS normalization. We show that CHAIN improves stability and generalization of GANs by reducing feature and weight gradients.

Pipeline



$p \in [0, 1]$: Bernoulli probability and ℓ^{OMR} strength, updated with $r(\mathbf{x}) = \mathbb{E}[\text{sign}(D(\mathbf{x}))]$, $p_{t+1} = p_t + \Delta_p \cdot \text{sign}(r(\mathbf{x}) - \tau)$.

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

\mathbf{x} : Real image. B_l : l -th block. C_s : Conv. in skip branch. D : Discriminator. τ : A predefined threshold. Δ_p : A small value. \mathcal{B} : Bernoulli noise. λ : Hyperparameter.

Method

Lowering real/fake discrepancy aids generalization:

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

$d_{\mathcal{H}}$: discrepancy of D . $\mu, \hat{\mu}_n$: unseen/seen real data. $\nu_n^*, \hat{\nu}_n$: ideal/seen fake. $\nu_n^* \approx \hat{\mu}_n$ implies lowering real/fake gap reduces ϵ_{gan} . μ inaccessible, thus we derive $\epsilon_{\text{gan}}^{\text{nn}}$.

Lowering weight gradient norms of D aids generalization:

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

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

Separate BN reduces discrepancy but enlarges gradient:

transform: $\mathbf{Y} = \mathbf{A}\mathbf{W}$ centering: $\tilde{\mathbf{Y}} = \mathbf{Y} - \boldsymbol{\mu}$ scaling: $\hat{\mathbf{Y}} = \tilde{\mathbf{Y}} / \boldsymbol{\sigma}$

(Centering) similarity dropping causes feature divergence:

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

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

(Scaling) unbounded Lipschitz causes gradient explosion:

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

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

CHAIN replaces centering/scaling with OMR/ARMS:

$$\text{mean } \boldsymbol{\mu}: \mu_c = \frac{1}{B \times H \times W} \sum_{b,h,w} Y_{b,c,h,w}$$

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

0-mean regularization: $\ell^{\text{OMR}}(\mathbf{Y}) = \lambda \cdot p \cdot \|\boldsymbol{\mu}\|_2^2$

Adaptive root mean square normalization:

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

$\epsilon = 10^{-5}$. λ : a hyperparameter. p controls ℓ^{OMR} and Bernoulli mask $\mathbf{M} \sim \mathcal{B}(p)$

CHAIN reduces gradients of features and weights in D :

$$\|\Delta \mathbf{y}_c\|_2^2 \leq \|\Delta \tilde{\mathbf{y}}_c\|_2^2 \left(\frac{(1-p)\psi_c + p\psi_{\min}}{\psi_c} \right)^2 - \frac{2(1-p)p\psi_{\min}}{B\psi_c} (\Delta \tilde{\mathbf{y}}_c^T \tilde{\mathbf{y}}_c)^2$$

$$\|\Delta \mathbf{w}_c\|_2^2 \leq \lambda_{\max}^2 \|\Delta \mathbf{y}_c\|_2^2$$

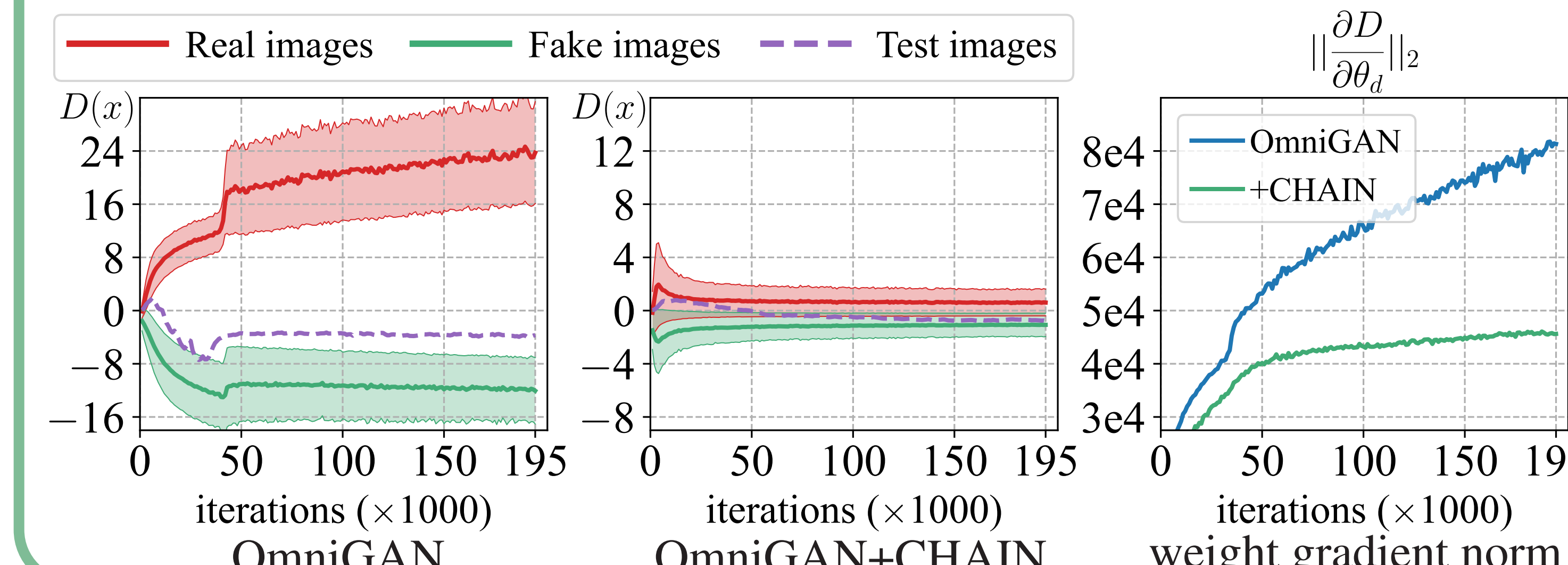
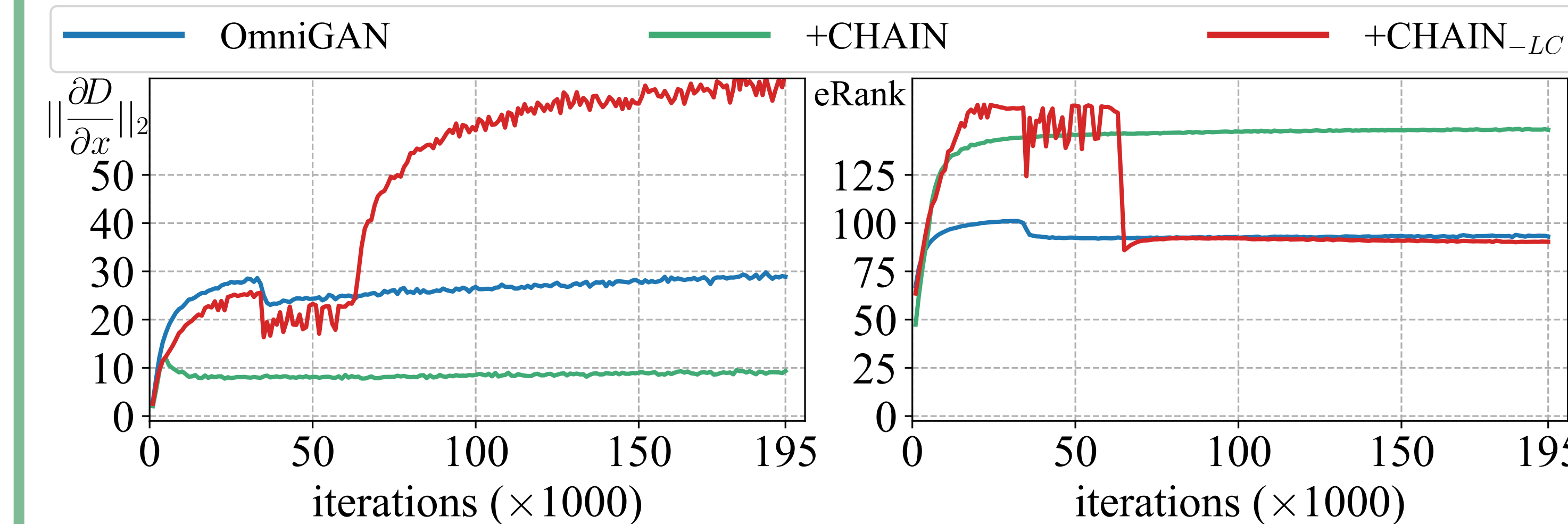
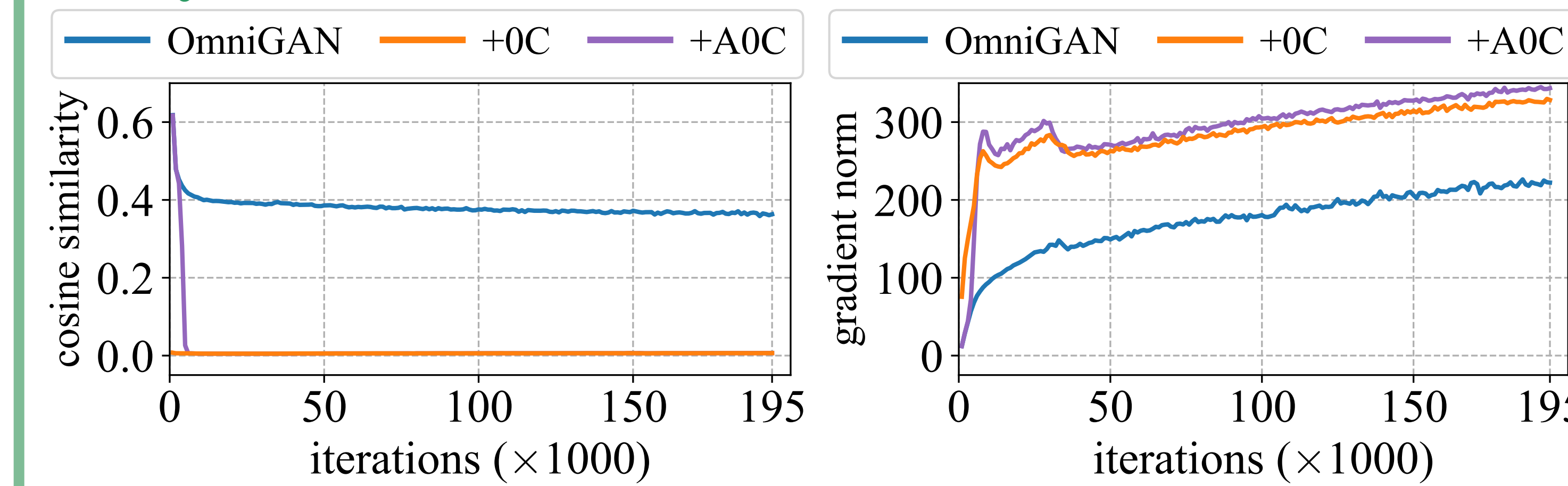
$\Delta \mathbf{y}_c, \Delta \tilde{\mathbf{y}}_c$: c -th column of gradient for CHAIN input/output $\mathbf{Y}, \tilde{\mathbf{Y}}$. λ_{\max} : top eigenvalue of \mathbf{A} . $\tilde{\mathbf{y}}_c$: c -th column of $\tilde{\mathbf{Y}} = \mathbf{Y} / \boldsymbol{\psi}$. $\Delta \mathbf{w}_c$: c -th column of grad for \mathbf{W} .

Experimental Results

Comparison with the state-of-the-art methods:

Method	CIFAR-10						CIFAR-100						Method	2.5% ImageNet			5% ImageNet			10% ImageNet		
	10% data		20% data		100% data		10% data		20% data		100% data			IS↑	tFID↓	vFID↓	IS↑	tFID↓	vFID↓	IS↑	tFID↓	vFID↓
	IS↑	tFID↓	IS↑	tFID↓	IS↑	tFID↓	IS↑	tFID↓	IS↑	tFID↓	IS↑	tFID↓		IS↑	tFID↓	vFID↓	IS↑	tFID↓	vFID↓	IS↑	tFID↓	vFID↓
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	BigGAN	8.61	101.62	100.09	6.27	90.32	88.01	12.44	50.75	49.84
+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	+CHAIN	14.68	30.66	29.32	17.34	21.13	19.95	20.45	14.70	13.84
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	ADA	7.93	67.84	66.55	11.56	47.56	46.25	14.82	31.75	30.68
+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	+CHAIN	16.57	23.01	21.90	19.15	16.14	15.17	22.04	12.91	12.17
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	Method (FID↓)			100-shot			Animal Face			
+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				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				
Method (FID↓)													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				
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															

Analyses on 10% CIFAR-10:



Generated Images:

0C: centering. A0C: adaptive centering. Centering reduces similarity and raises gradient.

–LC: without Lipschitz constraint. CHAIN reduces latent gradients while removing LC raises gradient and impairs feature eRank.

$D(x)$: discriminator output. CHAIN reduces D ’s weight gradient, lowering discrepancies among unseen/seen real data and fake data.

