

# Supplementary Material – An Internal Covariate Shift Bounding Algorithm for Deep Neural Networks by Unitizing Layers’ Outputs

You Huang

Yuanlong Yu

In this document, the proofs of all the theorems are given, with an example of the unbounded Earth Mover (EM) distance. Then, we propose an algorithm for estimating the EM distances that measure Internal Covariate Shift (ICS) in practice, based on the Wasserstein Generative Adversarial Network (WGAN) [1]. The estimated distances for the unitization in network training are reported. Besides, for micro-batches, the proposed unitization is compared against Batch Normalization (BN) [3] and Group Normalization (GN) [5]. Code is available at <https://github.com/unknown9567/unitization.git>.

## 1. Proofs of the Theorems

### 1.1. Notations

- $W(p_l^{(t+\Delta t)}, p_l^{(t)})$ : the EM distance between the distributions  $p_l^{(t+\Delta t)}$  and  $p_l^{(t)}$  at the  $(t + \Delta t)$ -th and  $t$ th iterations, respectively, for the  $l$ th layer’s outputs
- $\mu_i^{(t)}$ : the expectation of the  $i$ th element  $x_i$  of a random vector  $\mathbf{x} \sim p_l^{(t)}$
- $(\sigma_i^{(t)})^2$ : the variance of  $x_i$
- $g(\mathbf{x})$ : the vanilla unitization defined as

$$g(\mathbf{x}) = \begin{cases} \frac{\mathbf{x}}{\|\mathbf{x}\|_2} & , \|\mathbf{x}\|_2 \neq 0 \\ \mathbf{c} & , \|\mathbf{x}\|_2 = 0 \end{cases}$$

- $g(\mathbf{x}; \alpha)$ : the modified unitization with a parameter  $\alpha \in [0, 1]$ , defined as

$$g(\mathbf{x}; \alpha) = \begin{cases} \mathbf{c} & , \|\mathbf{x}\|_2 = 0, \alpha = 1 \\ \frac{\mathbf{x}}{\alpha\|\mathbf{x}\|_2 + (1 - \alpha)} & , other \end{cases}$$

- $g(\mathbf{x}; \boldsymbol{\alpha})$ : the modified unitization with a parameter vector  $\boldsymbol{\alpha} \in [0, 1]^d$ , defined as

$$g(\mathbf{x}; \boldsymbol{\alpha}) = \begin{cases} \mathbf{0} & , \|\mathbf{x}\|_2 = 0 \\ ((\|\mathbf{x}\|_2 - 1) \cdot \text{diag}(\boldsymbol{\alpha}) + E)^{-1} \mathbf{x} & , \|\mathbf{x}\|_2 > 0 \end{cases}$$

where  $\text{diag}(\boldsymbol{\alpha})$  is a diagonal matrix of  $\boldsymbol{\alpha}$  and  $E$  is an identity matrix

- $g(\mathbf{x}; \boldsymbol{\alpha}, \epsilon)$ : the practical unitization with a parameter vector  $\boldsymbol{\alpha} \in \mathbb{R}^d$  and  $\epsilon > 0$ , defined as

$$g(\mathbf{x}; \boldsymbol{\alpha}, \epsilon) = \left[ \frac{1}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \boldsymbol{\alpha} + (\mathbf{1} - \boldsymbol{\alpha}) \right] \odot \mathbf{x},$$

where  $\odot$  represents element-wise production

## 1.2. Useful facts

There are two facts that can be derived directly from the EM distance, and they'll be used in the proof of the theorems.

**Fact 1.**

$$\begin{aligned} W(p_i^{(t+\Delta t)}, p_i^{(t)}) &= \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})] \\ &= \sup_{\|f\|_L \leq 1} \left| \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})] \right|, \end{aligned}$$

since  $f$  is a 1-Lipschitz function iff  $-f$  is a 1-Lipschitz function, and

$$\left| \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})] \right| = \max \left\{ \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})], \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [-f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [-f(\mathbf{y})] \right\}.$$

**Fact 2.**

$$\begin{aligned} W(p_i^{(t+\Delta t)}, p_i^{(t)}) &= \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})] \\ &= \sup_{f: \|f\|_L \leq 1, f(\mathbf{0})=0} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})], \end{aligned}$$

since for any 1-Lipschitz function  $f$ ,

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})] &= \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - f(\mathbf{0}) - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})] + f(\mathbf{0}) \\ &= \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x}) - f(\mathbf{0})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y}) - f(\mathbf{0})] \\ &= \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [\tilde{f}(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [\tilde{f}(\mathbf{y})], \end{aligned}$$

where  $\tilde{f}(\mathbf{x}) = f(\mathbf{x}) - f(\mathbf{0})$  is still a 1-Lipschitz function, and satisfies  $\tilde{f}(\mathbf{0}) = 0$ .

### 1.3. Proofs

**Theorem 1.** Suppose that  $|\mu_i^{(t)}| < \infty, |\mu_i^{(t+\Delta t)}| < \infty, 1 \leq i \leq d$ . Then,

$$W(p_i^{(t+\Delta t)}, p_i^{(t)}) \leq \sum_{i=1}^d (\sigma_i^{(t+\Delta t)})^2 + \sum_{i=1}^d (\sigma_i^{(t)})^2 + \left( \sum_{i=1}^d (\mu_i^{(t+\Delta t)} - \mu_i^{(t)})^2 \right)^{\frac{1}{2}} + 2.$$

*Proof.* Denote by  $\mathbb{I}_A(x)$  an indicator function defined by  $\mathbb{I}_A(x) = 1$  if  $x \in A$ , otherwise  $\mathbb{I}_A(x) = 0$ . Then, according to **Fact 1**,

$$\begin{aligned} W(p_i^{(t+\Delta t)}, p_i^{(t)}) &= \sup_{\|f\|_L \leq 1} \left| \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})] \right| \\ &\leq \sup_{\|f\|_L \leq 1} \left| \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x}) - f(\boldsymbol{\mu}^{(t+\Delta t)})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y}) - f(\boldsymbol{\mu}^{(t)})] \right| + |f(\boldsymbol{\mu}^{(t+\Delta t)}) - f(\boldsymbol{\mu}^{(t)})| \\ &\leq \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [|f(\mathbf{x}) - f(\boldsymbol{\mu}^{(t+\Delta t)})|] + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [|f(\mathbf{y}) - f(\boldsymbol{\mu}^{(t)})|] + |f(\boldsymbol{\mu}^{(t+\Delta t)}) - f(\boldsymbol{\mu}^{(t)})| \\ &\leq \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [\|\mathbf{x} - \boldsymbol{\mu}^{(t+\Delta t)}\|_2] + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [\|\mathbf{y} - \boldsymbol{\mu}^{(t)}\|_2] + \|\boldsymbol{\mu}^{(t+\Delta t)} - \boldsymbol{\mu}^{(t)}\|_2 \\ &= \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [\mathbb{I}_{\|\mathbf{x} - \boldsymbol{\mu}^{(t+\Delta t)}\|_2 \leq 1}(\mathbf{x}) \|\mathbf{x} - \boldsymbol{\mu}^{(t+\Delta t)}\|_2] + \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [\mathbb{I}_{\|\mathbf{x} - \boldsymbol{\mu}^{(t+\Delta t)}\|_2 > 1}(\mathbf{x}) \|\mathbf{x} - \boldsymbol{\mu}^{(t+\Delta t)}\|_2] \\ &\quad + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [\mathbb{I}_{\|\mathbf{y} - \boldsymbol{\mu}^{(t)}\|_2 \leq 1}(\mathbf{y}) \|\mathbf{y} - \boldsymbol{\mu}^{(t)}\|_2] + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [\mathbb{I}_{\|\mathbf{y} - \boldsymbol{\mu}^{(t)}\|_2 > 1}(\mathbf{y}) \|\mathbf{y} - \boldsymbol{\mu}^{(t)}\|_2] + \|\boldsymbol{\mu}^{(t+\Delta t)} - \boldsymbol{\mu}^{(t)}\|_2 \\ &\leq \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [\mathbb{I}_{\|\mathbf{x} - \boldsymbol{\mu}^{(t+\Delta t)}\|_2 \leq 1}(\mathbf{x}) \cdot 1] + \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [\mathbb{I}_{\|\mathbf{x} - \boldsymbol{\mu}^{(t+\Delta t)}\|_2 > 1}(\mathbf{x}) \|\mathbf{x} - \boldsymbol{\mu}^{(t+\Delta t)}\|_2^2] \\ &\quad + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [\mathbb{I}_{\|\mathbf{y} - \boldsymbol{\mu}^{(t)}\|_2 \leq 1}(\mathbf{y}) \cdot 1] + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [\mathbb{I}_{\|\mathbf{y} - \boldsymbol{\mu}^{(t)}\|_2 > 1}(\mathbf{y}) \|\mathbf{y} - \boldsymbol{\mu}^{(t)}\|_2^2] + \|\boldsymbol{\mu}^{(t+\Delta t)} - \boldsymbol{\mu}^{(t)}\|_2 \\ &\leq 1 + \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} \left[ \sum_{i=1}^d (x_i - \mu_i^{(t+\Delta t)})^2 \right] + 1 + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} \left[ \sum_{i=1}^d (y_i - \mu_i^{(t)})^2 \right] + \|\boldsymbol{\mu}^{(t+\Delta t)} - \boldsymbol{\mu}^{(t)}\|_2 \\ &= \sum_{i=1}^d (\sigma_i^{(t+\Delta t)})^2 + \sum_{i=1}^d (\sigma_i^{(t)})^2 + \left( \sum_{i=1}^d (\mu_i^{(t+\Delta t)} - \mu_i^{(t)})^2 \right)^{\frac{1}{2}} + 2. \end{aligned}$$

□

**Theorem 2.** Suppose that  $C > 0$  is a real number, and  $n \geq 2$  is an integer. Then,

$$\begin{aligned} W(p_i^{(t+\Delta t)}, p_i^{(t)}) &= \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(\mathbf{y})] \\ &\geq |\mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f_{n,C}(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f_{n,C}(\mathbf{y})]|, \end{aligned}$$

where  $f_{n,C}$  is the 1-Lipschitz function defined as

$$f_{n,C}(\mathbf{x}) = \frac{1}{nC^{n-1}d^{\frac{1}{2}}} \left( \sum_{|x_i| \leq C} x_i^n + \sum_{x_i < -C} (-C)^n + \sum_{x_i > C} C^n \right).$$

*Proof.* According to **Fact 1**, it's obvious that the inequality holds if  $f_{n,C}$  is a 1-Lipschitz function. Thus, only the proof of Lipschitz continuity of  $f_{n,C}$  is required. For convenience, let  $f_{n,C} = f_b / (nC^{n-1}d^{\frac{1}{2}})$ , where

$$f_b(\mathbf{x}) = \sum_{|x_i| \leq C} x_i^n + \sum_{x_i < -C} (-C)^n + \sum_{x_i > C} C^n.$$

Then, we'll prove the Lipschitz continuity of  $f_b$ . In fact, for any  $\mathbf{v}, \mathbf{w} \in \mathbb{R}^d$ ,

$$\begin{aligned} |f_b(\mathbf{v}) - f_b(\mathbf{w})| &= \left| \sum_{|v_i| \leq C} v_i^n + \sum_{v_i < -C} (-C)^n + \sum_{v_i > C} C^n - \sum_{|w_i| \leq C} w_i^n - \sum_{w_i < -C} (-C)^n - \sum_{w_i > C} C^n \right| \\ &= \left| \sum_{i: |v_i| \leq C, |w_i| \leq C} (v_i^n - w_i^n) + \sum_{i: |v_i| \leq C, w_i < -C} [v_i^n - (-C)^n] + \sum_{i: |v_i| \leq C, w_i > C} (v_i^n - C^n) \right. \\ &\quad + \sum_{i: v_i < -C, |w_i| \leq C} [(-C)^n - w_i^n] + \sum_{i: v_i < -C, w_i < -C} [(-C)^n - (-C)^n] + \sum_{i: v_i < -C, w_i > C} [(-C)^n - C^n] \\ &\quad \left. + \sum_{i: v_i > C, |w_i| \leq C} (C^n - w_i^n) + \sum_{i: v_i > C, w_i < -C} [C^n - (-C)^n] + \sum_{i: v_i > C, w_i > C} (C^n - C^n) \right| \\ &= \left| \sum_{i: |v_i| \leq C, |w_i| \leq C} (v_i - w_i) \sum_{m=0}^{n-1} v_i^{n-1-m} w_i^m + \sum_{i: |v_i| \leq C, w_i < -C} [v_i - (-C)] \sum_{m=0}^{n-1} v_i^{n-1-m} (-C)^m \right. \\ &\quad + \sum_{i: |v_i| \leq C, w_i > C} (v_i - C) \sum_{m=0}^{n-1} v_i^{n-1-m} C^m + \sum_{i: v_i < -C, |w_i| \leq C} (-C - w_i) \sum_{m=0}^{n-1} (-C)^{n-1-m} w_i^m \\ &\quad + \sum_{i: v_i < -C, w_i > C} (-C - C) \sum_{m=0}^{n-1} (-C)^{n-1-m} C^m + \sum_{i: v_i > C, |w_i| \leq C} (C - w_i) \sum_{m=0}^{n-1} C^{n-1-m} w_i^m \\ &\quad \left. + \sum_{i: v_i > C, w_i < -C} [C - (-C)] \sum_{m=0}^{n-1} C^{n-1-m} (-C)^m \right| \\ &\leq nC^{n-1} \left( \sum_{i: |v_i| \leq C, |w_i| \leq C} |v_i - w_i| + \sum_{i: |v_i| \leq C, w_i < -C} |v_i - (-C)| + \sum_{i: |v_i| \leq C, w_i > C} |v_i - C| \right. \\ &\quad + \sum_{i: v_i < -C, |w_i| \leq C} |-C - w_i| + \sum_{i: v_i < -C, w_i > C} |-C - C| + \sum_{i: v_i > C, |w_i| \leq C} |C - w_i| \\ &\quad \left. + \sum_{i: v_i > C, w_i < -C} |C - (-C)| \right) \end{aligned}$$

$$\begin{aligned}
&\leq nC^{m-1} \left( \sum_{i:|v_i|\leq C,|w_i|\leq C} |v_i - w_i| + \sum_{i:|v_i|\leq C,w_i < -C} |v_i - w_i| + \sum_{i:|v_i|\leq C,w_i > C} |v_i - w_i| \right. \\
&\quad \left. + \sum_{i:v_i < -C,|w_i|\leq C} |v_i - w_i| + \sum_{i:v_i < -C,w_i > C} |v_i - w_i| + \sum_{i:v_i > C,|w_i|\leq C} |v_i - w_i| + \sum_{i:v_i > C,w_i < -C} |v_i - w_i| \right) \\
&\leq nC^{m-1} \sum_{i=1}^d |v_i - w_i| \\
&= nC^{m-1} \|\mathbf{v} - \mathbf{w}\|_1.
\end{aligned}$$

Note that Hölder's inequality implies that

$$\|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i| \leq \left( \sum_{j=1}^d |x_j|^2 \right)^{1/2} \left( \sum_{k=1}^d 1 \right)^{1/2} = d^{1/2} \|\mathbf{x}\|_2$$

for any  $\mathbf{x} \in \mathbb{R}^d$ . Thus,

$$|f_b(\mathbf{v}) - f_b(\mathbf{w})| \leq nC^{m-1} d^{1/2} \|\mathbf{v} - \mathbf{w}\|_2$$

Therefore,

$$|f_{n,C}(\mathbf{v}) - f_{n,C}(\mathbf{w})| = \frac{1}{nC^{m-1} d^{1/2}} |f_b(\mathbf{v}) - f_b(\mathbf{w})| \leq \|\mathbf{v} - \mathbf{w}\|_2$$

implying  $f_{n,C}$  is a 1-Lipschitz function. □

**Example of the Unbounded EM Distance:** Given  $C' > 0$ , consider the two distributions as follows:

$$p_l^{(t+\Delta t)}(\mathbf{x}) = \begin{cases} \frac{1}{\int_{[C'/2, C']^d} 1 d\mathbf{x}} & , \mathbf{x} \in [C'/2, C']^d \\ 0 & , \text{other} \end{cases}$$

and

$$p_l^{(t)}(\mathbf{y}) = \begin{cases} \frac{1}{\int_{[0, C'/4]^d} 1 d\mathbf{y}} & , \mathbf{y} \in [0, C'/4]^d \\ 0 & , \text{other} \end{cases}$$

Then, the supports of the distributions are subsets of  $[0, C']^d$  and for any  $i \in \{1, 2, \dots, d\}$ ,

$$\mathbb{E}_{\mathbf{x} \sim p_l^{(t+\Delta t)}} [x_i^n] > \frac{C'^n}{2^n}, \mathbb{E}_{\mathbf{y} \sim p_l^{(t)}} [y_i^n] < \frac{C'^n}{4^n}$$

Thus, according to **Theorem 2**, the lower bound on the EM distance between  $p_l^{(t+\Delta t)}$  and  $p_l^{(t)}$  is

$$\begin{aligned}
W(p_l^{(t+\Delta t)}, p_l^{(t)}) &\geq \frac{1}{nC'^{m-1} d^{1/2}} \left| \mathbb{E}_{\mathbf{x} \sim p_l^{(t+\Delta t)}} \left[ \sum_{i=1}^d x_i^n \right] - \mathbb{E}_{\mathbf{y} \sim p_l^{(t)}} \left[ \sum_{i=1}^d y_i^n \right] \right| \\
&= \frac{1}{nC'^{m-1} d^{1/2}} \left| \sum_{i=1}^d \mathbb{E}_{\mathbf{x} \sim p_l^{(t+\Delta t)}} [x_i^n] - \mathbb{E}_{\mathbf{y} \sim p_l^{(t)}} [y_i^n] \right| \\
&> \frac{(2^{-n} - 4^{-n}) d^{1/2}}{n} C'
\end{aligned}$$

Note that  $n$  and  $d$  are fixed, and then the lower bound is dominated by  $C'$ . Thus, the distance is unbounded and would go to infinity as  $C' \rightarrow \infty$ . Intuitively, the samples from  $p_l^{(t+\Delta t)}$  can be regarded as the scaled and shifted samples from  $p_l^{(t)}$ . The unstable scale and center of the distributions lead to the unbounded lower bound, implying the unbounded distance and upper bound. The distributions of deep layers' outputs might perform in the same way if they are not normalized. In contrast, the upper bound on the distance for the outputs processed by BN is relatively stable (though it depends on the moments with noise), and can constrain the distance to a reasonable range in some cases discussed in the paper.

**Theorem 3.1.** Suppose that for  $\mathbf{x} \sim p_i^{(t)}$ ,  $g(\mathbf{x}) \sim p_U^{(t)}$ . Then,

$$W(p_U^{(t+\Delta t)}, p_U^{(t)}) \leq 2.$$

*Proof.* Note that for any 1-Lipschitz function  $f$  such that  $f(\mathbf{0}) = 0$ ,

$$|f(g(\mathbf{x}))| = |f(g(\mathbf{x})) - f(\mathbf{0})| \leq \|g(\mathbf{x}) - \mathbf{0}\|_2 = 1, \forall \mathbf{x}$$

which yields

$$-1 \leq f(g(\mathbf{x})) \leq 1, \forall \mathbf{x}$$

Therefore, according to **Fact 2**,

$$\begin{aligned} W(p_U^{(t+\Delta t)}, p_U^{(t)}) &= \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(g(\mathbf{x}))] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(g(\mathbf{y}))] \\ &= \sup_{f: \|f\|_L \leq 1, f(\mathbf{0})=0} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(g(\mathbf{x}))] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(g(\mathbf{y}))] \\ &\leq \sup_{f: \|f\|_L \leq 1, f(\mathbf{0})=0} \mathbb{E}_{\mathbf{x} \sim p^{(t+\Delta t)}} [1] - \mathbb{E}_{\mathbf{y} \sim p^{(t)}} [-1] \\ &= 2. \end{aligned}$$

□

**Theorem 3.2.** Suppose that for  $\alpha \in [0, 1]$  and  $\mathbf{x} \sim p_i^{(t)}$ ,  $g(\mathbf{x}; \alpha) \sim p_U^{(t)}$ . Then,

$$W(p_U^{(t+\Delta t)}, p_U^{(t)}) \leq \mathbb{I}_{\alpha=0}(\alpha) \cdot (\mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [\|\mathbf{x}\|_2] + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [\|\mathbf{y}\|_2]) + \mathbb{I}_{\alpha>0}(\alpha) \cdot \frac{2}{\alpha}.$$

*Proof.* Note that given a 1-Lipschitz function  $f$  that satisfies  $f(\mathbf{0}) = 0$ , for any  $\mathbf{x}$  and  $\alpha \in [0, 1]$ ,

$$\begin{aligned} |f(g(\mathbf{x}; \alpha))| &= |f(g(\mathbf{x}; \alpha)) - f(\mathbf{0})| \\ &\leq \|g(\mathbf{x}; \alpha) - \mathbf{0}\|_2 \\ &= \mathbb{I}_{\|\mathbf{x}\|_2 > 0}(\mathbf{x}) \cdot \frac{\|\mathbf{x}\|_2}{\alpha \|\mathbf{x}\|_2 + (1 - \alpha) \times 1} + \mathbb{I}_{\|\mathbf{x}\|_2 = 0, \alpha = 1}(\mathbf{x}, \alpha) \cdot 1 + \mathbb{I}_{\|\mathbf{x}\|_2 = 0, \alpha < 1}(\mathbf{x}, \alpha) \cdot 0 \\ &= \mathbb{I}_{\|\mathbf{x}\|_2 > 0, \alpha = 0}(\mathbf{x}, \alpha) \cdot \frac{\|\mathbf{x}\|_2}{\alpha \|\mathbf{x}\|_2 + (1 - \alpha) \times 1} + \mathbb{I}_{\|\mathbf{x}\|_2 > 0, \alpha > 0}(\mathbf{x}, \alpha) \cdot \frac{1}{\alpha + (1 - \alpha) / \|\mathbf{x}\|_2} + \mathbb{I}_{\|\mathbf{x}\|_2 = 0, \alpha = 1}(\mathbf{x}, \alpha) \cdot 1 \\ &\leq \mathbb{I}_{\|\mathbf{x}\|_2 > 0, \alpha = 0}(\mathbf{x}, \alpha) \cdot \|\mathbf{x}\|_2 + \mathbb{I}_{\|\mathbf{x}\|_2 > 0, \alpha > 0}(\mathbf{x}, \alpha) \cdot \frac{1}{\alpha} + \mathbb{I}_{\|\mathbf{x}\|_2 = 0, \alpha > 0}(\mathbf{x}, \alpha) \cdot 1 \\ &\leq \mathbb{I}_{\alpha=0}(\alpha) \cdot \|\mathbf{x}\|_2 + \mathbb{I}_{\alpha>0}(\alpha) \cdot \frac{1}{\alpha} \end{aligned}$$

which yields

$$-\mathbb{I}_{\alpha=0}(\alpha) \cdot \|\mathbf{x}\|_2 - \mathbb{I}_{\alpha>0}(\alpha) \cdot \frac{1}{\alpha} \leq f(g(\mathbf{x}; \alpha)) \leq \mathbb{I}_{\alpha=0}(\alpha) \cdot \|\mathbf{x}\|_2 + \mathbb{I}_{\alpha>0}(\alpha) \cdot \frac{1}{\alpha}.$$

Therefore, according to **Fact 2**,

$$\begin{aligned} W(p_U^{(t+\Delta t)}, p_U^{(t)}) &= \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(g(\mathbf{x}; \alpha))] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(g(\mathbf{y}; \alpha))] \\ &= \sup_{f: \|f\|_L \leq 1, f(\mathbf{0})=0} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(g(\mathbf{x}; \alpha))] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(g(\mathbf{y}; \alpha))] \\ &\leq \mathbb{I}_{\alpha=0}(\alpha) \cdot (\mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [\|\mathbf{x}\|_2] + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [\|\mathbf{y}\|_2]) + \mathbb{I}_{\alpha>0}(\alpha) \cdot \frac{2}{\alpha}. \end{aligned}$$

□

**Theorem 3.3.** Suppose that for  $\alpha \in [0, 1]^d$  and  $\mathbf{x} \sim p_t^{(t)}$ ,  $g(\mathbf{x}; \alpha) \sim p_U^{(t)}$ . Then,

$$W(p_U^{(t+\Delta t)}, p_U^{(t)}) \leq \mathbb{I}_{\min_j \alpha_j > 0}(\alpha) \cdot \frac{2}{\min_j \alpha_j} + \mathbb{I}_{\min_j \alpha_j = 0}(\alpha) \cdot (\mathbb{E}_{\mathbf{x} \sim p_t^{(t+\Delta)}}[|\mathbf{x}|_2] + \mathbb{E}_{\mathbf{y} \sim p_t^{(t)}}[|\mathbf{y}|_2] + 2).$$

*Proof.* Given a 1-Lipschitz function  $f$  satisfying  $f(\mathbf{0}) = 0$ , for any  $\mathbf{x}$  and  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_d) \in [0, 1]^d$ ,

$$\begin{aligned} |f(g(\mathbf{x}; \alpha))| &\leq \|g(\mathbf{x}; \alpha)\|_2 = \mathbb{I}_{\|\mathbf{x}\|_2=0}(\mathbf{x}) \cdot 0 + \mathbb{I}_{\|\mathbf{x}\|_2>0}(\mathbf{x}) \cdot \left[ \sum_{i=1}^d \left( \frac{x_i}{\alpha_i \|\mathbf{x}\|_2 + (1 - \alpha_i) \cdot 1} \right)^2 \right]^{\frac{1}{2}} \\ &\leq \mathbb{I}_{\|\mathbf{x}\|_2>0}(\mathbf{x}) \cdot \left[ \sum_{i=1}^d \left( \frac{x_i}{\min_j \alpha_j \|\mathbf{x}\|_2 + (1 - \alpha_j)} \right)^2 \right]^{\frac{1}{2}} \\ &\leq \mathbb{I}_{\|\mathbf{x}\|_2>0}(\mathbf{x}) \cdot \frac{\|\mathbf{x}\|_2}{\min_j \alpha_j \|\mathbf{x}\|_2 + (1 - \alpha_j)} \\ &= \mathbb{I}_{\|\mathbf{x}\|_2>0}(\mathbf{x}) \cdot \left( \mathbb{I}_{\min_j \alpha_j > 0}(\alpha) \cdot \frac{1}{\min_j \alpha_j + (1 - \alpha_j)/\|\mathbf{x}\|_2} + \mathbb{I}_{\min_j \alpha_j = 0}(\alpha) \cdot \frac{\|\mathbf{x}\|_2}{\min_j \alpha_j \|\mathbf{x}\|_2 + (1 - \alpha_j)} \right) \\ &\leq \mathbb{I}_{\|\mathbf{x}\|_2>0}(\mathbf{x}) \cdot \left( \mathbb{I}_{\min_j \alpha_j > 0}(\alpha) \cdot \frac{1}{\min_j \alpha_j} + \mathbb{I}_{\min_j \alpha_j = 0}(\alpha) \cdot \max\{\|\mathbf{x}\|_2, 1\} \right) \\ &\leq \mathbb{I}_{\min_j \alpha_j > 0}(\alpha) \cdot \frac{1}{\min_j \alpha_j} + \mathbb{I}_{\min_j \alpha_j = 0}(\alpha) \cdot \max\{\|\mathbf{x}\|_2, 1\}. \end{aligned}$$

which yields

$$\begin{aligned} f(g(\mathbf{x}; \alpha)) &\geq -\mathbb{I}_{\min_j \alpha_j > 0}(\alpha) \cdot \frac{1}{\min_j \alpha_j} - \mathbb{I}_{\min_j \alpha_j = 0}(\alpha) \cdot \max\{\|\mathbf{x}\|_2, 1\} \\ f(g(\mathbf{x}; \alpha)) &\leq \mathbb{I}_{\min_j \alpha_j > 0}(\alpha) \cdot \frac{1}{\min_j \alpha_j} + \mathbb{I}_{\min_j \alpha_j = 0}(\alpha) \cdot \max\{\|\mathbf{x}\|_2, 1\}. \end{aligned}$$

Therefore, according to **Fact 2**,

$$\begin{aligned} W(p_U^{(t+\Delta t)}, p_U^{(t)}) &= \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_t^{(t+\Delta)}}[f(g(\mathbf{x}; \alpha))] - \mathbb{E}_{\mathbf{y} \sim p_t^{(t)}}[f(g(\mathbf{y}; \alpha))] \\ &= \sup_{f: \|f\|_L \leq 1, f(\mathbf{0})=0} \mathbb{E}_{\mathbf{x} \sim p_t^{(t+\Delta)}}[f(g(\mathbf{x}; \alpha))] - \mathbb{E}_{\mathbf{y} \sim p_t^{(t)}}[f(g(\mathbf{y}; \alpha))] \\ &\leq \mathbb{I}_{\min_j \alpha_j > 0}(\alpha) \cdot \frac{2}{\min_j \alpha_j} + \mathbb{I}_{\min_j \alpha_j = 0}(\alpha) \cdot \left( \mathbb{E}_{\mathbf{x} \sim p_t^{(t+\Delta)}}[\max\{\|\mathbf{x}\|_2, 1\}] + \mathbb{E}_{\mathbf{y} \sim p_t^{(t)}}[\max\{\|\mathbf{y}\|_2, 1\}] \right) \\ &\leq \mathbb{I}_{\min_j \alpha_j > 0}(\alpha) \cdot \frac{2}{\min_j \alpha_j} + \mathbb{I}_{\min_j \alpha_j = 0}(\alpha) \cdot \left( \mathbb{E}_{\mathbf{x} \sim p_t^{(t+\Delta)}}[|\mathbf{x}|_2] + \mathbb{E}_{\mathbf{x} \sim p_t^{(t+\Delta)}}[1] + \mathbb{E}_{\mathbf{y} \sim p_t^{(t)}}[|\mathbf{y}|_2] + \mathbb{E}_{\mathbf{y} \sim p_t^{(t)}}[1] \right) \\ &\leq \mathbb{I}_{\min_j \alpha_j > 0}(\alpha) \cdot \frac{2}{\min_j \alpha_j} + \mathbb{I}_{\min_j \alpha_j = 0}(\alpha) \cdot (\mathbb{E}_{\mathbf{x} \sim p_t^{(t+\Delta)}}[|\mathbf{x}|_2] + \mathbb{E}_{\mathbf{y} \sim p_t^{(t)}}[|\mathbf{y}|_2] + 2). \end{aligned}$$

□

**Theorem 3.4.** Suppose that for  $\alpha \in \mathbb{R}^d, \epsilon > 0$  and  $\mathbf{x} \sim p_i^{(t)}, g(\mathbf{x}; \alpha, \epsilon) \sim p_g^{(t)}$ . Then,

$$W(p_g^{(t+\Delta t)}, p_g^{(t)}) \leq 2\|\alpha\|_\infty + \|\mathbf{1} - \alpha\|_\infty (\mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}}[\|\mathbf{x}\|_2] + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}}[\|\mathbf{y}\|_2]).$$

*Proof.* Given a 1-Lipschitz function  $f$  satisfying  $f(\mathbf{0}) = 0$ , for any  $\mathbf{x}$  and  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_d) \in \mathbb{R}^d$ ,

$$\begin{aligned} |f(g(\mathbf{x}; \alpha, \epsilon))| &\leq \|g(\mathbf{x}; \alpha, \epsilon)\|_2 \\ &= \left\| \frac{\mathbf{x}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \odot \alpha + (\mathbf{1} - \alpha) \odot \mathbf{x} \right\|_2 \\ &\leq \left\| \frac{\mathbf{x}}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \odot \alpha \right\|_2 + \left\| (\mathbf{1} - \alpha) \odot \mathbf{x} \right\|_2 \\ &= \frac{1}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \left( \sum_{i=1}^d \alpha_i^2 x_i^2 \right)^{1/2} + \left( \sum_{i=1}^d (1 - \alpha_i)^2 x_i^2 \right)^{1/2} \\ &\leq \frac{1}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} \left( \max_{1 \leq i \leq d} \alpha_i^2 \sum_{i=1}^d x_i^2 \right)^{1/2} + \left( \max_{1 \leq i \leq d} (1 - \alpha_i)^2 \sum_{i=1}^d x_i^2 \right)^{1/2} \\ &= \|\alpha\|_\infty \cdot \frac{\|\mathbf{x}\|_2}{\sqrt{\|\mathbf{x}\|_2^2 + \epsilon}} + \|\mathbf{1} - \alpha\|_\infty \|\mathbf{x}\|_2 \\ &\leq \|\alpha\|_\infty + \|\mathbf{1} - \alpha\|_\infty \|\mathbf{x}\|_2 \end{aligned}$$

which yields

$$-\|\alpha\|_\infty - \|\mathbf{1} - \alpha\|_\infty \|\mathbf{x}\|_2 \leq f(g(\mathbf{x}; \alpha, \epsilon)) \leq \|\alpha\|_\infty + \|\mathbf{1} - \alpha\|_\infty \|\mathbf{x}\|_2.$$

Therefore, according to **Fact 2**,

$$\begin{aligned} W(p_g^{(t+\Delta t)}, p_g^{(t)}) &= \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(g(\mathbf{x}; \alpha, \epsilon))] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(g(\mathbf{y}; \alpha, \epsilon))] \\ &= \sup_{f: \|f\|_L \leq 1, f(\mathbf{0})=0} \mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}} [f(g(\mathbf{x}; \alpha, \epsilon))] - \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}} [f(g(\mathbf{y}; \alpha, \epsilon))] \\ &\leq 2\|\alpha\|_\infty + \|\mathbf{1} - \alpha\|_\infty (\mathbb{E}_{\mathbf{x} \sim p_i^{(t+\Delta t)}}[\|\mathbf{x}\|_2] + \mathbb{E}_{\mathbf{y} \sim p_i^{(t)}}[\|\mathbf{y}\|_2]). \end{aligned}$$

□



## 2. Earth Mover Distance Estimation

In the theoretical analysis, the EM distance has been used to measure ICS. However, whether there is a correlation between the performance of networks and ICS in practice is unclear. Thus, an algorithm is proposed to estimate the EM distance in network training, and the results of the following experiments yield evidence that ICS is related to the networks' performance.

### 2.1. Algorithm

By definition, the accurate EM distance is obtained by optimizing a function  $f$  over the 1-Lipschitz function space:

$$W(p_l^{(t+\Delta t)}, p_l^{(t)}) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p_l^{(t+\Delta t)}}[f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim p_l^{(t)}}[f(\mathbf{y})].$$

Thus, the proposed algorithm aims to get the optimal 1-Lipschitz function  $f^*$  and use  $f^*$  to compute the difference between the expectations. In WGAN [1], the 1-Lipschitz function  $f$  is parameterized by a neural network  $f_w$  with a scalar output in the algorithm. Then, the parameters of  $f_w$  are bounded to satisfy the constraint  $\|f_w\|_L \leq 1$ . Furthermore, the constraint can be relaxed such that  $w$  is required to satisfy  $\|f_w\|_L \leq K$  for some  $K > 0$ , and then the EM distance is  $K \cdot W(p_l^{(t+\Delta t)}, p_l^{(t)})$ . For convenience, the parameter space is denoted by  $\mathcal{W} = \{w \mid \|f_w\|_L \leq K\}$ . Then, a suboptimal  $K$ -Lipschitz function  $f_w^*$  is obtained by training  $f_w$  on the dataset consisting of the real and generated images. The EM distance is computed by  $f_w^*$ 's outputs with the real and generated images as the inputs.

In a convolutional neural network  $f_{1:L}$  with  $L$  layers, the EM distance for measuring ICS is estimated in the same way. For the  $l$ th ( $l < L$ ) convolutional layer of  $f_{1:L}$ , the multi-channel outputs can be treated as images, though there are more than 3 channels. Thus, analogously, the proposed algorithm trains  $f_w$  to maximize the same objective function of WGANs, with the training samples generated by the local networks  $f_{1:l}^{(t)}$  and  $f_{1:l}^{(t+\Delta t)}$  formed by the previous  $l$  layers of  $f_{1:L}$  at the  $t$ th and  $t + \Delta t$ th iterations, respectively. Then, the algorithm estimates the EM distance by  $f_w^*$  similarly.

To formulate the process, the approximated EM distance, which replaces the 1-Lipschitz function space with  $\mathcal{W}$ , is defined as

$$\widetilde{W}(p_l^{(t+\Delta t)}, p_l^{(t)}) = \sup_{w \in \mathcal{W}} E_{\mathbf{x} \sim p_l^{(t+\Delta t)}}[f_w(\mathbf{x})] - E_{\mathbf{y} \sim p_l^{(t)}}[f_w(\mathbf{y})].$$

Then, in practice,  $\widetilde{W}(p_l^{(t+\Delta t)}, p_l^{(t)})$  is further approximated by the empirical estimation:

$$\widetilde{W}(p_l^{(t+\Delta t)}, p_l^{(t)}) \approx \sup_{w \in \mathcal{W}} \frac{1}{N} \sum_{i=1}^N f_w(f_{1:l}^{(t+\Delta t)}(\mathbf{x}_i)) - \frac{1}{N} \sum_{i=1}^N f_w(f_{1:l}^{(t)}(\mathbf{x}_i)),$$

where  $N$  is the number of samples and  $\mathbf{x}_i$  is a sample from the dataset for  $f_{1:L}$ . The algorithm is presented in Algorithm 3.

---

#### Algorithm 3 EM Distance Estimation Algorithm

---

**Input:** datasets  $\mathcal{D}_{train}$  and  $\mathcal{D}_{test}$ , local networks  $f_{1:l}^{(t+\Delta t)}$  and  $f_{1:l}^{(t)}$ , initialized  $f_w$ , parameters  $T$ ,  $n$  and  $c$

**Output:** estimated EM distance  $d$

- 1: **for**  $i \leftarrow 1$  to  $T$  **do**
  - 2:   Sample  $\{\mathbf{x}_j\}_{j=1}^n$  randomly from  $\mathcal{D}_{train}$
  - 3:    $g_w \leftarrow \nabla_w \left[ \frac{1}{n} \sum_{j=1}^n (f_w(f_{1:l}^{(t+\Delta t)}(\mathbf{x}_j)) - f_w(f_{1:l}^{(t)}(\mathbf{x}_j))) \right]$
  - 4:   Update  $w$  by  $g_w$
  - 5:    $w \leftarrow clip(w, -c, c)$
  - 6: **end for**
  - 7:  $d \leftarrow \frac{1}{|\mathcal{D}_{test}|} \sum_{\mathbf{x} \in \mathcal{D}_{test}} (f_w(f_{1:l}^{(t+\Delta t)}(\mathbf{x})) - f_w(f_{1:l}^{(t)}(\mathbf{x})))$
- 

### 2.2. Experiments

The EM distances for some specific unitization layers of ResNet-110s are estimated. The ResNet-110s, *i.e.*,  $f_{1:L}$ s, are trained on CIFAR-10 and CIFAR-100 datasets with the data augmentation method [4], while the  $f_w$ s corresponding to the specific layers are trained on the same dataset without data augmentation.

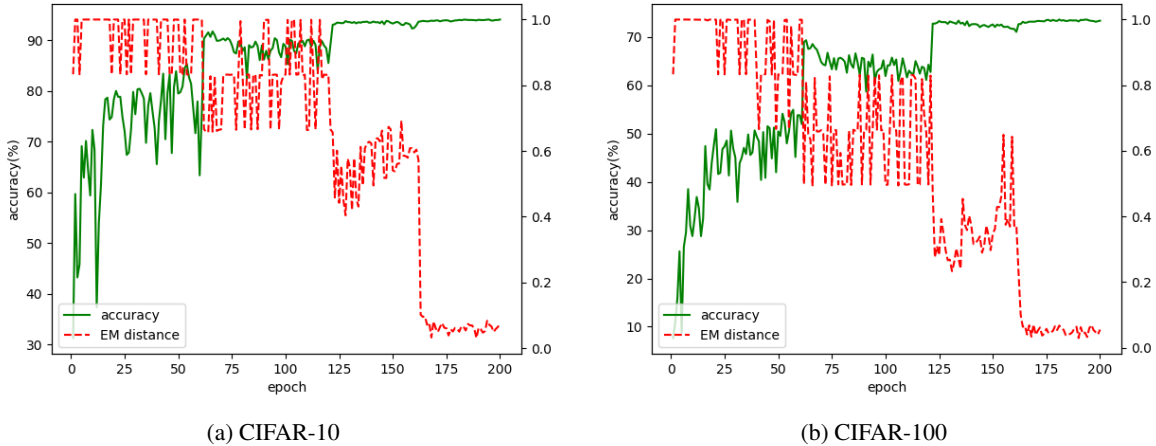


Figure 1: Illustration of test accuracy of ResNet-110s on the CIFAR datasets and the average EM distances. On both the CIFAR-10 (a) and CIFAR-100 (b) datasets, the accuracy increases significantly in about 60th, 120th and 160th epochs, and the average EM distance sharply drops concurrently.

**Network Architectures:** Each  $f_w$  is parameterized as a vanilla CNN. The inputs of  $f_w$  are the outputs of the corresponding layer, followed by four  $3 \times 3$  convolutions on the feature maps of sizes  $\{64, 64, 128, 128\}$  with strides of  $\{2, 1, 2, 1\}$ , respectively. The convolutional outputs are activated using ReLU but not normalized/unitized. A fully-connected layer with a sigmoid activation is applied to generate the network’s outputs.

**Implementation Details:** The experiments of training ResNet-110s on the CIFAR datasets have been described in the paper, and the experiments of estimating the distances are described as follows. The EM distances are computed per epoch. At the end of each epoch, the weights of the current local network  $f_{1:l}^{(t+\Delta t)}$  are first saved for the next epoch. Then another local network  $f_{1:l}^{(t)}$  is initialized by the weights saved in the previous epoch.  $f_w$  is trained by Algorithm 3 with these two local networks as the algorithm’s inputs. For each estimated EM distance, the same hyperparameters are used: the iteration number  $T$  is  $1.5k$ ; the mini-batch size  $n$  is 128; the bound  $c$  is 0.01.  $f_w$  is initialized by the same weights that are generated using the method [2]. Finally, the average EM distance for the deep layers including the 60th, 69th, 78th, 87th, 96th and 105th layers, along with the test accuracy of the ResNet-110s, are reported.

**Results:** As is shown in Figure 1, the accuracy of each experiment significantly increases in about 60th, 120th and 160th epoch as the learning rate decreases. Meanwhile, the average EM distance dramatically drops in these epochs. The results demonstrate that the EM distance is a suitable ICS measure since it drops immediately as the learning rate decreases, in which ICS obviously decreases. Then, based on the effectiveness of the EM distance in measuring ICS, the results further substantiate that ICS is related to the performance of networks.

### 3. Experiments of the Unitization for Micro-Batches

The unitization is evaluated in the case of micro-batches to further demonstrate the performance. For comparison, BN and GN [5] techniques are also evaluated in this case. Only the CIFAR-10 dataset with the same data augmentation in the previous experiments are used.

**Network Architectures:** The unitization, BN and GN are evaluated with the same ResNet-110s in the previous experiments. For GN, there are different ways of group division, according to the experiments in [5]. For convenience, denote by  $G$  the number of the groups in GN, with the value in  $\{1, 2, 4, 8, 16\}$ . Furthermore, for a network, there are two types of settings for  $G$  in each convolutional layer: fixing the numbers of (1) the groups or (2) the channels per group. Thus, there are 10 ways of group division in total.

**Implementation Details:** The ResNet-110s are trained with the same settings as the previous experiments, except for the

batch size, denoted by  $s$ . Unlike [5], we experiment on **only one GPU**, and report the median accuracy of 3 runs for each experiment of  $s = 2$ . In fact, for the cases of  $s \in \{4, 8, 16\}$ , BN still outperforms GN. Thus, only the case of  $s = 2$ , in which BN degrades, is considered.

Table 1: Classification accuracy (%) on the CIFAR-10 dataset

Groups ( $G$ )					Channels per group					BatchNorm	Unitization
1	2	4	8	16	1	2	4	8	16		
10.00	10.00	10.00	71.63	73.76	79.54	73.58	70.68	72.04	10.00	71.50	<b>81.88</b>

**Results:** As is shown in Table 1, the unitization still outperforms the other techniques for micro-batches. Some networks with GN cannot even converge, where the accuracy (of only 10%) degrades, and BN suffers from highly noisy estimations in this case, with the almost lowest accuracy among all the comparable results ( $> 10\%$ ). The results further demonstrate the effectiveness of the unitization in controlling ICS.

## References

[1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pages 214–223, 2017.

[2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.

[3] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

[4] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial Intelligence and Statistics*, pages 562–570, 2015.

[5] Yuxin Wu and Kaiming He. Group normalization. *arXiv preprint arXiv:1803.08494*, 2018.