

# Fine-grained Optimization of Deep Neural Networks

Mete Ozay  
meteozy@gmail.com

## Problems

We conjecture that if we can impose multiple constraints on weights of DNNs to set upper bounds of the norms of the weight matrices, and train the DNNs with these weights, then the DNNs can achieve empirical generalization errors closer to the proposed theoretical bounds. We pose two problems in order to achieve this goal;

1. Renormalization of weights to upper bound norms of their matrices.
2. Training DNNs with renormalized weights with assurance to convergence to minima.

## Weight Manifolds in DNNs

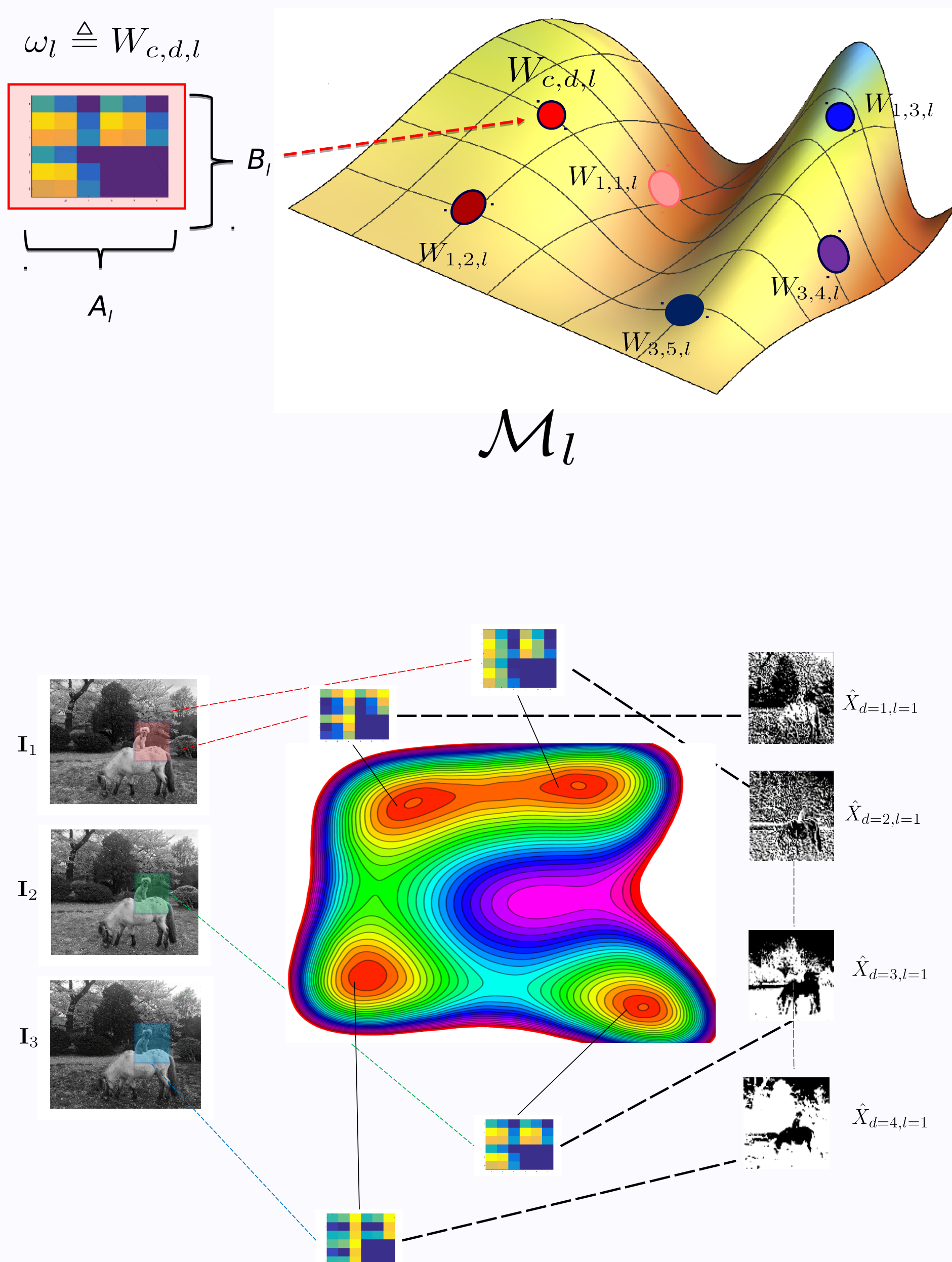
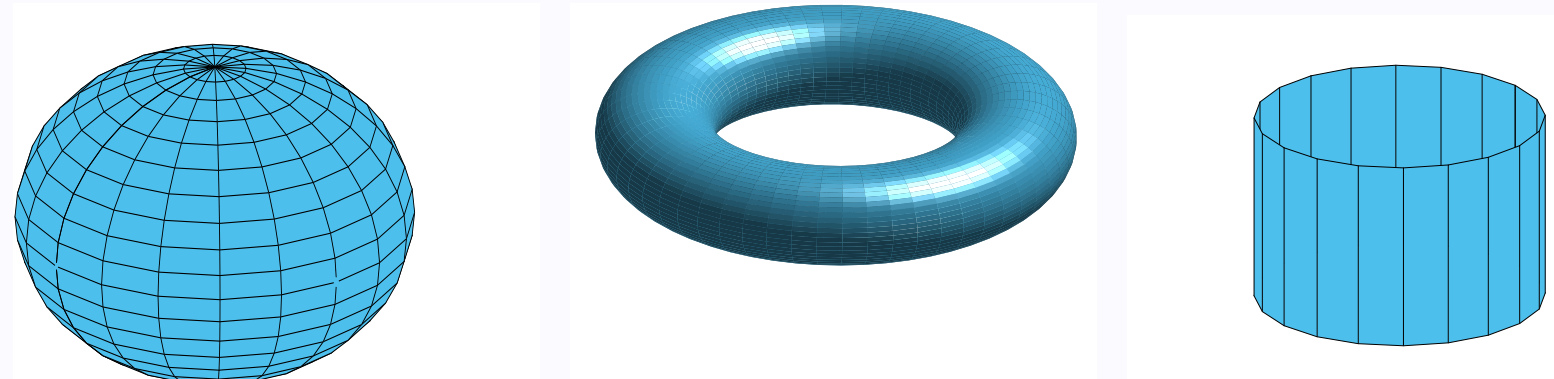


Table 1: Weights  $\omega_{g,l}^i \in \mathbb{R}^{A_l \times B_l}$  belonging to the  $g^{th}$  group of size  $|g|$ ,  $g = 1, 2, \dots, G_l$ ,  $\forall l$  have the same size  $A_l \times B_l$  for simplicity, and  $\sigma(\omega_{g,l}^i)$  denotes the top singular value of  $\omega_{g,l}^i$ .  $\|\omega_{g,l}^i\|_F$ ,  $\|\omega_{g,l}^i\|_2$ , and  $\|\omega_{g,l}^i\|_{2 \rightarrow 1}$ , denotes respectively the Frobenius, spectral and  $\ell_{2 \rightarrow 1}$  norms of the weight  $\omega_{g,l}^i$ .

Norms	Sphere	Stiefel	Oblique
$\ \omega_{g,l}^i\ _2$	$\sigma(\omega_{g,l}^i)$	1.0	$\sigma(\omega_{g,l}^i)$
$\ \omega_{g,l}^i\ _F$	1.0	$(B_l)^{1/2}$	$(B_l)^{1/2}$
$\ \omega_{g,l}^i\ _{2 \rightarrow 1}$	1.0	$(B_l)^{1/4}$	$(B_l)^{1/4}$



(a) An orthonormalized weight  $\omega \in \mathbb{R}^{3 \times 1}$  ( $\omega \in \mathbb{R}^{A \times B}$ ) resides on a two-sphere  $\mathbb{S}^2$  ( $\mathbb{S}^{AB-1}$ ) which has constant positive sectional curvature, 1. (b) A weight  $\omega = (\omega_1, \omega_2)$ , where each  $\omega_i \in \mathbb{R}^{2 \times 1}$ ,  $i = 1, 2$ , belongs to a circle  $\mathbb{S}^1$ , resides on a two-torus  $\mathbb{T}^2$  with varying curvature. (c) If  $\omega_1 \in \mathbb{S}^1$  ( $\omega_1 \in \mathbb{S}^p$ ) and  $\omega_2 \in \mathbb{R}$  ( $\omega_2 \in \mathbb{R}^{q-p}$ ), then  $\omega$  resides on a cylinder  $\mathbb{S}^1 \times \mathbb{R}$  with varying curvature.

## Products of Weight Manifolds

**Definition 1**  $\mathcal{G}_l = \{\mathcal{M}_{\iota,l} : \iota \in \mathcal{I}_{\mathcal{G}_l}\}$  is a set of **weight manifolds**  $\mathcal{M}_{\iota,l}$  of dimension  $n_{\iota,l}$ , which is identified by a set of indices  $\mathcal{I}_{\mathcal{G}_l}$ ,  $\forall l = 1, 2, \dots, L$ .

$\mathcal{I}_{\mathcal{G}_l}$  contains indices each of which represents an identity number ( $\iota$ ) of a weight that resides on a manifold  $\mathcal{M}_{\iota,l}$  at the  $l^{th}$  layer.

A subset  $\mathcal{I}_l^g \subseteq \mathcal{I}_{\mathcal{G}_l}$ ,  $g = 1, 2, \dots, G_l$ , is used to determine a subset  $\mathcal{G}_l^g \subseteq \mathcal{G}_l$  of weight manifolds which will be aggregated to construct a **product of weight manifolds (POM)**.

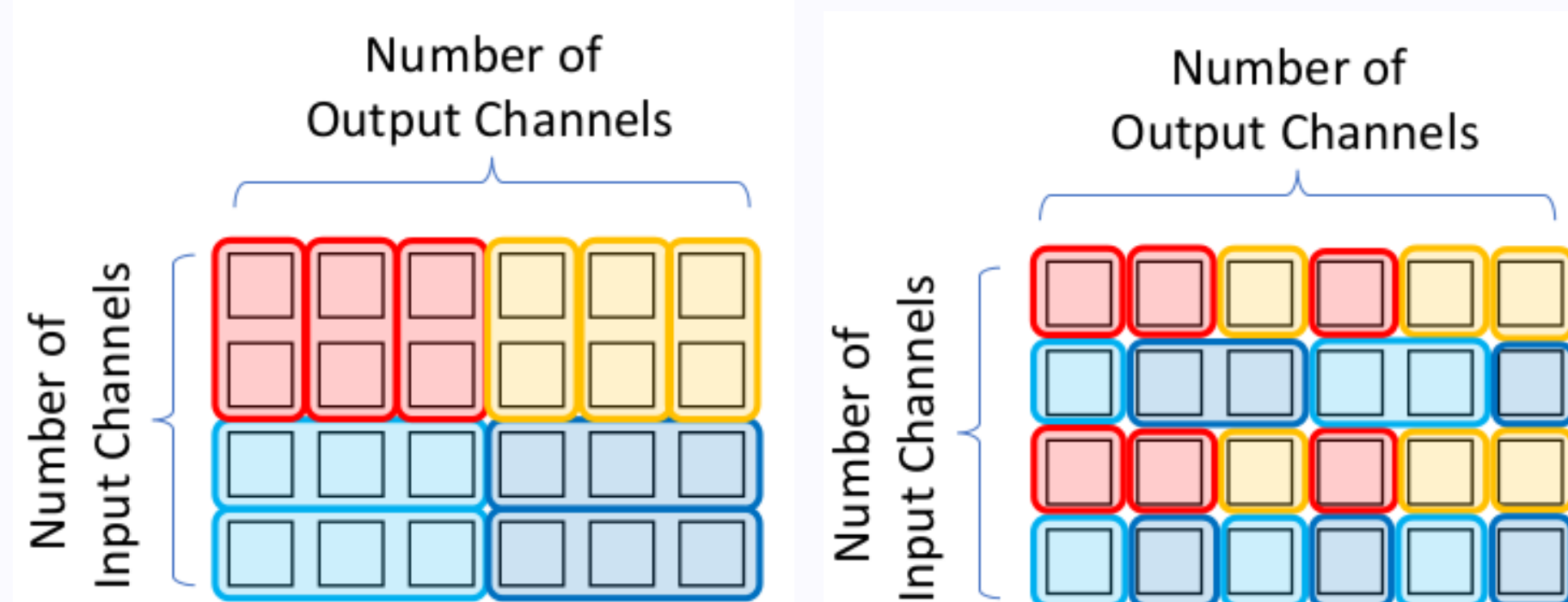
Each  $\mathcal{M}_{\iota,l} \in \mathcal{G}_l^g$  is called a **component manifold** of a product of weight manifolds which is denoted by  $\mathbb{M}_{g,l}$ . A weight  $\omega_{g,l} \in \mathbb{M}_{g,l}$  is obtained by concatenating weights belonging to  $\mathcal{M}_{\iota,l}$ ,  $\forall \iota \in \mathcal{I}_l^g$ , using  $\omega_{g,l} = (\omega_1, \omega_2, \dots, \omega_{|\mathcal{I}_l^g|})$ , where  $|\mathcal{I}_l^g|$  is the cardinality of  $\mathcal{I}_l^g$ . A  $\mathcal{G}_l$  is called a **collection of POMs**.

**Example 1** We have a weight tensor of size  $3 \times 3 \times 4 \times 6$  where the number of input and output channels is 4 and 6. In total, we have  $4 \times 6 = 24$  weight matrices of size  $3 \times 3$ .

**POMs for input and output channels (PIO):** We split the set of 24 weights into 10 subsets. For 6 output channels, we split the set of weights corresponding to 4 input channels into 3 subsets.

We choose the sphere (Sp) for 2 subsets each containing 3 weights (depicted by light blue rectangles), and 3 subsets each containing 2 weights (depicted by red rectangles).

We choose the Stiefel manifold (St) similarly for the remaining subsets. Then, our ensemble contains 5 POMs of St and 5 POMs of Sp.



**Deterministic and random selection of channels:** In the experiments, indices of the sets are selected randomly using a hypergeometric distribution without replacement at the initialization of a training step, and fixed in the rest of the training.

## Geometry of POMs

**Computation of metrics:** A metric defined on a product weight manifold  $\mathbb{M}_{g,l}$  can be computed by superposition (i.e. linear combination) of Riemannian metrics of its component manifolds.

**Lower bounds of sectional curvatures:** Sectional curvature of a product weight manifold  $\mathbb{M}_{g,l}$  is lower bounded by 0.

## Optimization using FG-SGD

**input** :  $T$  (number of iterations),  
 $S$  (training set),  
 $\Theta$  (set of hyperparameters),  
 $\mathcal{L}$  (a loss function),

$\mathcal{I}_g^l \subseteq \mathcal{I}_{\mathcal{G}_l}$ ,  $\forall g, l$ .

**output:** A set of weights  $\{\omega_{g,l}^T\}_{l=1}^L, \forall g$ .

- 1 **Initialization:** Construct a collection of products of weight manifolds  $\mathcal{G}_l$ , initialize re-scaling parameters  $\mathcal{R}_l^t$  and initialize weights  $\omega_{g,l}^t \in \mathbb{M}_{g,l}$  with  $\mathcal{I}_g^l \subseteq \mathcal{I}_{\mathcal{G}_l}$ ,  $\forall m, l$ .
- 2 **for each iteration**  $t = 1, 2, \dots, T$  **do**
- 3     **for each layer**  $l = 1, 2, \dots, L$  **do**
- 4          $\mathcal{L}(\omega_{g,l}^t) := \Pi_{\omega_{g,l}^t} \left( \mathcal{L}(\omega_{g,l}^t), \Theta, \mathcal{R}_l^t \right), \forall \mathcal{G}_l$ .
- 5          $v_t := h(\mathcal{L}(\omega_{g,l}^t), r(t, \Theta)), \forall \mathcal{G}_l$ .
- 6          $\omega_{g,l}^{t+1} := \phi_{\omega_{g,l}^t}(v_t, \mathcal{R}_l^t), \forall \omega_{g,l}^t, \forall \mathcal{G}_l$ .
- 7     **end for**
- 8 **end for**

## Generalization Bounds

Norms of concatenated weights  $\omega_{g,l}, \forall g$ , are lower bounded by products of norms of component weights  $\omega_{g,l}^i, \forall i$ . Weights are rescaled at each  $t^{th}$  epoch using  $\mathcal{R}_{i,l}^t = \frac{\gamma_{i,l}}{\lambda_{i,l}^t}$ , where  $\gamma_{i,l}$  is a geometric scaling parameter and  $\lambda_{i,l}^t$  is the standard deviation of features input to the  $i^{th}$  weight in the  $g^{th}$  group  $\omega_{g,l}^i, \forall i, g$ .

Table 2: **Comparison of generalization bounds.**  $\delta_{l,F}$ ,  $\delta_{l,2}$ , and  $\delta_{l,2 \rightarrow 1}$  denotes upper bounds of the Frobenius norm  $\|\omega_l\|_F \leq \delta_{l,F}$ , spectral norm  $\|\omega_l\|_2 \leq \delta_{l,2}$  and the sum of the Euclidean norms for all rows  $\|\omega_l\|_{2 \rightarrow 1} \leq \delta_{l,2 \rightarrow 1}$  ( $\ell_{2 \rightarrow 1}$ ) of weights  $\omega_l$  at the  $l^{th}$  layer of an  $L$  layer DNN using  $N$  samples.

DNNs (dynamic group scaling)
$\mathcal{O}\left(\frac{2^L \prod_{l=1}^L \prod_{g=1}^{G_l} \delta_{g,l,F}}{\sqrt{N}}\right)$
$\tilde{\mathcal{O}}\left(\frac{\prod_{l=1}^L \prod_{g=1}^{G_l} \delta_{g,l,2}}{\sqrt{N}} \left(\sum_{l=1}^L \prod_{g=1}^{G_l} \left(\frac{\delta_{g,l,2 \rightarrow 1}}{\delta_{g,l,2}}\right)^{\frac{2}{3}}\right)^{\frac{3}{2}}\right)$
$\tilde{\mathcal{O}}\left(\frac{\prod_{l=1}^L \prod_{g=1}^{G_l} \delta_{g,l,2}}{\sqrt{N}} \sqrt{L^2 \varpi \sum_{l=1}^L \prod_{g=1}^{G_l} \frac{\delta_{g,l,F}^2}{\delta_{g,l,2}^2}}\right)$

## Convergence of FG-SGD

**Convergence to local minima:** The loss function of a non-linear DNN converges to a local minimum, and the corresponding gradient converges to zero almost surely (a.s.).

**Convergence to global minima:** Loss functions of particular nonlinear DNNs (e.g. pyramidal), converge to a global minimum a.s.

## Experimental Analysis

	Imagenet(Resnet-101)	Imagenet(SENNet-Resnet-101)
Euc.	23.15 $\pm$ 0.09	22.38 $\pm$ 0.30
PIO (Sp+Ob+St)	22.83 $\pm$ 0.06	21.93 $\pm$ 0.12
PIO (Sp+Ob+St+Euc.)	22.75 $\pm$ 0.02	21.76 $\pm$ 0.09