

Norm-Preservation: Why Residual Networks Can Become Extremely Deep?

Alireza Zaeemzadeh, *Student Member, IEEE*, Nazanin Rahnavard, *Senior Member, IEEE*,
and Mubarak Shah, *Fellow, IEEE*

Abstract—Augmenting neural networks with skip connections, as introduced in the so-called ResNet architecture, surprised the community by enabling the training of networks of more than 1,000 layers with significant performance gains. This paper deciphers ResNet by analyzing the effect of skip connections, and puts forward new theoretical results on the advantages of identity skip connections in neural networks. We prove that the skip connections in the residual blocks facilitate preserving the norm of the gradient, and lead to stable back-propagation, which is desirable from optimization perspective. We also show that, perhaps surprisingly, as more residual blocks are stacked, the norm-preservation of the network is enhanced. Our theoretical arguments are supported by extensive empirical evidence.

Can we push for extra norm-preservation? We answer this question by proposing an efficient method to regularize the singular values of the convolution operator and making the ResNet's transition layers extra norm-preserving. Our numerical investigations demonstrate that the learning dynamics and the classification performance of ResNet can be improved by making it even more norm preserving. Our results and the introduced modification for ResNet, referred to as Procrustes ResNets, can be used as a guide for training deeper networks and can also inspire new deeper architectures.

Index Terms—Residual Networks, Convolutional Neural Networks, Optimization Stability, Norm Preservation, Spectral Regularization.



1 INTRODUCTION

Deep neural networks have progressed rapidly during the last few years, achieving outstanding, sometimes super human, performance [1]. It is known that the depth of the network, i.e., number of stacked layers, is of decisive significance. It is shown that as the networks become deeper, they are capable of representing more complex mappings [2]. However, deeper networks are notoriously harder to train. As the number of layers is increased, optimization issues arise and, in particular, avoiding vanishing/exploding gradients is essential to optimization stability of such networks. Batch normalization, regularization, and initialization techniques have shown to be useful remedies for this problem [3], [4].

Furthermore, it has been observed that as the networks become increasingly deep, the performance gets saturated or even deteriorates [5]. This problem has been addressed by many recent network designs [5], [6], [7], [8]. All of these approaches use the same design principle: skip connections. This simple trick makes the information flow across the layers easier, by bypassing the activations from one layer to the next using skip connections. Highway Networks [7], ResNets [5], [6], and DenseNets [8] have consistently achieved state-of-the-art performances by using skip connections in different network topologies. The main goal of skip connection is to enable the information to flow through many layers without attenuation. In all of these efforts, it is observed empirically that it is crucial to keep the information path *clean* by using identity mapping in the skip connection. It is also observed that more complicated transformations in the

skip connection lead to more difficulty in optimization, even though such transformations have more representational capabilities [6]. This observation implies that *identity* skip connection, while provides adequate representational ability, has a great feature of optimization stability, enabling deeper well-behaved networks.

Since the introduction of Residual Networks (ResNets) [5], [6], there have been some efforts on understanding how the residual blocks may help the optimization process and how they improve the representational ability of the networks. Authors in [9] showed that skip connection eliminates the singularities caused by the model non-identifiability. This makes the optimization of deeper networks feasible and faster. Similarly, to understand the optimization landscape of ResNets, authors in [10] prove that linear residual networks have no critical points other than the global minimum. This is in contrast to plain linear networks, in which other critical points may exist [11]. Furthermore, authors in [12] show that as depth increases, gradients of plain networks resemble white noise and become less correlated. This phenomenon, which is referred to as *shattered gradient* problem, makes training more difficult. Then, it is demonstrated that residual networks reduce shattering, compared to plain networks, leading to numerical stability and easier optimization.

In this paper, we present and analytically study another desirable effect of identity skip connection: *the norm preservation of error gradient*, as it propagates in the backward path. We show theoretically and empirically that each residual block in ResNets is *increasingly norm-preserving*, as the network becomes *deeper*. This interesting result is in contrast to hypothesis provided in [13], which states that residual networks avoid vanishing gradient *solely* by shortening the effective path of the gradient.

The authors are with the School of Electrical Engineering and Computer Science, University of Central Florida, Orlando, FL 32816 USA (e-mail: zaeemzadeh@eecs.ucf.edu; nazanin@eecs.ucf.edu; shah@crco.ucf.edu).

Furthermore, we show that identity skip connection enforces the norm-preservation during the training, leading to well-conditioning and easier training. This is in contrast to the initialization techniques, in which the initialization distribution is modified to make the training easier [3], [14]. This is done by keeping the variance of weights gradient the same across layers. However, as observed in [14] and verified by our experiments, using such initialization methods, although the network is initially fairly norm-preserving, the norms of the gradients diverge as training progresses.

We analyze the role of identity mapping as skip connection in the ResNet architecture from a theoretical perspective. Moreover, we use the insight gained from our theoretical analysis to propose modifications to some of the building blocks of the ResNet architecture. Two main contributions of this paper are as follows.

- **Proof of the Norm Preservation of ResNets:** We show that having identity mapping in the shortcut path leads to norm-preserving building blocks. Specifically, identity mapping shifts all the singular values of the transformations towards 1. This makes the optimization of the network much easier by preserving the magnitude of the gradient across the layers. Furthermore, we show that, perhaps surprisingly, *as the network becomes deeper, its building blocks become more norm-preserving*. Hence, the gradients can flow smoothly through very deep networks, making it possible to train such networks. Our experiments validate our theoretical findings.
- **Enhancing Norm Preservation:** Using insights from our theoretical investigation, we propose important modifications to the *transition blocks* in the ResNet architecture. The transition blocks are used to change the number of channels and feature map size of the activations. Since these blocks do not use identity mapping as the skip connection, in general, they do not preserve the norm of the gradient. We propose to change the dimension of the activations in a norm preserving manner, such that the network becomes even more norm-preserving. For that, we propose a computationally efficient method to set the nonzero singular values of the convolution operator, without using singular value decomposition. We refer to the proposed architecture as Procrustes ResNet (ProcResNet). Our experiments demonstrate that the proposed norm-preserving blocks are able to improve the optimization stability and performance of ResNets.

The rest of the paper is organized as follows. In Section 2, the theoretical results and the bounds for norm-preservation of linear and nonlinear residual networks are presented. Then, in Section 3, we show how to enhance the norm preservation of the residual networks by introducing a new computationally efficient regularization of convolutions. To verify our theoretical investigation and to demonstrate the effectiveness of the proposed regularization, we provide our experiments in Section 4. Finally, Section 5 draws conclusions.

2 NORM-PRESERVATION OF RESIDUAL NETWORKS

Our following main theorem states that, under certain conditions, a deep residual network representing a nonlinear mapping is norm-preserving in the backward path. We show that, at each residual block, the norm of the gradient with respect to the input is close to the norm of gradient with respect to the output. In other words, *the residual block with identity mapping, as the skip connection, preserves the norm of the gradient in the backward path*. This results in several useful characteristics such as avoiding vanishing/exploding gradient, stable optimization, and performance gain.

Suppose we want to represent a nonlinear mapping $\mathcal{F} : \mathbb{R}^N \rightarrow \mathbb{R}^N$ with a sequence of L non-linear residual blocks of form:

$$\mathbf{x}_{l+1} = \mathbf{x}_l + F_l(\mathbf{x}_l). \quad (1)$$

As illustrated in Figure 1(b), \mathbf{x}_l and \mathbf{x}_{l+1} represent respectively the input and output at l^{th} layer. $F_l(\mathbf{x}_l)$ is the residual transformation learned by the l^{th} layer. Before presenting the theorem, we lay out the following assumptions on \mathcal{F} .

Assumption 1. *The function $\mathcal{F} : \mathbb{R}^N \rightarrow \mathbb{R}^N$ is differentiable, invertible, and satisfies the following conditions:*

- (i) $\forall \mathbf{x}, \mathbf{y}, \mathbf{z}$ with bounded norm, $\exists \alpha > 0$, $\|(\mathcal{F}'(\mathbf{x}) - \mathcal{F}'(\mathbf{y}))\mathbf{z}\| \leq \alpha \|\mathbf{x} - \mathbf{y}\| \|\mathbf{z}\|$,
- (ii) $\forall \mathbf{x}, \mathbf{y}$ with bounded norm, $\exists \beta > 0$, $\|\mathcal{F}^{-1}(\mathbf{x}) - \mathcal{F}^{-1}(\mathbf{y})\| \leq \beta \|\mathbf{x} - \mathbf{y}\|$, and
- (iii) $\exists \mathbf{x}$ with bounded norm such that $\text{Det}(\mathcal{F}'(\mathbf{x})) > 0$,

α and β are constants, independent of network size and architecture. Also, we assume that the domain of inputs is bounded. By rescaling inputs, we can assume, without loss of generality, that $\|\mathbf{x}_1\|_2 \leq 1$ for any input \mathbf{x}_1 .

We would like to emphasize the point that these assumptions are on the mapping that we are trying to represent by the network, not the network itself. Thus, assumptions are independent of architecture. Assumptions (i) and (ii) mean that the function \mathcal{F} is smooth, Lipschitz continuous, and its inverse is differentiable. The practical relevance of invertibility assumption is justified by the success of reversible networks [15], [16], [17]. Reversible architectures look for the true mapping \mathcal{F} *only* in the space of invertible functions and it is shown that imposing such strict assumption on the architecture does not hurt its representation ability [16]. Thus, the mapping \mathcal{F} is either invertible or can be well approximated by an invertible function, in many scenarios. However, unlike the reversible architectures, we do not assume residual blocks or the residual transformations, $F_l(\cdot)$, are invertible, which makes the assumption less strict. Furthermore, our extensive experiments in Section 4 show that our theoretical analysis, which is based on these assumptions, hold. This is further empirical justification that these assumptions are relevant in practice. Finally, Assumption (iii) is without loss of generality [10], [18].

Theorem 1. *Suppose we want to represent a nonlinear mapping $\mathcal{F} : \mathbb{R}^N \rightarrow \mathbb{R}^N$, satisfying Assumption 1, with a sequence of L non-linear residual blocks of form $\mathbf{x}_{l+1} = \mathbf{x}_l + F_l(\mathbf{x}_l)$. There exists a solution such that for all residual blocks we have:*

$$(1 - \delta) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2 \leq \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_l} \right\|_2 \leq (1 + \delta) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2,$$

where $\delta = c \frac{\log(2L)}{L}$, $\mathcal{E}(\cdot)$ is the cost function, and $c = c_1 \max\{\alpha\beta(1 + \beta), \beta(2 + \alpha) + \alpha\}$ for some $c_1 > 0$. α and β are constants defined in Assumption 1.

Proof. See Section A.1. \square

This theorem shows that the mapping \mathcal{F} can be represented by a sequence of L non-linear residual blocks, such that the norm of the gradient does not change significantly, as it is backpropagated through the layers. *One interesting implication of Theorem 1 is that as L , the number of layers, increases, δ becomes smaller and the solution becomes more norm-preserving.* This is a very desirable feature because vanishing or exploding gradient often occurs in deeper network architectures. However, by utilizing residual blocks, as more blocks are stacked, the solution becomes extra norm-preserving.

Now that we proved such a solution exists, we show why residual networks can remain norm preserving throughout the training. For that, we consider the case where $F_i(\mathbf{x}_i)$ consists of two layers of convolution and nonlinearity. The following corollary shows the bound on norm preservation of the residual block depends on the norm of the weights. Therefore, if we bound the optimizer to search only in the space of functions with small norms, we can ensure that the network will remain norm preserving throughout the training. Therefore, any critical point in this space is also norm-preserving. On the other hand, based on Theorem 1, we know that at least one norm preserving solution exists. It is also known that, under certain conditions, any critical point achieved during optimization of ResNets is a global minimizer, meaning that it achieves the same loss function value as the global minimum [10], [18], [19]. Thus, this result implies that ResNets are able to maintain norm-preservation throughout the training and if they converge, the solution is a norm-preserving global minimizer. The conclusions of the corollary can be easily generalized for residual block with more than two layers.

Corollary 1. *Suppose a network contains non-linear residual blocks of form $\mathbf{x}_{l+1} = \mathbf{x}_l + \mathbf{W}_l^{(2)} \rho(\mathbf{W}_l^{(1)} \rho(\mathbf{x}_l))$, where $\rho(\cdot)$ is an element-wise non-linear operator with bounded derivative, i.e., $0 \leq \frac{\partial \rho_n(\mathbf{x}_l)}{\partial x_{l,n}} \leq c_\rho, \forall n = 1, \dots, N$. Then, we have:*

$$(1 - \delta) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2 \leq \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_l} \right\|_2 \leq (1 + \delta) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2$$

$$\text{and } \delta = c_\rho^2 \|\mathbf{W}_l^{(1)}\|_2 \|\mathbf{W}_l^{(2)}\|_2.$$

Proof. See Section A.3 \square

Here, $\|\cdot\|_2$ is the induced matrix norm and is the largest singular value of the matrix, which is known to be upper bounded by the entry-wise ℓ_2 norm. This means that norm-preservation is enforced throughout the training process, as long as the norm of the weights are small, not just at the beginning of the training by good initialization. This is the case in practice, since the weights of the network are regularized either explicitly using ℓ_2 regularization, also known as weight decay, or implicitly by the optimization algorithm [20], [21]. Thus, the gradients will have very similar magnitudes at different layers, and this leads to well-conditioning and faster convergence [14].

Although Theorem 1 holds for linear blocks as well, we can derive tighter bounds for linear residual blocks by taking a slightly different approach. For that, we model each linear residual block as:

$$\mathbf{x}_{l+1} = \mathbf{x}_l + \mathbf{W}_l \mathbf{x}_l, \quad (2)$$

where, $\mathbf{x}_l, \mathbf{x}_{l+1} \in \mathbb{R}^N$ are respectively the input and output of the l^{th} residual block, with dimension N . The weight matrix $\mathbf{W}_l \in \mathbb{R}^{N \times N}$ is the tunable linear transformation. The goal of learning is to compute a function $\mathbf{y} = \mathcal{M}(\mathbf{x}, \mathcal{W})$, where $\mathbf{x} = \mathbf{x}_1$ is the input, $\mathbf{y} = \mathbf{x}_{L+1}$ is its corresponding output, and \mathcal{W} is the collection of all adjustable linear transformations, i.e., $\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L$. In the case of simplified linear residual networks, function $\mathcal{M}(\mathbf{x}, \mathcal{W})$ is a stack of L residual blocks, as formulated in (2). Mathematically speaking, we have:

$$\mathbf{y} = \mathcal{M}(\mathbf{x}, \mathcal{W}) = \prod_{l=1}^L (\mathbf{I} + \mathbf{W}_l) \mathbf{x}, \quad (3)$$

where \mathbf{I} is an $N \times N$ identity matrix. $\mathcal{M}(\mathbf{x}, \mathcal{W})$ is used to learn a linear mapping $\mathbf{R} \in \mathbb{R}^{N \times N}$ from its inputs and outputs. Furthermore, assume that \mathbf{y} is contaminated with independent identically distributed (i.i.d) Gaussian noise, i.e., $\hat{\mathbf{y}} = \mathbf{R}\mathbf{x} + \epsilon$, where ϵ is a zero mean noise vector with covariance matrix \mathbf{I} . Hence, our objective is to minimize the expected error of the maximum likelihood estimator as:

$$\min_{\mathcal{W}} \mathcal{E}(\mathcal{W}) = \mathbb{E} \left\{ \frac{1}{2} \|\hat{\mathbf{y}} - \mathcal{M}(\mathbf{x}, \mathcal{W})\|_2^2 \right\}, \quad (4)$$

where the expectation \mathbb{E} is with respect to the population (\mathbf{x}, \mathbf{y}) . The following theorem states the bound on the norm preservation of the linear residual blocks.

Theorem 2. *For learning a linear map, $\mathbf{R} \in \mathbb{R}^{N \times N}$, between its input \mathbf{x} and output \mathbf{y} contaminated with i.i.d Gaussian noise, using a network consisting of L linear residual blocks of form $\mathbf{x}_{l+1} = \mathbf{x}_l + \mathbf{W}_l \mathbf{x}_l$, there exists a global optimum for $\mathcal{E}(\cdot)$, as defined in (4), such that for all residual blocks we have*

$$(1 - \delta) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2 \leq \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_l} \right\|_2 \leq (1 + \delta) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2$$

for $L \geq 3\gamma$, where $\delta = \frac{c}{L}$, $c = 2(\sqrt{\pi} + \sqrt{3\gamma})^2$, and $\gamma = \max(|\log \sigma_{\max}(\mathbf{R})|, |\log \sigma_{\min}(\mathbf{R})|)$, where $\sigma_{\max}(\mathbf{R})$ and $\sigma_{\min}(\mathbf{R})$, respectively, are maximum and minimum singular values of \mathbf{R} .

Proof. See Section A.2 \square

Similar to the nonlinear residual blocks, the linear blocks become more norm-preserving as we increase the depth. However, the linear blocks become norm-preserving at a faster rate. The gradient norm ratio for the linear blocks approaches 1 with a rate of $\mathcal{O}(\frac{1}{L})$, while this ratio for nonlinear blocks approaches 1 with a rate of $\mathcal{O}(\frac{\log(L)}{L})$.

3 PROCRUSTES RESIDUAL NETWORK

As depicted in Figure 1(a), residual networks contain four different types of blocks: (i) convolution layer (first layer), (ii) fully connected layer (last layer), (iii) transition blocks (which change the dimension) as depicted in Figure 1(c), and

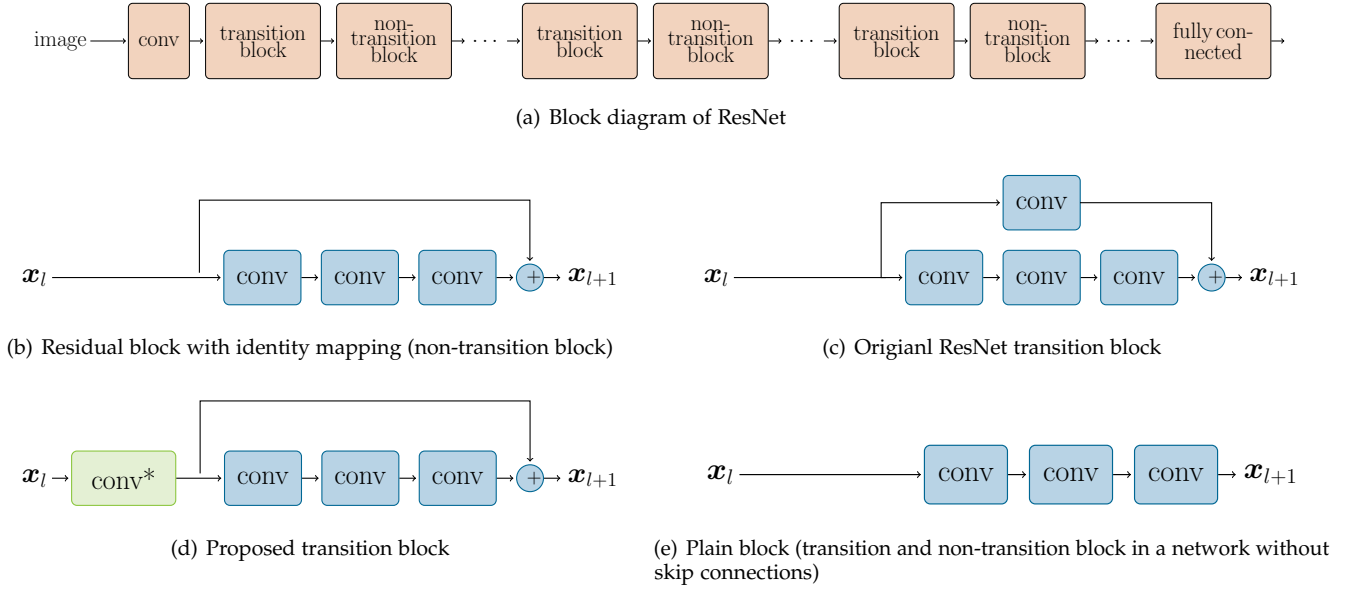


Figure 1: ResNet architecture and its building blocks. Each conv block represents a sequence of batch normalization, ReLU, and convolution layers. conv* block represents the regularized convolution layer.

(iv) residual blocks with identity skip connection, as illustrated in Figure 1(b), which we also refer to as non-transition blocks. Theoretical investigation presented in Section 2 holds only for residual blocks with identity mapping as the skip connection. Such identity skip connection cannot be used in the transition blocks, since the size of the input is not the same as the size of output. If the benefits of residual networks can be explained, at least partly, by norm-preservation, then one can improve them by alternative methods for preserving the norm. In this section, we propose to modify the transition blocks of ResNet architecture, to make them norm-preserving. Due to multiplicative effect through the layers, making these layers norm-preserving may be important, although they make up a small portion of the network. In the following, we discuss how to preserve the norm of the back-propagated gradients across all the blocks of the network.

As depicted in Figure 1(c), in the original ResNet architecture, the dimension changing blocks, also known as transition blocks, use 1×1 convolution with stride of 2 in their skip connections to match the dimension of input and output activations. Such transition blocks are not norm-preserving in general.

Figure 1(d) shows the block diagram of the proposed norm-preserving transition block. To change the dimension in a norm-preserving manner, we utilize a norm preserving convolution layer, conv*. For that, we project the convolution kernel onto the set of norm preserving kernels by setting its singular values. Here, we show how we can make the convolution layer norm preserving by regularizing the singular values, without using singular value decomposition. Specifically, the gradient of a convolution layer with kernel of size k , with c input channels, and d output channels can be formulated as:

$$\Delta_{\mathbf{x}} = \hat{\mathbf{W}} \Delta_{\mathbf{y}}, \quad (5)$$

where $\Delta_{\mathbf{x}}$ and $\Delta_{\mathbf{y}}$ respectively are the gradients with respect to the input and output of the convolution. $\Delta_{\mathbf{y}}$ is

an $n^2 d$ dimensional vector, representing $n \times n$ pixels in d output channels, and $\Delta_{\mathbf{x}}$ is an $n^2 c$ dimensional vector, representing the gradient at the input. Furthermore, $\hat{\mathbf{W}}$ is an $n^2 c \times n^2 d$ dimensional matrix embedding the back-propagation operation for the convolution layer. We can represent this linear transformation as:

$$\Delta_{\mathbf{x}} = \sum_{j=1}^{n^2 c} \sigma_j \mathbf{u}_j \langle \Delta_{\mathbf{y}}, \mathbf{v}_j \rangle, \quad (6)$$

where $\{\sigma_j, \mathbf{u}_j, \mathbf{v}_j\}$ is the set of singular values and singular vectors of $\hat{\mathbf{W}}$. Furthermore, since the set of the right singular vectors, i.e., $\{\mathbf{v}_j\}$, is an orthonormal basis set for $\Delta_{\mathbf{y}}$, we can write the gradient at the output as:

$$\Delta_{\mathbf{y}} = \sum_{j=1}^{n^2 d} \mathbf{v}_j \langle \Delta_{\mathbf{y}}, \mathbf{v}_j \rangle.$$

Thus, we can compute the expected value of the norm of the gradients as:

$$\begin{aligned} \mathbb{E}[\|\Delta_{\mathbf{x}}\|_2^2] &= \sum_{j=1}^{n^2 c} \sigma_j^2 \mathbb{E}[\langle \Delta_{\mathbf{y}}, \mathbf{v}_j \rangle^2], \\ \mathbb{E}[\|\Delta_{\mathbf{y}}\|_2^2] &= \sum_{j=1}^{n^2 d} \mathbb{E}[\langle \Delta_{\mathbf{y}}, \mathbf{v}_j \rangle^2], \end{aligned}$$

where we use the fact that $\mathbf{u}_i^T \mathbf{u}_j = \mathbf{v}_i^T \mathbf{v}_j = 0$ for $i \neq j$ and $\mathbf{u}_j^T \mathbf{u}_j = \mathbf{v}_j^T \mathbf{v}_j = 1$ and the expectation is over the data population. We propose to preserve the norm of the gradient, i.e., $\mathbb{E}[\|\Delta_{\mathbf{x}}\|_2^2] = \mathbb{E}[\|\Delta_{\mathbf{y}}\|_2^2]$, by setting all the non-zero singular values to σ . It is easy to show that we can achieve this by setting

$$\sigma^2 = \frac{\sum_{j=1}^{n^2 d} \mathbb{E}[\langle \Delta_{\mathbf{y}}, \mathbf{v}_j \rangle^2]}{\sum_{j, \sigma_j \neq 0} \mathbb{E}[\langle \Delta_{\mathbf{y}}, \mathbf{v}_j \rangle^2]}, \quad (7)$$

where the summation in the denominator is over the singular vectors \mathbf{v}_j corresponding to the nonzero singular values, i.e., $\sigma_j \neq 0$. The ratio in (7) is the ratio of expected energy of $\Delta_{\mathbf{y}}$, i.e. $\mathbb{E}[\|\Delta_{\mathbf{y}}\|_2^2]$, divided by the portion of energy that does not lie in the null space of $\hat{\mathbf{W}}$. We make the assumption that this ratio can be approximated by $\frac{n^2 d}{n^2 \min(d, c)}$. This assumption implies that about $\frac{n^2 \min(d, c)}{n^2 d}$ of the total energy of $\Delta_{\mathbf{y}}$ will lie in the $n^2 \min(d, c)$ -dimensional subspace, corresponding to orthonormal basis set $\{\mathbf{v}_j | \sigma_j \neq 0\}$, of our $n^2 d$ -dimensional space. It is easy to notice that the assumption holds if the energy of $\Delta_{\mathbf{y}}$ is distributed uniformly among the directions in the basis set $\{\mathbf{v}_j\}$. But, since we are taking the sum over a large number of bases, it can also hold with high probability in cases where there is some variation in the distribution of energies along different directions. This is not a strict assumption in high dimensional spaces and we will investigate the practical relevance of this assumption in a real-world setting shortly. Thus, we can achieve norm preservation by setting all the nonzero singular values to $\sqrt{\frac{d}{\min(d, c)}}$. We can enforce this equality without using singular value decomposition. For that, we use the following theorem from [22]. This theorem states that the singular values of the convolution operator can be calculated by finding the singular values of the Fourier transform of the slices of the convolution kernels.

Theorem 3. (Theorem 6 from [22]) For any convolution kernel $\mathbf{K} \in \mathbb{R}^{k \times k \times d \times c}$ acting on an $n \times n \times d$ input, let $\hat{\mathbf{W}}$ be the matrix encoding the linear transformation computed by a convolutional layer parameterized by \mathbf{K} . Also, for each $u, v \in [n] \times [n]$, let $\mathbf{P}^{(u, v)} \in \mathbb{C}^{d \times c}$ be the matrix given by $\mathbf{P}_{i, j}^{(u, v)} = (\mathcal{F}_n(\mathbf{K}_{:::, i, j}))_{u, v}$, where $\mathcal{F}_n(\cdot)$ is the operator describing an $n \times n$ 2D Fourier transform. Then, the set of singular values of $\hat{\mathbf{W}}$ is the union (allowing repetitions) of all the singular values of $\mathbf{P}^{(u, v)}$ matrices $\forall u, v$.

Proof. See [22]. \square

Hence, to satisfy the condition (7), we can set all the nonzero singular values of $\mathbf{P}^{(u, v)}$ to $\sqrt{\frac{d}{\min(d, c)}}$ for all u and v . This can be done by finding the matrix $\hat{\mathbf{P}}^{(u, v)}$ that minimizes $\|\mathbf{P}^{(u, v)} - \hat{\mathbf{P}}^{(u, v)}\|_F^2$, such that $\hat{\mathbf{P}}^{(u, v)T} \hat{\mathbf{P}}^{(u, v)} = \frac{d}{\min(d, c)} \mathbf{I}$, where $\|\cdot\|_F$ denotes the Frobenius norm and \mathbf{I} is a $c \times c$ identity matrix. It can be shown that the solution to this problem is given by

$$\hat{\mathbf{P}}^{(u, v)} = \sqrt{\frac{d}{\min(d, c)}} \mathbf{P}^{(u, v)} (\mathbf{P}^{(u, v)T} \mathbf{P}^{(u, v)})^{-\frac{1}{2}}. \quad (8)$$

This is closely related to *Procrustes* problems, in which the goal is to find the closest orthogonal matrix to a given matrix [23]. Finding the inverse of the square root of product $\mathbf{P}^{(u, v)T} \mathbf{P}^{(u, v)}$ can be computationally expensive, specifically for large number of channels c . Thus, we exploit an iterative algorithm that computes the inverse of the square root using

Algorithm 1 Update rules for transition kernels at each iteration

Input: Convolution kernel \mathbf{K} at the current iteration

- 1: Perform the gradient descent step on the kernel \mathbf{K} .
- 2: Calculate $\mathbf{P}^{(u, v)}$ for each $u, v \in [n] \times [n]$ as $\mathbf{P}_{i, j}^{(u, v)} = (\mathcal{F}_n(\mathbf{K}_{:::, i, j}))_{u, v}$.
- 3: Compute $(\mathbf{P}^{(u, v)T} \mathbf{P}^{(u, v)})^{-\frac{1}{2}}$ using (9).
- 4: Calculate $\hat{\mathbf{P}}^{(u, v)}$ using (8).
- 5: Update \mathbf{K} using the inverse 2D Fourier transform of $\hat{\mathbf{P}}^{(u, v)}$.

only matrix multiplications. Specifically, one can use the following iterations to compute $(\mathbf{P}^{(u, v)T} \mathbf{P}^{(u, v)})^{-\frac{1}{2}}$ [24]:

$$\begin{aligned} \mathbf{T}_k &= 3\mathbf{I} - \mathbf{Z}_k \mathbf{Y}_k, \\ \mathbf{Y}_{k+1} &= \frac{1}{2} \mathbf{Y}_k \mathbf{T}_k, \\ \mathbf{Z}_{k+1} &= \frac{1}{2} \mathbf{T}_k \mathbf{Z}_k, \end{aligned} \quad (9)$$

for $k = 0, 1, \dots$ and the iterators are initialized as:

$$\mathbf{Y}_0 = \mathbf{P}^{(u, v)T} \mathbf{P}^{(u, v)}, \mathbf{Z}_0 = \mathbf{I}.$$

It has been shown that \mathbf{Z}_k converges to $(\mathbf{P}^{(u, v)T} \mathbf{P}^{(u, v)})^{-\frac{1}{2}}$ quadratically [24]. Since the iterations only involve matrix multiplication, they can be implemented efficiently on GPUs.

Thus, to keep the convolution kernels norm preserving throughout the training, at each iteration, we compute the matrices $\mathbf{P}^{(u, v)}$ and set the nonzero singular values using (8). Algorithm 1 summarizes the operations performed at each iteration on the kernels of the regularized convolution layers. To keep the desired norm-preservation property after performing the gradient descent step, such as SGD, Adam, etc, the proposed scheme is used to re-enforce norm-preservation on the updated kernel. In this manner, we can maintain norm-preservation, while updating the kernel during the training. Our experiments in Section 4 show that performing the proposed projection on the transition block of deep ResNets increases the training time by less than 8%. Also, since the number of transition blocks are independent of depth, the deeper the network gets, the computational overhead of the proposed modification becomes less significant. Figure 1(d) shows the diagram of the proposed transition block, where a regularized convolution layer, conv^* , is used to change the dimension. Hence, we are able to exploit a regular residual block with identity mapping, which is norm preserving.

Similar to [3], to take into the account the effect of a ReLU nonlinearity and to make a Conv-ReLU layer norm-preserving, we just need to add a factor of $\sqrt{2}$ to the singular values and set them to $\sqrt{\frac{2d}{\min(d, c)}}$. Intuitively, the element-wise ReLU sets half of the units to zero on average, making the expected value of the energy of the gradient equal to $\mathbb{E}[\|\Delta_{\mathbf{x}}\|_2^2] = \frac{1}{2} \sum_{j=1}^{n^2 c} \sigma_j^2 \mathbb{E}[|\langle \Delta_{\mathbf{y}}, \mathbf{v}_j \rangle|^2]$. Therefore, to compensate this, we need to satisfy this condition:

$$\frac{1}{2} \sum_{j=1}^{n^2 c} \sigma_j^2 \mathbb{E}[|\langle \Delta_{\mathbf{y}}, \mathbf{v}_j \rangle|^2] = \sum_{j=1}^{n^2 d} \mathbb{E}[|\langle \Delta_{\mathbf{y}}, \mathbf{v}_j \rangle|^2]$$

It is also worthwhile to mention that since we are trying to preserve the norm of the backward signal, the variable n

in Theorem 3 represents the size of feature map size at the output of the convolution.

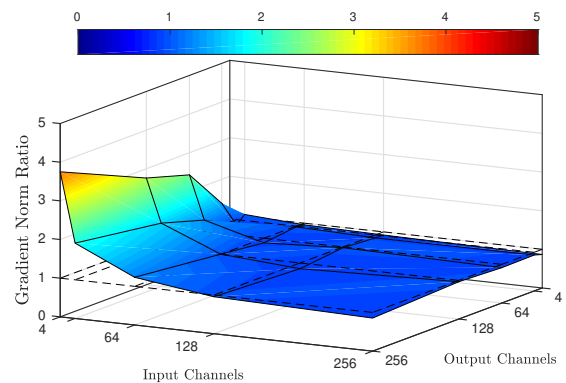
To evaluate the effectiveness of the proposed projection, we design the following experiment. We perform the projection on the convolution layers of a small 3-layer network. The network consists of 3 convolutional layers, followed by ReLU non-linearity. To examine the gradient norm ratio for different number of input and output channels, the second layer is a 3×3 convolution with c input channels and d output channels. The first and third layers are 1×1 convolutions to change the number of channels and to match the size of the input and output layers. Figure 2 shows the gradient norm ratio, i.e., $\|\frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}}\|_2$ to $\|\frac{\partial \mathcal{E}}{\partial \mathbf{x}_l}\|_2$, for different values of c and d at 10th training epoch on CIFAR-10, with and without the proposed projection. The values are averaged over 10 different runs.

It is evident that the proposed projection enhances the norm preservation of the Conv-ReLU layer, as it moves the gradient norm ratios toward 1. The only failure case is for networks with very small c and $c \ll d$. This is because, due to the smaller size of the space, our assumption that the energy of the signal in the n^2c dimensional subspace, corresponding to the n^2c non-zero singular values, is approximately $\frac{n^2c}{n^2d}$ of the total energy of the signal, is violated with higher probability. However, in more practical settings, where the number of channels is large and the assumption is held, the proposed projection performs as expected. This experiment illustrates the validity of our analysis as well as the effectiveness of the proposed projection for such practical scenarios. In the next section, we demonstrate the advantages of the proposed method for image classification task.

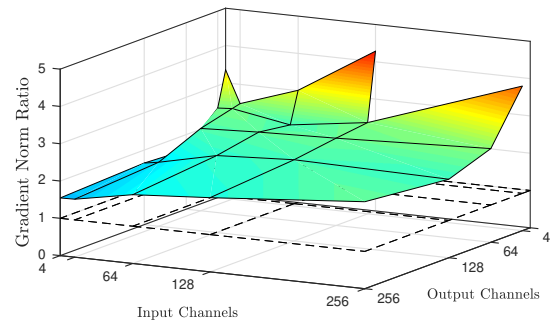
4 EXPERIMENTS

To validate our theoretical investigation, presented in Section 2, and to empirically demonstrate the behavior and effectiveness of the proposed modifications, we experimented with Residual Network (ResNet) and the proposed Procrustes Residual Network (ProcResNet) architectures on CIFAR10 and CIFAR100 datasets. Training and testing datasets contain 50,000 and 10,000 images of visual classes, respectively [25]. Standard data augmentation (flipping and shifting), same as [5], [6], [8], is adopted. Furthermore, channel means and standard deviations are used to normalize the images. The network is trained using stochastic gradient descent. The weights are initialized using the method proposed in [3] and the initial learning rate is 0.1. Batch size of 128 is used for all the networks. The weight decay is 10^{-4} and momentum is 0.9. The results are based on the top-1 classification accuracy.

Experiments are performed on three different network architectures: 1) **ResNet** contains one convolution layer, L residual blocks, three of which are transition blocks, and one fully connected layer. Each residual block consists of three convolution layers, as depicted in Figure 1(b) and Figure 1(c), resulting in a network of depth $3L + 2$. This is the same architecture as in [6]. 2) **ProcResNet** has the same architecture as ResNet, except the transition layers are modified, as explained in Section 3. In this design, 3 extra convolution layers are added to the network. However, we can use the first convolution layer of the original ResNet design to match the dimensions and only add two extra



(a) With the Proposed Regularization



(b) Without Regularization

Figure 2: The ratio of gradient norm at output to gradient norm at input, i.e., $\|\frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}}\|_2$ to $\|\frac{\partial \mathcal{E}}{\partial \mathbf{x}_l}\|_2$, of a convolution layer for different number of input and output channels at 10th training epoch (a) with, and (b) without the proposed regularization on the singular values of the convolution.

layers. This leads to a network of depth $3L + 4$. 3) **Plain** network is also same as ResNet without the skip connection in all the L residual blocks, as shown in Figure 1(e).

Furthermore, to decrease the computational burden of the proposed regularization, we perform the projection, as described in Section 3, every 2 iterations. This reduces the computation time significantly without hurting the performance much. In this setting, performing the proposed regularization increases the training time for ResNet164 about 7.6%. However, since we perform the regularization only on three blocks, regardless of the depth, as the network becomes deeper the computational overhead becomes less significant. For example, implementing the same projections on ResNet1001 increases the training time by only 3.5%. This is significantly less computation compared to regularization using SVD, which leads to 53% and 23% training time overhead for ResNet164 and ResNet1001, respectively¹.

4.1 Norm-Preservation

In the first set of experiments, the behavior of different architectures is studied as the function of network depth. To this end, the ratio of gradient norm at output to gradient norm at input, i.e., $\|\frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}}\|_2$ to $\|\frac{\partial \mathcal{E}}{\partial \mathbf{x}_l}\|_2$, is captured for all the

1. An implementation of ProcResNet is provided here: <https://github.com/zaemzadeh/ProcResNet>

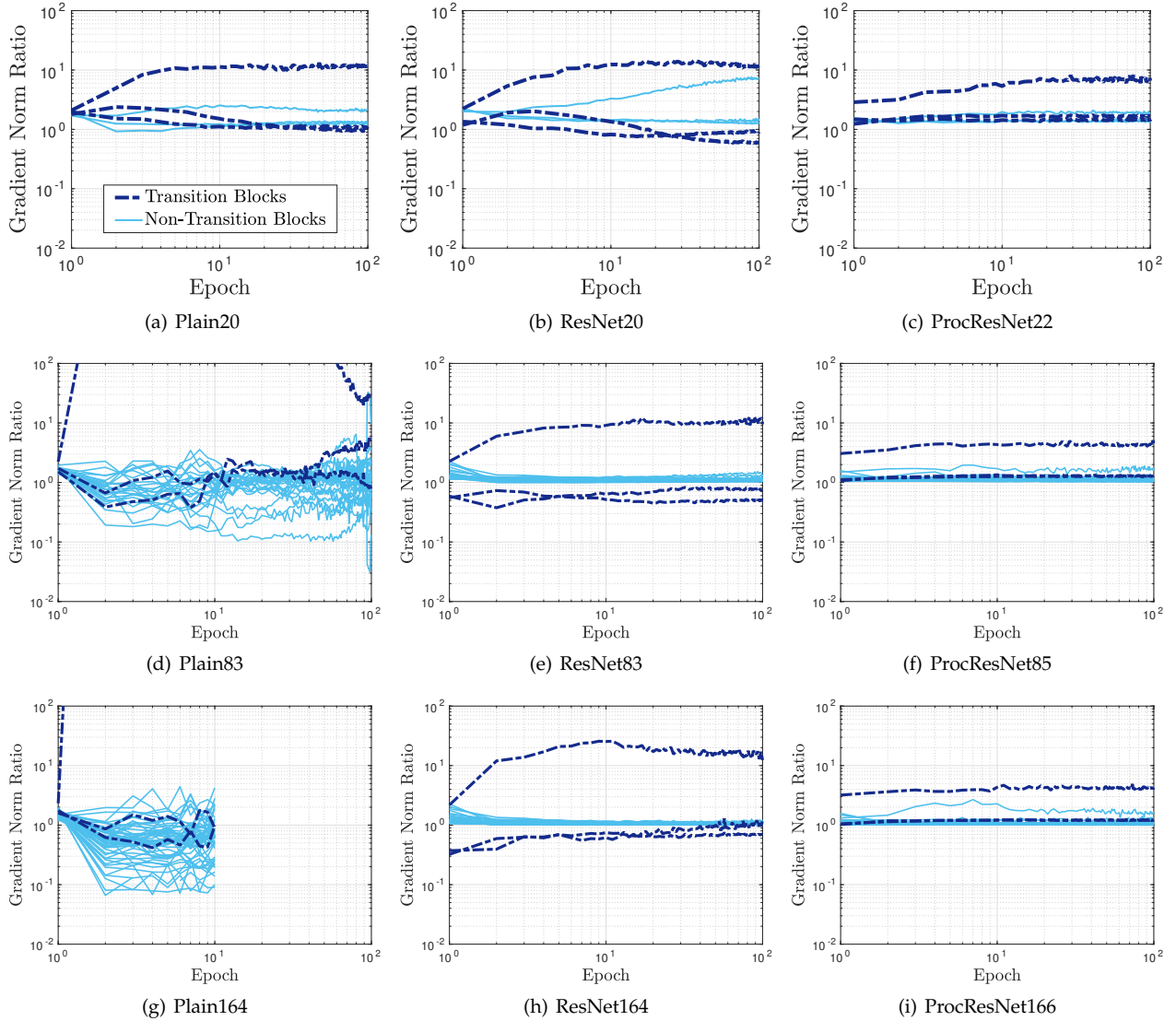


Figure 3: Training on CIFAR10. Gradient norm ratio over the first 100 epochs for transition blocks (blocks that change the dimension) and non-transition blocks (blocks that do not change the dimension). The darker color lines represent the transition blocks and the lighter color lines represent the non-transition blocks. The proposed regularization enhances the norm-preservation of the transition blocks effectively.

residual blocks², both transition and non-transition. Figure 3 shows the ratios for different blocks over training epochs. We ran the training for 100 epochs, without decaying the learning rate. Plain network (Figure 3.(g)) with 164 layers became numerically unstable and the training procedure stopped after 10 epochs.

Several interesting observations can be made from this experiment:

- This experiment emphasizes the fact that one needs more than careful initialization to make the network norm-preserving. Although the plain network is initially norm-preserving, the range of the gradient norm ratios becomes very large and diverges from 1, as the parameters are updated. However, ResNet and ProcResNet are able to enforce the norm-preservation during training procedure by using identity skip connection.

2. In Plain architecture, which does not have skip connections, the gradient norm ratio is obtained at the input and output of its building blocks as depicted in Figure 1(e).

- As the networks become deeper, the plain network becomes less norm preserving, which leads to numerical instability, optimization difficulty, and performance degradation. On the contrary, the non-transition blocks, the blocks with identity mapping as skip connection, of ResNet and ProcResNet become extra norm preserving. This is in line with our theoretical investigation for linear residual networks, which states that as we stack more residual blocks the network becomes extra norm-preserving.
- Comparing Plain83 (Figure 3(d)) and Plain164 (Figure 3(h)) networks, it can be observed that most of the blocks behave fairly similar, except one transition block. Specifically, in Plain83, the gradient norm ratio of the first transition block goes up to 100 in the first few epochs. But it eventually decreases and the network is able to converge. On the other hand, in Plain164, the gradient norm ratio of the same block becomes too large, which makes the network unable to converge. Hence, a single block is enough to make the optimization difficult and numerically unstable.

This highlights the fact that it is necessary to enforce norm-preservation on all the blocks.

- In ResNet83 (Figure 3(e)) and ResNet164 (Figure 3(h)), it is easy to notice that only 3 transition blocks are not norm preserving. As mentioned earlier, due to multiplicative effect, the magnitude of the gradient will not be preserved because of these few blocks.
- The behaviors of ResNet and Plain architectures are fairly similar for depth of 20. This was somehow expected, since it is known that the performance gain achieved by ResNet is more significant in deeper architectures [5]. However, even for depth of 20, ProcResNet architecture is more norm preserving.
- In ProcResNet, the only block that is less norm preserving is the first transition block, where the 3 RGB channels are transformed into 64 channels. This is because, as we have shown in Figure 2, under such condition, where the number of input channels is very small, the assumption that energy of the gradient signal in the low-dimensional subspace, corresponding to the few non-zero singular values, is approximately proportional to the size of the subspace is violated with higher probability.
- The ratios of the gradients for all networks, even the Plain network, are roughly concentrated around 1, while training is stable. This shows that some degree of norm preservation exists in any stable network. However, as clear in the Plain network, such biases of the optimizer is not enough and we need skip connections to enforce norm preservation throughout training and to enjoy its desirable properties. Furthermore, although the transition blocks of ResNet tend to converge to be more norm preserving, our proposed modification enforces this property for all the epochs, which leads to stability and performance gain, as will be discussed shortly.

This experiment both validates our theoretical arguments and clarifies some of the inner workings of ResNet architecture, and also shows the effectiveness of the proposed modifications in ProcResNet. It is evident that, as stated in Theorem 1, addition of identity skip connection makes the blocks increasingly extra norm-preserving, as the network becomes deeper. Furthermore, we have been able to enhance norm-preserving property by applying the changes proposed in Section 3.

4.2 Optimization Stability and Learning Dynamics

In the next set of experiments, numerical stability and learning dynamics of different architectures is examined. For that, loss and classification error, in both training and testing phases, are depicted in Figure 4. This experiment illustrates that how optimization stability of deep networks is improved significantly, and how it can be further improved by having norm preservation in mind during the design procedure.

As depicted in Figure 4, unlike the plain network, training error and loss curves corresponding to ResNet and ProcResNet architectures are consistently decreasing as the number of layers increases, which was the main motivation behind proposing residual blocks [5]. Moreover, Figure 4(a) and Figure 4(d) show that the plain networks have a poor generalization performance. The fluctuations in testing error shows that the points along the optimization path of the

TABLE 1: Mean and maximum generalization gap (%) during the first 100 epochs of training on CIFA10 for different network architectures, averaged over 10 runs.

| Depth | Plain | | ResNet | | ProcResNet | |
|-------|-------|------|--------|------|------------|------------|
| | mean | max | mean | max | mean | max |
| 20 | 6.7 | 20.0 | 5.5 | 23.1 | 2.3 | 8.3 |
| 83 | 7.5 | 30.1 | 5.1 | 12.5 | 2.0 | 7.7 |
| 164 | - | - | 5.2 | 18.7 | 3.3 | 8.7 |

plain networks do not generalize well. This issue is also present, to a lesser extent, in ResNet architecture. Comparing Figure 4(h) and 4(b), we can see that the fluctuations are more apparent in deeper ResNet networks. However, in proposed ProcResNet architecture, the amplitude of the fluctuations is smaller and does not change as the depth of the network is increased. This indicates that ProcResNet architecture is taking a better path toward the optimum and has better generalization performance.

To quantify this, we repeated the training 10 times with different random seeds and measured the generalization gap, which is the difference between training and testing classification error, for the first 100 epochs. Table 1 shows the mean and max generalization gap, averaged over 10 different runs. This results indicate that generalization gap of ProcResNet is smaller. Furthermore, the generalization gap fluctuates far less significantly for ProcResNet, as quantified by the difference between mean generalization gap and maximum generalization gap.

The implication of this is that by modifying only a few blocks in an extremely deep network, it is possible to make the network more stable and improve the learning dynamics. This emphasizes the utmost importance of norm-preservation of all blocks in avoiding optimization difficulties of very deep networks. Moreover, this sheds light on the reasons why architectures using residual blocks, or identity skip connection in general, perform so well and are easier to optimize.

4.3 Classification Performance

In this section, we show the impact of the proposed norm-preserving transition blocks on the classification performance of ResNet. Table 2 compares the performance of ResNet and its EraseReLU version, as proposed in [26], with and without the proposed transition blocks. The results for standard ResNet are the best results reported by [6] and [26] and the results of ProcResNet are obtained by making the proposed changes to standard ResNet implementation.

Table 2 shows that the proposed network performs better than the standard ResNet. This performance gain comes with a slight increase the number of parameters (under 1%) and by changing only 3 blocks. The total number of residual blocks for ResNet164 and ResNet1001 are 54 and 333, respectively. Furthermore, Figure 5 compares the parameter efficiency of ResNet and ProcResNet architectures. The results indicate that the proposed modification can improve the parameter efficiency significantly. For example, ProcResNet274 (with 2.82M parameter) slightly outperforms ResNet1001 (with 10.32M parameters). This translates into about 4x reduction in the number of parameters to achieve the same classification accuracy. This illustrates that we are able to improve the

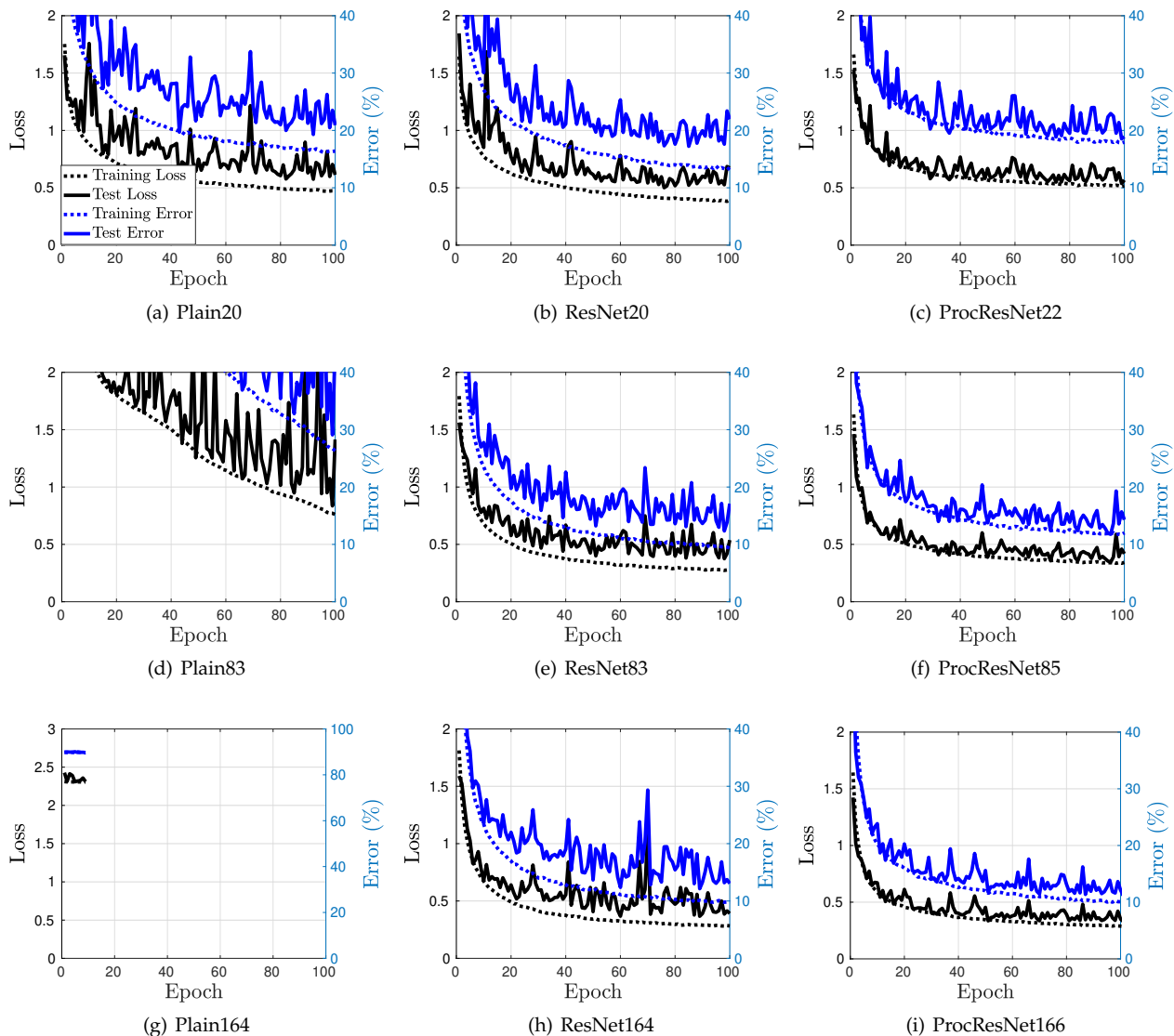


Figure 4: Loss (black lines) and error (blue lines) during training procedure on CIFAR10. Solid lines represent the test values and dotted lines represent the training values. This experiments shows how the residual connections enhance the stability of the optimization and how the proposed regularization enhances the stability even further.

performance by changing a tiny portion of the network and emphasizes the importance of norm-preservation in the performance of neural networks.

Finally, Table 3 investigates the impact of changing the architecture, i.e., moving the convolution layer from the skip connection to before the skip connection, and performing the proposed regularization, separately. Each of these design components have positive impact on the performance of the network, as both of them enhance the norm preservation of the transition block, which further highlights the impact of norm preservation on the performance of the network.

5 CONCLUSIONS

This paper theoretically analyzes building blocks of residual networks and demonstrates that adding identity skip connection makes the residual blocks norm-preserving. Furthermore, the norm-preservation is enforced during the

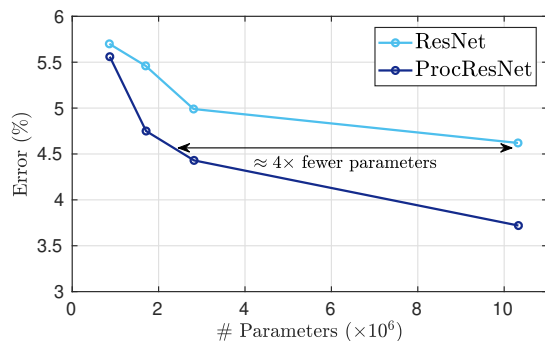


Figure 5: Comparison of the parameter efficiency on CIFAR10 between ResNet and ProcResNet.

training procedure, which makes the optimization stable and improves the performance. This is in contrast to ini-

TABLE 2: Performance of different methods on CIFAR-10 and CIFAR-100 using moderate data augmentation (flip/translation). The modified transition blocks in ProcResNet can improve the accuracy of ResNet significantly.

| Architecture | Setting | # Params | Depth | Error (%) | |
|--------------|-----------------|----------|-------|-------------|--------------|
| | | | | CIFAR10 | CIFAR100 |
| ResNet [6] | pre-activation | 1.71M | 164 | 5.46 | 24.33 |
| | | 10.32M | 1001 | 4.62 | 22.71 |
| | ErasedReLU [26] | 1.70M | 164 | 4.65 | 22.41 |
| | | 10.32M | 1001 | 4.10 | 20.63 |
| ProcResNet | pre-activation | 1.72M | 166 | 4.75 | 22.61 |
| | | 10.33M | 1003 | 3.72 | 19.99 |
| | ErasedReLU [26] | 1.72M | 166 | 4.53 | 21.91 |
| | | 10.33M | 1003 | 3.42 | 18.12 |

TABLE 3: Ablation study on ResNet with 164 layers on CIFAR100.

| Transition Block | Projection | Error (%) |
|------------------|------------|-----------|
| Original | No | 24.33 |
| Modified | No | 23.06 |
| Modified | Yes | 22.61 |

tialization techniques, such as [14], which ensure norm-preservation only at the beginning of the training. Our experiments validate our theoretical investigation by showing that (i) identity skip connection results in norm preservation, (ii) residual blocks become extra norm-preserving as the network becomes deeper, and (iii) the training can become more stable through enhancing the norm preservation of the network. Our proposed modification of ResNet, Procrustes ResNet, enforces norm-preservation on the transition blocks of the network and is able to achieve better optimization stability and performance. For that we propose an efficient regularization technique to set the nonzero singular values of the convolution operator, without performing singular value decomposition. Our findings can be seen as design guidelines for very deep architectures. By having norm-preservation in mind, we will be able to train extremely deep networks and alleviate the optimization difficulties of such networks.

6 ACKNOWLEDGEMENTS

This research is based upon work supported in parts by the National Science Foundation under Grants No. 1741431 and CCF-1718195 and the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via IARPA R&D Contract No. D17PC00345. The views, findings, opinions, and conclusions or recommendations contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the NSF, ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon.

REFERENCES

- [1] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis, "Mastering the game of Go without human knowledge," *Nature*, vol. 550, pp. 354–359, 10 2017.
- [2] G. F. Montufar, R. Pascanu, K. Cho, and Y. Bengio, "On the Number of Linear Regions of Deep Neural Networks," in *Advances in Neural Information Processing Systems 27* (Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, eds.), pp. 2924–2932, Curran Associates, Inc., 2014.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification," in *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 1026–1034, IEEE, 12 2015.
- [4] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," 6 2015.
- [5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, IEEE, 6 2016.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 9908 LNCS, pp. 630–645, Springer, Cham, 10 2016.
- [7] R. K. Srivastava, K. Greff, and J. Schmidhuber, "Training Very Deep Networks," 2015.
- [8] G. Huang, Z. Liu, L. v. d. Maaten, and K. Q. Weinberger, "Densely Connected Convolutional Networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2261–2269, IEEE, 7 2017.
- [9] A. E. Orhan and X. Pitkow, "Skip connections eliminate singularities," in *6th International Conference on Learning Representations, ICLR 2018 - Conference Track Proceedings*, 1 2018.
- [10] M. Hardt and T. Ma, "Identity matters in deep learning," in *5th International Conference on Learning Representations, ICLR 2017 - Conference Track Proceedings*, 11 2017.
- [11] K. Kawaguchi, "Deep Learning without Poor Local Minima," 2016.
- [12] D. Balduzzi, M. Frean, L. Leary, J. P. Lewis, K. W.-D. Ma, and B. McWilliams, "The Shattered Gradients Problem: If resnets are the answer, then what is the question?," 7 2017.
- [13] A. Veit, M. J. Wilber, and S. Belongie, "Residual Networks Behave Like Ensembles of Relatively Shallow Networks," in *Advances in Neural Information Processing Systems 29* (D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, eds.), pp. 550–558, Curran Associates, Inc., 2016.
- [14] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," 3 2010.
- [15] L. Dinh, J. Sohl-Dickstein, Google, B. Samy, and B. Google Brain, "Density estimation using Real NVP," in *ICLR*, 2017.
- [16] A. N. Gomez, M. Ren, R. Urtasun, and R. B. Grosse, "The reversible residual network: Backpropagation without storing activations," in *Advances in Neural Information Processing Systems*, 2017.
- [17] J. Behrmann, W. Grathwohl, R. T. Chen, D. Duvenaud, and J. H. Jacobsen, "Invertible residual networks," in *36th International Conference on Machine Learning, ICML 2019*, 2019.
- [18] P. L. Bartlett, S. N. Evans, and P. M. Long, "Representing smooth functions as compositions of near-identity functions with implications for deep network optimization," *arXiv preprint arXiv:1804.05012*, 4 2018.
- [19] K. Kawaguchi and Y. Bengio, "Depth with nonlinearity creates no bad local minima in ResNets," *Neural Networks*, 2019.
- [20] S. Gunasekar, J. D. Lee, N. Srebro, and D. Soudry, "Implicit bias of gradient descent on linear convolutional networks," in *Advances in Neural Information Processing Systems*, 2018.

- [21] E. Hoffer, I. Hubara, and D. Soudry, "Train longer, generalize better: Closing the generalization gap in large batch training of neural networks," in *Advances in Neural Information Processing Systems*, 2017.
- [22] H. Sedghi, V. Gupta, and P. M. Long, "The Singular Values of Convolutional Layers," in *International Conference on Learning Representations*, 2019.
- [23] J. C. Gower, G. B. Dijkstra, and others, *Procrustes problems*, vol. 30. Oxford University Press on Demand, 2004.
- [24] N. J. Higham, "Stable iterations for the matrix square root," *Numerical Algorithms*, vol. 15, no. 2, pp. 227–242, 1997.
- [25] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images.(2009)," tech. rep., 2009.
- [26] X. Dong, G. Kang, K. Zhan, and Y. Yang, "EraseReLU: A Simple Way to Ease the Training of Deep Convolution Neural Networks," *arXiv preprint arXiv:1709.07634*, 2017.
- [27] N. A. Derzko and A. M. Pfeffer, "Bounds for the Spectral Radius of a Matrix," *Mathematics of Computation*, vol. 19, p. 62, 4 1965.



Alireza Zaeemzadeh (S'11) received the B.S. degree in electrical engineering from the University of Tehran, Tehran, Iran, in 2014. He is currently working toward the Ph.D. degree in electrical engineering at the University of Central Florida. His current research interests lie in the areas of machine learning, linear algebra, and optimization. Alireza's awards and honors include University of Central Florida Multidisciplinary Doctoral Fellowship and Graduate Dean's Fellowship.



Nazanin Rahnavard (S'97-M'10-SM'19) received her Ph.D. in the School of Electrical and Computer Engineering at the Georgia Institute of Technology, Atlanta, in 2007. She is currently an Associate Professor in the Department of Electrical and Computer Engineering at the University of Central Florida, Orlando, Florida. Dr. Rahnavard is the recipient of NSF CAREER award in 2011 and 2020 UCF's College of Engineering and Computer Science Excellence in Research Award. She has interest and expertise in a variety of research topics in the communications, networking, signal processing, and machine learning areas. She serves on the editorial board of the Elsevier Journal on Computer Networks (COMNET) and on the Technical Program Committee of several prestigious international conferences.



Mubarak Shah, the UCF Trustee chair professor, is the founding director of the Center for Research in Computer Vision at the University of Central Florida (UCF). He is a fellow of the NAI, IEEE, AAAS, IAPR, and SPIE. He is an editor of an international book series on video computing, was editor-in-chief of Machine Vision and Applications journal, and an associate editor of ACM Computing Surveys journal. He was the program cochair of CVPR 2008, an associate editor of the IEEE T-PAMI, and a guest editor of the special issue of the International Journal of Computer Vision on Video Computing. His research interests include video surveillance, visual tracking, human activity recognition, visual analysis of crowded scenes, video registration, UAV video analysis, and so on. He has served as an ACM distinguished speaker and IEEE distinguished visitor speaker. He is a recipient of ACM SIGMM Technical Achievement award; IEEE Outstanding Engineering Educator Award; Harris Corporation Engineering Achievement Award; an honorable mention for the ICCV 2005 Where Am I? Challenge Problem; 2013 NGA Best Research Poster Presentation; 2nd place in Grand Challenge at the ACM Multimedia 2013 conference; and runner up for the best paper award in ACM Multimedia Conference in 2005 and 2010. At UCF he has received Pegasus Professor Award; University Distinguished Research Award; Faculty Excellence in Mentoring Doctoral Students; Scholarship of Teaching and Learning award; Teaching Incentive Program award; Research Incentive Award.

APPENDIX A PROOFS

A.1 Proof of Theorem 1

For a cost function $\mathcal{E}(\cdot)$ and the Jacobian \mathbf{J} of \mathbf{x}_{l+1} with respect to \mathbf{x}_l , applying chain rule, following is true:

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{x}_l} &= \mathbf{J} \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}}, \\ \mathbf{J} &= \frac{\partial \mathbf{x}_{l+1}}{\partial \mathbf{x}_l} = \mathbf{I} + D\mathbf{F}_l(\mathbf{x}_l),\end{aligned}$$

where D is the differential operator and for any \mathbf{v} with bounded norm we have:

$$D\mathbf{F}_l(\mathbf{x}_l)\mathbf{v} = \lim_{t \rightarrow 0^+} \frac{F_l(\mathbf{x}_l + t\mathbf{v}) - F_l(\mathbf{x}_l)}{t}$$

To prove Theorem 1, we first state a lemma.

Lemma 1. For any non-singular matrix $\mathbf{I} + \mathbf{M}$, we have:

$1 - \sigma_{\max}(\mathbf{M}) \leq \sigma_{\min}(\mathbf{I} + \mathbf{M}) \leq \sigma_{\max}(\mathbf{I} + \mathbf{M}) \leq 1 + \sigma_{\max}(\mathbf{M})$, where $\sigma_{\min}(\mathbf{M})$ and $\sigma_{\max}(\mathbf{M})$ represent the minimum and maximum singular values of \mathbf{M} , respectively.

Proof. Since $\sigma_{\min}(\mathbf{I} + \mathbf{M}) > 0$, the lower bound is trivial for $\sigma_{\max}(\mathbf{M}) \geq 1$. For $\sigma_{\max}(\mathbf{M}) < 1$, it is known that $|\lambda_{\max}(\mathbf{M})| < 1$, where $\lambda_{\max}(\mathbf{M})$ is the maximum eigenvalue of \mathbf{M} [27]. Thus, we can show that:

$$\begin{aligned}\sigma_{\min}(\mathbf{I} + \mathbf{M}) &= (\sigma_{\max}((\mathbf{I} + \mathbf{M})^{-1}))^{-1} \\ &= \|(\mathbf{I} + \mathbf{M})^{-1}\|_2^{-1} \\ &\stackrel{(a)}{=} \left\| \sum_{k=0}^{\infty} (-1)^k \mathbf{M}^k \right\|_2^{-1} \\ &\geq \left(\sum_{k=0}^{\infty} \|(-1)^k \mathbf{M}^k\|_2 \right)^{-1} \geq \left(\sum_{k=0}^{\infty} \|\mathbf{M}\|_2^k \right)^{-1} \\ &= \left(\frac{1}{1 - \|\mathbf{M}\|_2} \right)^{-1} = 1 - \sigma_{\max}(\mathbf{M}).\end{aligned}$$

Identity (a) is known as Neuman series of a matrix, which holds when $|\lambda_{\max}(\mathbf{M})| < 1$ and $\|\cdot\|_2$ represents the l_2 -norm of a matrix.

The upper bound is easier to show. Due to triangle inequality:

$$\sigma_{\max}(\mathbf{I} + \mathbf{M}) = \|\mathbf{I} + \mathbf{M}\|_2 \leq \|\mathbf{I}\|_2 + \|\mathbf{M}\|_2 = 1 + \sigma_{\max}(\mathbf{M}). \quad \square$$

Thus, knowing that

$$\sigma_{\min}(\mathbf{J}) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2 \leq \left\| \mathbf{J} \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2 \leq \sigma_{\max}(\mathbf{J}) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2,$$

using Lemma 1, we conclude that

$$(1 - \delta') \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2 \leq \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_l} \right\|_2 \leq (1 + \delta') \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2,$$

where, $\delta' = \sigma_{\max}(D\mathbf{F}_l(\mathbf{x}_l))$. Furthermore, we have:

$$\begin{aligned}\sigma_{\max}(D\mathbf{F}_l(\mathbf{x}_l)) &= \sup_{\mathbf{v}} \frac{\|D\mathbf{F}_l(\mathbf{x}_l)\mathbf{v}\|_2}{\|\mathbf{v}\|_2} \\ &= \sup_{\mathbf{v}} \frac{1}{\|\mathbf{v}\|_2} \left\| \lim_{t \rightarrow 0^+} \frac{F_l(\mathbf{x}_l + t\mathbf{v}) - F_l(\mathbf{x}_l)}{t} \right\|_2 \\ &= \lim_{t \rightarrow 0^+} \sup_{\mathbf{v}} \frac{1}{\|\mathbf{v}\|_2} \left\| \frac{F_l(\mathbf{x}_l + t\mathbf{v}) - F_l(\mathbf{x}_l)}{t} \right\|_2 \\ &= \lim_{t \rightarrow 0^+} \sup_{\mathbf{v}} \frac{\|F_l(\mathbf{x}_l + t\mathbf{v}) - F_l(\mathbf{x}_l)\|_2}{t\|\mathbf{v}\|_2} \\ &\leq \|F_l\|_L,\end{aligned}$$

where $\|f\|_L$ is the Lipschitz seminorm of function f and is defined as

$$\|f\|_L := \sup_{\mathbf{x} \neq \mathbf{y}} \frac{\|f(\mathbf{x}) - f(\mathbf{y})\|_2}{\|\mathbf{x} - \mathbf{y}\|_2}.$$

To conclude the proof, we use the following lemma:

Lemma 2. (Theorem 1 in [18]) Suppose we want to represent a nonlinear mapping $F : \mathbb{R}^N \rightarrow \mathbb{R}^N$, satisfying Assumption 1, with a sequence of L non-linear residual blocks of form $\mathbf{x}_{l+1} = \mathbf{x}_l + F_l(\mathbf{x}_l)$. There exists a solution such that for all residual blocks we have $\|F_l\|_L \leq c \frac{\log(2L)}{L}$.

Therefore, $\delta' = \sigma_{\max}(D\mathbf{F}_l(\mathbf{x}_l)) \leq \|F_l\|_L \leq c \frac{\log(2L)}{L} = \delta$, which concludes the proof.

A.2 Proof of Theorem 2

In the classical back-propagation equation, for a cost function $\mathcal{E}(\cdot)$ and the Jacobian \mathbf{J} of \mathbf{x}_{l+1} with respect to \mathbf{x}_l , applying chain rule, following is true:

$$\begin{aligned}\frac{\partial \mathcal{E}}{\partial \mathbf{x}_l} &= \mathbf{J} \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}}, \\ \mathbf{J} &= \frac{\partial \mathbf{x}_{l+1}}{\partial \mathbf{x}_l} = \mathbf{I} + \mathbf{W}_l^T,\end{aligned} \quad (10)$$

To prove the theorem, using Lemma 1 and knowing that

$$\sigma_{\min}(\mathbf{J}) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2 \leq \left\| \mathbf{J} \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2 \leq \sigma_{\max}(\mathbf{J}) \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2,$$

we conclude that

$$(1 - \delta') \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2 \leq \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_l} \right\|_2 \leq (1 + \delta') \left\| \frac{\partial \mathcal{E}}{\partial \mathbf{x}_{l+1}} \right\|_2.$$

where, $\delta' = \sigma_{\max}(\mathbf{W}_l)$. To conclude the proof, we use the following lemma.

Lemma 3. (Theorem 2.1 in [10]) Suppose $L \geq 3\gamma$. Then, there exists a global optimum for $\mathcal{E}(\mathcal{W})$, such that we have

$$\sigma_{\max}(\mathbf{W}_l) \leq \frac{2(\sqrt{\pi} + \sqrt{3\gamma})^2}{L}, \forall l = 1, 2, \dots, L,$$

where γ is $\max(|\log \sigma_{\max}(\mathbf{R})|, |\log \sigma_{\min}(\mathbf{R})|)$.

Using the results from this lemma and setting $\delta = \frac{2(\sqrt{\pi} + \sqrt{3\gamma})^2}{L}$, Theorem 2 follows immediately.

A.3 Proof for Corollary 1

Here, Jacobian matrix is $\mathbf{J} = \mathbf{I} + \mathbf{F}'\mathbf{W}^{(1)T}\mathbf{F}'\mathbf{W}_l^{(2)T}$, where \mathbf{F}' is the Jacobian of $\rho(\cdot)$ with respect to its input. Since we know that $0 \leq \frac{\partial \rho_n(\mathbf{x})}{\partial x_{n'}} \leq c_\rho, \forall n = n'$ and $\frac{\partial \rho_n(\mathbf{x})}{\partial x_{n'}} = 0, \forall n \neq n'$, we have $\|\mathbf{F}'\|_2 \leq c_\rho$. Therefore:

$$\begin{aligned}\|\mathbf{F}'\mathbf{W}_l^{(1)T}\mathbf{F}'\mathbf{W}_l^{(2)T}\|_2 &\leq \|\mathbf{F}'\|_2 \|\mathbf{W}_l^{(1)T}\|_2 \|\mathbf{F}'\|_2 \|\mathbf{W}_l^{(2)T}\|_2 \\ &\leq c_\rho^2 \|\mathbf{W}_l^{(1)}\|_2 \|\mathbf{W}_l^{(2)}\|_2\end{aligned}$$

and using Lemma 1 and setting $\delta = c_\rho^2 \|\mathbf{W}_l^{(1)}\|_2 \|\mathbf{W}_l^{(2)}\|_2$, Corollary 1 follows immediately.