# Convergence Theory of Learning Over-parameterized ResNet: A Full Characterization \*

Huishuai Zhang<sup>†</sup> Da Yu<sup>‡</sup> Mingyang Yi<sup>\*</sup> Wei Chen<sup>†</sup> Tie-Yan Liu<sup>†</sup>

†Microsoft Research Asia {huzhang, wche, tyliu}@microsoft.com ‡Sun Yat-Sen University yuda3@mail2.sysu.edu.cn \*University of Chinese Academy of Sciences yimingyang17@mails.ucas.edu.cn

July 15, 2019

#### Abstract

ResNet structure has achieved great empirical success since its debut. Recent work established the convergence of learning over-parameterized ResNet with a scaling factor  $\tau=1/L$  on the residual branch where L is the network depth. However, it is not clear how learning ResNet behaves for other values of  $\tau$ . In this paper, we fully characterize the convergence theory of gradient descent for learning over-parameterized ResNet with different values of  $\tau$ . Specifically, with hiding logarithmic factor and constant coefficients, we show that for  $\tau \leq 1/\sqrt{L}$  gradient descent is guaranteed to converge to the global minma, and especially when  $\tau \leq 1/L$  the convergence is irrelevant of the network depth. Conversely, we show that for  $\tau > L^{-\frac{1}{2}+c}$ , the forward output grows at least with rate  $L^c$  in expectation and then the learning fails because of gradient explosion for large L. This means the bound  $\tau \leq 1/\sqrt{L}$  is sharp for learning ResNet with arbitrary depth. To the best of our knowledge, this is the first work that studies learning ResNet with full range of  $\tau$ .

#### 1 Introduction

Residual Network (ResNet) has achieved great success in computer vision tasks since the seminal paper (He et al., 2016). Moreover, the ResNet structure has also been extended to natural language processing and achieved the state-of-the-art performance (Vaswani et al., 2017; Devlin et al., 2018). Empirically, it is widely observed that ResNet can be trained with thousands of layers easily while vanilla feedforward network can seldomly be trained with more than thirty layers (He et al., 2016). It has caught a lot of attention to study the benefit of ResNet. Several papers (Veit et al., 2016; Orhan and Pitkow, 2018; Balduzzi et al., 2017; Zhang et al., 2018) have argued the benefit of ResNet using empirical evidence and informal intuitions but do not have rigorous theoretical understanding. Recently, over-parameterization has been used as a hammer to tackle the optimization (Allen-Zhu et al., 2018; Du et al., 2018; Zou et al., 2018) and generalization (Brutzkus et al., 2017; Li and Liang, 2018; Allen-Zhu et al., 2018a; Arora et al., 2019; Cao and Gu, 2019; Nevshabur et al., 2019)

properties of neural network. Over-parameterization considers the case when the neural network is

<sup>\*</sup>The work of D. Yu and M. Yi was done when visiting MSRA.

very wide (often more neurons than samples) at each layer and enables many theoretical analysis by taking into account the statistical concentration property of the parameter matrix.

Especially, Allen-Zhu et al. (2018b); Du et al. (2018) establish the global convergence of gradient descent algorithm for learning ResNet via over-parameterization. Specifically consider the ResNet with the following residual block,

$$h_l = \phi(h_{l-1} + \tau \mathbf{W}_l h_{l-1}), \tag{1}$$

where  $\phi(\cdot)$  is the ReLU activation,  $h_l$  is the output of layer l,  $\boldsymbol{W}_l$  is the parameter of layer l and  $\tau$  is a scale factor on the parametric branch in a residual block. They both show that the width/overparameterization requirement and the convergence time for learning ResNet are polynomial with the depth of the network. However, their result holds only for a specific  $\tau = 1/\Omega(L)$ . A full picture of learning ResNet lacks. It is natural to ask

"What is the range of  $\tau$  that guarantees that ResNet can be learned with arbitrary layers?"

In this paper, we also use over-parameterization technique to tackle the above question. We fully characterize the convergence theory of learning ResNet with arbitrary depth for all values of  $\tau$ . Specifically, our contribution can be summarized as follows.

- For  $\tau \leq 1/\Omega(\sqrt{L}\log m)$ , we establish the convergence of gradient descent to global minima for learning over-parameterized ResNet with arbitrary depth.
  - For  $\tau = 1/\Omega(L\log m)$ , the same setting as Du et al. (2018); Allen-Zhu et al. (2018b), we show that over-parameterization requirement and the provable training steps for ResNet is independent of the network depth, in sharp contrast with the polynomial dependence on the network depth in previous work.
- On the converse side, we show for  $\tau > L^{-\frac{1}{2}+c}$  the network output explodes in expectation as L becomes large even at the initialization. This indicates that our achievable range  $\tau \leq 1/\Omega(\sqrt{L}\log m)$  is tight up to a logarithmic factor.

Compared with previous work, we establish the convergence of gradient descent learning ResNet with a much larger value of the scaling factor  $\tau$ . We also show that  $\tau = 1/\sqrt{L}$  is an upper bound for learning ResNet with arbitrary depth. To the best of our knowledge, such discussion on the range of  $\tau$  has not been considered before. We also present empirical evidence to illustrate our theoretical claim.

The key technical steps for proving global convergence with  $\tau = 1/\Omega(\sqrt{L}\log m)$  are as follows. We first show that the spectral norm of the ResNet mapping is upper bounded when  $\tau \leq 1/\Omega(\sqrt{L})$  by innovatively applying martingale inequalities. This bound is a bit surprising as a natural spectral bound on the ResNet mapping  $(1+1/\sqrt{L})^L$  explodes as L becomes large. This bound helps us establish that the forward/backward process is well-behaved. We then establish the upper/lower bounds on the gradient based on the forward/backward stability. Our gradient upper bound is  $\tau$ -related, much tighter than that in previous works, which enables the weak depth-dependent argument of learning ResNet. Especially this tighter gradient upper bound helps us achieve a depth-independent convergence result for the case  $\tau = 1/\Omega(L\log m)$ . We then show the gradient bounds do not change much after perturbation as long as the perturbation is relatively small. Putting the above properties together, we show the smoothness property of the objective function and establish the convergence of gradient descent with  $\tau \leq 1/\Omega(\sqrt{L}\log m)$  for learning over-parameterized ResNet.

#### 1.1 Related works

Several papers have argued the benefit of ResNet but they are either lack of rigorous theory or study the ResNet with linear activation. Specifically, Veit et al. (2016) interpreted ResNet as an ensemble of shallower networks, which is imprecise because the shallower networks are trained jointly, not independently (Xie et al., 2017). Zhang et al. (2018) argued the benefit of skip connection from the perspective of improving the local Hessian and Hardt and Ma (2016) showed that deep linear residual networks have no spurious local optima.

Our paper is closely related to recent work on over-parameterization technique. Specifically on the optimization side, Allen-Zhu et al. (2018b); Du et al. (2018); Chizat and Bach (2018) showed that gradient descent converges linearly to the global minima for training over-parameterized deep neural network. They also extended the analysis for training over-parameterized ResNet. Du et al. (2018) showed that the width/over-parameterization requirement for training ResNet is polynomial with the number of layers in contrast with exponential depth dependence for vanilla feedforward network. Nonetheless, Allen-Zhu et al. (2018b) showed that the width requirements for both feedforward network and ResNet are polynomial with the depth.

Moreover, on the generalization side, Brutzkus et al. (2017) provided the optimization and generalization guarantees of the SGD solution for over-parameterized two-layer networks given that the data is linear separable. Li and Liang (2018); Allen-Zhu et al. (2018a); Arora et al. (2019); Cao and Gu (2019) showed that the over-parameterized neural network provably generalizes for two-layer or multiple-layer networks. Neyshabur et al. (2019) used unit-wise capacity and obtained a bound on the empirical Rademacher complexity, which gave an explanation (not rigorous argument) of the generalization for over-parameterized two-layer ReLU networks.

#### 1.2 Paper Organization

The rest of this paper is organized as follows. Section 2 introduces the model and notations. Section 3 presents the main results. Section 4 discusses the proof roadmap and the main challenges. Section 5 gives some experiments to further illustrate our theory. Finally, we conclude in Section 6.

#### 2 Preliminaries

There are many residual network models since the seminal paper He et al. (2016). Here we study a simple version<sup>1</sup> to ease the presentation. It is sufficient to understand how ResNet helps the optimization. The ResNet model is described as follows,

- Input layer:  $h_0 = \phi(\mathbf{A}x)$ ;
- L-1 residual layers:  $h_l = \phi(h_{l-1} + \tau W_l h_{l-1})$ , for l = 1, 2, ..., L-1;
- A fully-connected layer:  $h_L = \phi(\boldsymbol{W}_L h_{L-1});$
- Output layer:  $y = Bh_L$ ;

<sup>&</sup>lt;sup>1</sup> The same ResNet model has been used in Allen-Zhu et al. (2018b) and Du et al. (2018). We borrow some notations from Allen-Zhu et al. (2018b).

where  $\phi(\cdot)$  is the ReLU activation, i.e.,  $\phi(\cdot) := \max\{0, \cdot\}$ , and  $\tau$  is the scaling factor on the parametric branch in a residual block which will be specified later. We assume the input dimension is p and hence  $x \in \mathbb{R}^p$ , the intermediate layers have the same width m, and hence  $h_l \in \mathbb{R}^m$  for l = 0, 1, ..., L, and the output has dimension d and hence  $y \in \mathbb{R}^d$ . Denote the values before activation by  $g_0 = Ax, g_l = h_{l-1} + \tau W_l h_{l-1}$  for l = 1, 2, ..., L-1 and  $g_L = W_L h_{L-1}$ . Use  $h_{i,l}$  and  $g_{i,l}$  to denote the value of  $h_l$  and  $g_l$ , respectively, when the input vector is  $x_i$ . Let  $D_{i,l}$  be the diagonal sign matrix where  $[D_{i,l}]_{k,k} = \mathbf{1}_{\{(g_{i,l})_k \geq 0\}}$ .

We adopt the following initialization scheme:

- Each entry of  $\mathbf{A} \in \mathbb{R}^{m \times p}$  is sampled independently from  $\mathcal{N}(0, \frac{2}{m})$ ;
- Each entry of  $W_l \in \mathbb{R}^{m \times m}$  is sampled independently from  $\mathcal{N}(0, \frac{2}{m})$  for l = 1, 2, ..., L;
- Each entry of  $\mathbf{B} \in \mathbb{R}^{d \times m}$  is sampled independently from  $\mathcal{N}(0, \frac{2}{d})$ .

The training data set is  $\{(x_i, y_i^*)\}_{i=1}^n$ , where  $x_i$  is the feature vector and  $y_i^*$  is the target signal for all i = 1, ..., n. We make the following assumption on the training data.

**Assumption 1.** For every  $i \in [n]$ ,  $||x_i|| = 1$ . For every pair  $i, j \in [n]$ ,  $||x_i - x_j|| \ge \delta$ .

We consider  $\ell_2$  regression task and the objective function is

$$F(\overrightarrow{\boldsymbol{W}}) := \sum_{i=1}^{n} F_i(\overrightarrow{\boldsymbol{W}}) \text{ where } F_i(\overrightarrow{\boldsymbol{W}}) := \frac{1}{2} \|\boldsymbol{B}h_{i,L} - y_i^*\|^2,$$

where  $\overrightarrow{\boldsymbol{W}}:=(\boldsymbol{W}_1,\boldsymbol{W}_2,\ldots,\boldsymbol{W}_L)$  are the trainable parameters. We clarify some notations here. We use  $\|v\|$  to denote the  $l_2$  norm of the vector v. We further use  $\|\boldsymbol{M}\|_2$  and  $\|\boldsymbol{M}\|_F$  to denote the spectral norm and the Frobenius norm of the matrix  $\boldsymbol{M}$ , respectively. Denote  $\|\overrightarrow{\boldsymbol{W}}\|_2:=\max_{l\in[L]}\|\boldsymbol{W}_l\|_2$  and  $\|\boldsymbol{W}_{[L-1]}\|_2:=\max_{l\in[L-1]}\|\boldsymbol{W}_l\|_2$ .

We note that the initialization scheme and the assumption on the data are the same as those in Allen-Zhu et al. (2018b) so the result is comparable. The model is trained by running the gradient descent algorithm. The gradient is computed through back-propagation.

## 3 Full Characterization of Learning ResNet

Given the model introduced in Section 2, we state the main result for gradient descent as follows.

**Theorem 1.** For the ResNet defined and initialized as in Section 2 with  $\tau \leq 1/\Omega(\sqrt{L}\log m)$ , if the network width  $m \geq \max\{L, \Omega(n^{24} \max\{(\tau L)^{14}, 1\}\delta^{-8}d\log^2 m)\}$ , then with probability at least  $1 - \exp(-\Omega(\log^2 m))$ , gradient descent with learning rate  $\eta = \Theta(\frac{d\delta}{n^4 m})$  finds a point  $F(\overline{\mathbf{W}}) \leq \varepsilon$  in  $T = \Omega(n^6 \delta^{-2} \log \frac{n \log^2 m}{\varepsilon})$  iterations.

*Proof.* A proof roadmap is given in Section 4 and full proof is deferred to the supplemental material in Appendix F.

This theorem establishes the linear convergence of gradient descent for learning ResNet with the scaling factor  $\tau \leq 1/\Omega(\sqrt{L}\log m)$ . We have two corollaries with specific values of  $\tau$ .

Corollary 1. For  $\tau = 1/\Omega(\sqrt{L}\log m)$ ,  $m \ge \Omega(n^{24}L^7\delta^{-8}d\log^2 m)$ } is sufficient to guarantee global convergence of gradient descent for learning ResNet with arbitrary depth.

The scaling factor  $\tau$  that guarantees convergence of learning over-parameterized ResNet in our result is much larger than that  $\tau = 1/\Omega(L\log m)$  in previous work Allen-Zhu et al. (2018b); Du et al. (2018). The range of  $\tau$  for learning ResNet with arbitrary depth has been greatly enlarged.

Corollary 2. For the case of  $\tau \leq 1/\Omega(L\log m)$ ,  $m \geq \max\{L, \Omega(n^{24}\delta^{-8}d\log m)\}$  is sufficient to guarantee global convergence of gradient descent for learning ResNet with arbitrary depth.

This indicates that for  $\tau \leq 1/\Omega(L\log m)$ , the convergence and over-parameterization requirement is almost depth-independent for learning ResNet. Thus training deep residual network can be as easy as training a two-layer network, for the case of  $\tau \leq 1/\Omega(L\log m)$ . This result is much stronger than that under the same assumption in previous work Allen-Zhu et al. (2018b) and Du et al. (2018) which only establish the polynomial dependence of the over-parameterization on the network depth. This also theoretically justifies the advantage of ResNet over vanilla feedforward network in terms of facilitating the convergence of gradient descent. Moreover, we will argue that even for  $\tau = 1/\Omega(\sqrt{L}\log m)$ , the depth dependence for learning over-parameterized ResNet is weak compared with that for learning feedforward network in next section.

We note that the proof of Theorem 1 relies on the value of  $\tau$ . The larger  $\tau$ , the harder the convergence, the better the expressiveness of the network. However, one may wonder if  $\tau = 1/\Omega(\sqrt{L}\log m)$  is the largest allowable value to guarantee the convergence of learning ResNet.

We next argue the tightness by showing that if  $\tau \geq L^{-\frac{1}{2}+c}$ , the network output will highly likely explode for large L. We note that L represents the number of residual blocks for other residual models.

**Theorem 2.** For the ResNet defined and initialized as in Section 2, if  $\tau \geq L^{-\frac{1}{2}+c}$ , then in expectation we have

$$\mathbb{E}||h_L||^2 > L^{2c}.\tag{2}$$

*Proof.* The proof is relegated to the supplemental material in Appendix G.  $\Box$ 

This indicates the bound of  $\tau$  in Theorem 1 for successful learning ResNet is tight up to a logarithmic factor. In the following, we first explain the roadmap and the challenges for proving Theorem 1. We then give key steps to tackle the challenges. Finally we present a proof outline for Theorem 1.

## 4 Proof Roadmap and Key Steps

#### 4.1 Proof Roadmap and Main Challenges

Obviously, it is not expected that the objective function of ResNet has any good global property for optimization because of the nonconvexity and non-smoothness of the objective. However, it is possible to characterize some good properties locally or at least over the optimization path. In the sequel, all the optimization properties built for ResNet holds locally within a neighborhood of initialization. In order to show the linear convergence of gradient descent for learning ResNet, three properties are established locally that gradient is upper bounded, gradient is lower bounded, i.e., the gradient is large when the objective is large, and the objective satisfies certain smoothness property.

For ResNet, the gradient with respect to the parameter is computed through back-propagation, e.g.,  $\partial \mathbf{W}_l = \partial h_l \cdot h_{l-1}^T$ , where  $\partial \cdot$  represents the gradient of the objective with respect to  $\cdot$ . Thus, the gradient upper bound is guaranteed if  $h_l$  and  $\partial h_l$  are bounded across layers and iterations. Therefore, we first show the forward/backward process is bounded at the initialization stage and after small perturbation. The **gradient upper bound** is presented in Theorem 3 in Section 4.2. Based on bounded forward/backward process, the gradient for each individual sample can be shown to be lower bounded. The sum of individual gradients still being large can be argued by the smoothed analysis as in (Allen-Zhu et al., 2018b). The **gradient lower bound** is presented as Theorem 5 in Appendix D.2 in the supplemental material. The **objective smoothness** can also be argued by using the gradient upper bound and the bounded forward/backward process, which is presented in Theorem 4.

One main challenge is to show the forward/backward process of ResNet is bounded even when  $\tau$  is as large as  $1/\Omega(\sqrt{L})$ . We note that previous work establishing the convergence of learning ResNet for a much smaller  $\tau = 1/\Omega(L\log m)$  heavily relies on a natural spectral norm bound  $\|(\boldsymbol{I} + \tau \boldsymbol{W}_L) \cdots (\boldsymbol{I} + \tau \boldsymbol{W}_1)\| \le (1 + \frac{1}{L})^L$ . However, for  $\tau = 1/\Omega(\sqrt{L})$  the natural spectral bound  $(1 + \frac{1}{\sqrt{L}})^L$  explodes. Moreover, it also needs to show the  $\|h_{i,l}\|$  is lower bounded for  $\tau = 1/\Omega(\sqrt{L})$  with high probability. We have detailed discussion on these bounds in Section 4.2.

The other challenge is to theoretically justify the advantage of learning ResNet: weak dependence on the network depth, which is echoed by empirical evidence that deep ResNet is much easier to train than deep feedforward network. This advantage is not reflected clearly in previous work. In this paper, we show that the stability of the forward/backward process of ResNet is depth-independent for with  $\tau \leq 1/\Omega(\sqrt{L})$ . Based on that, we show the gradient of the parameter in the residual block is scaled down by the factor  $\tau$ . Utilizing this fact, we derive a new smoothness property for ResNet, where the dominant term is depth-independent. Specifically for  $\tau \leq 1/L$ , we show the whole smoothness property is depth-independent which is essential to the claim in Corollary 2. We give detailed argument on the weak depth dependence of learning ResNet in Section 4.2.

#### 4.2 Key Steps to Solve Main Challenges

Here we show the key steps that solve the two challenges on the way of establishing Theorem 1. Forward/backward process is bounded for  $\tau = 1/\Omega(\sqrt{L})$ .

The first and important step is the following lemma on the spectral norm bound at initialization.

**Lemma 1.** Suppose that  $\overrightarrow{W}^{(0)}$ , A are randomly generated as in the initialization step, and  $D_0, \ldots, D_L$  are diagonal matrices such that  $\|D_l\|_2 \leq 1$  for all  $l \in [L]$  and  $D_l$  is independent from the randomness  $W_a^{(0)}$  for all a > l. Then with probability at least  $1 - L^2 \cdot \exp(-\Omega(mc^2))$  over the initialization randomness we have

$$\left\| \boldsymbol{D}_{b} \left( \boldsymbol{I} + \tau \boldsymbol{W}_{b}^{(0)} \right) \boldsymbol{D}_{b-1} \cdots \boldsymbol{D}_{a} \left( \boldsymbol{I} + \tau \boldsymbol{W}_{a}^{(0)} \right) \right\|_{2} \leq 1 + c, \tag{3}$$

where c is some small constant determined by the value of  $\tau \leq 1/\Omega(\sqrt{L})$ .

The above result is a bit surprising since for  $\tau = 1/\Omega(\sqrt{L})$  a natural spectral bound  $\|(\boldsymbol{I} + \tau \boldsymbol{W}_L^{(0)}) \cdots (\boldsymbol{I} + \tau \boldsymbol{W}_1^{(0)})\| \le (1 + \frac{1}{\sqrt{L}})^L$  explodes. Here the main idea is to utilize the independent randomness of matrices  $\boldsymbol{W}_l^{(0)}$  so that the cross-product term has mean 0. We develop an argument by constructing a martingale sequence.

Proof Outline. Suppose we have  $||h_{a-1}|| = 1$ . Then we abuse notations  $g_l = h_{l-1} + \tau \mathbf{W}_l^{(0)} h_{l-1}$  and  $h_l = \mathbf{D}_l g_l$  for  $a \le l \le b$ , and we have

$$||h_b||^2 = \frac{||h_b||^2}{||h_{b-1}||^2} \cdots \frac{||h_a||^2}{||h_{a-1}||^2} ||h_{a-1}||^2 \le \frac{||g_b||^2}{||h_{b-1}||^2} \cdots \frac{||g_a||^2}{||h_{a-1}||^2} ||h_{a-1}||^2.$$

Taking logarithm at both side, we have

$$\log \|h_b\|^2 \le \sum_{l=a}^b \log \Delta_l$$
, where  $\Delta_l := \frac{\|g_l\|^2}{\|h_{l-1}\|^2}$ .

If let  $\tilde{h}_{l-1} := \frac{h_{l-1}}{\|h_{l-1}\|}$ , then we obtain that

$$\log \Delta_{l} = \log \left( 1 + 2\tau \left\langle \tilde{h}_{l-1}, \boldsymbol{W}_{l}^{(0)} \tilde{h}_{l-1} \right\rangle + \tau^{2} \|\boldsymbol{W}_{l}^{(0)} \tilde{h}_{l-1}\|^{2} \right) \leq 2\tau \left\langle \tilde{h}_{l-1}, \boldsymbol{W}_{l}^{(0)} \tilde{h}_{l-1} \right\rangle + \tau^{2} \|\boldsymbol{W}_{l}^{(0)} \tilde{h}_{l-1}\|^{2},$$

where the inequality is because  $\log(1+x) < x$  for x > -1. Let  $\xi_l := 2\tau \left\langle \tilde{h}_{l-1}, \boldsymbol{W}_l^{(0)} \tilde{h}_{l-1} \right\rangle$  and  $\zeta_l := \tau^2 \|\boldsymbol{W}_l^{(0)} \tilde{h}_{l-1}\|^2$ , for given  $\tilde{h}_{l-1}, \xi_l \sim \mathcal{N}\left(0, \frac{8\tau^2}{m}\right), \zeta_l \sim \frac{2\tau^2}{m} \chi_m^2$ .

Without rigor, we could say  $\sum_{l=a}^{b} \xi_l \sim \mathcal{N}\left(0, \frac{8(b-a)\tau^2}{m}\right)$  and  $\sum_{l=a}^{b} \zeta_l \sim \frac{2(b-a)\tau^2}{m}\chi_m^2$ . Hence we have  $\sum_{l=a}^{b} \log \Delta_l \leq c$  with probability at least  $1 - \exp(-\Omega(m)c^2)$ . Taking  $\varepsilon$ -net argument, we can establish the spectral norm bound for all vector  $h_{a-1}$ . Let a and b vary from 1 to L-1 and taking the union bound gives the claim.

We next show the norm of the representation is close to 1 at every layer.

**Lemma 2.** With probability at least  $1 - O(nL) \cdot e^{-\Omega(m)}$  over the randomness of  $\mathbf{A} \in \mathbb{R}^{m \times p}$  and  $\overrightarrow{\mathbf{W}}^{(0)} \in (\mathbb{R}^{m \times m})^L$ , we have

$$\forall i \in [n], l \in \{0, 1, \dots, L\} : \|h_{i,l}^{(0)}\| \in [1 - c, 1 + c], \tag{4}$$

where c can be arbitrarily small for  $\tau \leq 1/\Omega(\sqrt{L})$ .

The proof is relegated to Appendix C.1. We note that Lemma 2 is much stronger compared with the result in Allen-Zhu et al. (2018b) which is only showed for the case  $\tau = 1/\Omega(L \log m)$  and cannot guarantee  $||h_{i,l}||$  arbitrarily close to 1. The property of  $||h_{i,l}||$  arbitrarily close to 1 is required for down-streaming bounding tasks. Moreover we note that the above two lemmas also holds for  $\overrightarrow{W}$  that is within the neighborhood of  $\overrightarrow{W}^{(0)}$  and the result is presented in Appendix C.2.

#### Convergence weakly depends on network depth for learning ResNet.

We next justify that the convergence of learning ResNet weakly depends on the network depth. We first present a new gradient upper bound that is  $\tau$ -related and specific for ResNet.

**Theorem 3.** With probability at least  $1 - \exp(-\Omega(m))$  over the randomness of  $\overrightarrow{W}^{(0)}$ , A, B, it satisfies for every  $l \in [L-1]$ , every  $i \in [n]$ , and every  $\overrightarrow{W}$  with  $\|\overrightarrow{W} - \overrightarrow{W}^{(0)}\|_2 \le \omega$  for  $\omega \in [0,1]$ ,

$$\|\nabla_{\boldsymbol{W}_{l}}F(\overrightarrow{\boldsymbol{W}})\|_{F}^{2} \leq O\left(\frac{F(\overrightarrow{\boldsymbol{W}})}{d} \times \tau^{2}mn\right), \qquad \|\nabla_{\boldsymbol{W}_{L}}F(\overrightarrow{\boldsymbol{W}})\|_{F}^{2} \leq O\left(\frac{F(\overrightarrow{\boldsymbol{W}})}{d} \times mn\right). \tag{5}$$

The full proof is relegated to Appendix D.1. Here we give an outline.

*Proof Outline*. The argument is based on the bounded forward/backward process at  $\overrightarrow{W}$  and the back-propagation formula. We focus only the residual layer here. For each  $i \in [n]$  and  $l \in [L-1]$ , we have

$$\|\nabla_{\boldsymbol{W}_{l}}F_{i}(\overrightarrow{\boldsymbol{W}})\|_{F} = \tau \left(\boldsymbol{D}_{i,l}(\boldsymbol{I} + \tau \boldsymbol{W}_{l+1})^{T} \cdots \boldsymbol{D}_{i,L-1} \boldsymbol{W}_{L}^{T} \boldsymbol{D}_{i,L} \boldsymbol{B}^{T} \left(\boldsymbol{B} h_{i,L} - y_{i}^{*}\right)\right) h_{i,l-1}^{T}$$

$$\leq O(\tau \sqrt{m/d}) \sqrt{F_{i}(\overrightarrow{\boldsymbol{W}})}.$$

The above upper bound holds because of the bounded forward/backward process for all the  $\overrightarrow{\boldsymbol{W}}$  such that  $\|\overrightarrow{\boldsymbol{W}} - \overrightarrow{\boldsymbol{W}}^{(0)}\|_2 \leq \omega$  for some  $\omega$  (this result is Lemma 5 in supplemental material).

This gradient upper bound is tighter than Allen-Zhu et al. (2018b) by involving  $\tau$  for the residual layers. It treats the top layer  $\mathbf{W}_L$  and the residual layers  $\mathbf{W}_l$  for  $l \in [L-1]$  separately. This upper bound helps us improve the smoothness property and is central to show a weak dependence on the network depth for learning ResNet.

We next present an informal semi-smoothness<sup>2</sup> result to illustrate why the convergence of learning ResNet could weakly depend on the network depth.

**Theorem 4** (Informal semi-smoothness result). Let  $\omega < 1$  and  $\tau^2 L \leq 1$ . With high probability, we have for every  $\overrightarrow{\mathbf{W}} \in (\mathbb{R}^{m \times m})^L$  with  $\| \widecheck{\mathbf{W}}_L - \mathbf{W}_L^{(0)} \|_2 \leq \omega$  and  $\| \widecheck{\mathbf{W}}_l - \mathbf{W}_l^{(0)} \|_2 \leq \tau \omega$  for  $l \in [L-1]$ , and for every  $\overrightarrow{\mathbf{W}}' \in (\mathbb{R}^{m \times m})^L$  with  $\| \mathbf{W}'_L \|_2 \leq \omega$  and  $\| \mathbf{W}'_l \|_2 \leq \tau \omega$  for  $l \in [L-1]$ , we have

$$F(\overset{\smile}{\overrightarrow{W}} + \overset{\smile}{\overrightarrow{W}}') \leq F(\overset{\smile}{\overrightarrow{W}}) + \langle \nabla F(\overset{\smile}{\overrightarrow{W}}), \overset{\smile}{\overrightarrow{W}}' \rangle + O\left(\frac{nm}{d}\right) \left( \| \boldsymbol{W}_{L}' \|_{2} + \tau \sum_{l=1}^{L-1} \| \boldsymbol{W}_{l}' \|_{2} \right)^{2} + \sqrt{\frac{mn\omega^{2/3}}{d}} \cdot O\left( \| \boldsymbol{W}_{L}' \|_{2} + \max\{(\tau L)^{4/3}, 1\} \sum_{l=1}^{L-1} \| \boldsymbol{W}_{l}' \|_{2} \right) \sqrt{F(\overset{\smile}{\overrightarrow{W}})}.$$
(6)

We note that apart from the second-order term (the third term on the right hand side of (6)) in classical smoothness, the semi-smoothness (6) also has a first-order term (the last term on the right hand side of (6)). One can see that as m becomes large in the over-parameterization regime the effect of the first-order term becomes small comparing to the second-order term. Interestingly, if one replaces  $\mathbf{W}'_L$  and  $\mathbf{W}'_l$  with the gradient upper bounds in Theorem 3, only the first-order term depends on L while the second-order term, which is dominant when m is large, is depth-independent. We note that the first-order term is the only source where the depth dependence in Theorem 1 comes from, which renders the depth dependence is weak for learning ResNet when m is large, in contrast with the case of learning deep feedforward network Allen-Zhu et al. (2018b). Furthermore, we note that if  $\tau < 1/L$  the whole right hand side of (6) is depth-independent and hence we have a depth-independent convergence in Corollary 2.

#### 4.3 Outline Proof of Theorem 1

For better presentation, we hide the exposure of intermediate results and state the gradient lower bound as follows, for  $\|\overrightarrow{W} - \overrightarrow{W}^{(0)}\|_2 \le \omega$  and  $\|\overrightarrow{W}'\|_2 \le \omega$ ,

$$\|\nabla_{\boldsymbol{W}}F(\overrightarrow{\boldsymbol{W}})\|_F^2 \ge \Omega\left(\frac{m\delta}{dn^2}\right)F(\overrightarrow{\boldsymbol{W}}). \tag{7}$$

<sup>&</sup>lt;sup>2</sup> The smoothness is compromised from the usual sense because of the non-smoothness of ReLU.

Based on the forward stability and the randomness of  $\boldsymbol{B}$ , we can show that  $\|\boldsymbol{B}h_{i,L}^{(0)} - y_i^*\|^2 \leq O(\log^2 m)$  with high probability and therefore  $F(\overrightarrow{\boldsymbol{W}}^{(0)}) \leq O(n\log^2 m)$ .

Assume that for every t = 0, 1, ..., T - 1,  $\overrightarrow{\boldsymbol{W}}^{(t)}$  stay within a neighborhood of the initialization,

$$\|\boldsymbol{W}_{L}^{(t)} - \boldsymbol{W}_{L}^{(0)}\|_{F} \le \omega \stackrel{\Delta}{=} O\left(\frac{\delta^{3}}{n^{9}(\tau L)^{7}}\right), \quad \|\boldsymbol{W}_{l}^{(t)} - \boldsymbol{W}_{l}^{(0)}\|_{F} \le \tau \omega.$$
 (8)

Then for one gradient descent step  $\overrightarrow{\boldsymbol{W}}^{(t+1)} = \overrightarrow{\boldsymbol{W}}^{(t)} - \eta \nabla F(\overrightarrow{\boldsymbol{W}}^{(t)})$ , based on (6) we have

$$F(\overrightarrow{\boldsymbol{W}}^{(t+1)}) \le \left(1 - \Omega\left(\frac{\eta\delta m}{dn^2}\right)\right) F(\overrightarrow{\boldsymbol{W}}^{(t)}),$$
 (9)

where the inequality uses the gradient lower bound in (7) and the choice of  $\eta = \frac{d\delta}{n^4m}$  and the assumption on  $\omega$  in (8). That is, after  $T = \Omega(n^6\delta^{-2})\log\frac{n\log^2 m}{\varepsilon}$  iterations,  $F(\overrightarrow{\boldsymbol{W}}^{(T)}) \leq \varepsilon$ .

We need to verify for each t, the iterate  $\overrightarrow{W}^{(t)}$  stays in the region given by (8). Therefore, we calculate

$$\|\boldsymbol{W}_{L}^{(t)} - \boldsymbol{W}_{L}^{(0)}\|_{F} \leq \sum_{i=1}^{t-1} \|\eta \nabla_{\boldsymbol{W}_{L}} F(\overrightarrow{\boldsymbol{W}}^{(i)})\|_{F} \stackrel{(a)}{\leq} O(\eta \sqrt{\frac{nm}{d}}) \sum_{i=1}^{t-1} \sqrt{F(\overrightarrow{\boldsymbol{W}}^{(i)})} \stackrel{(b)}{\leq} O\left(\frac{n^{3} \sqrt{d}}{\delta \sqrt{m}}\right)$$
(10)

where (a) is due to Theorem 3 and (b) is due to an upper bound of the sum of a geometric sequence. Similarly, we have for  $l \in [L-1]$ ,  $\|\boldsymbol{W}_{l}^{(t)} - \boldsymbol{W}_{l}^{(0)}\|_{F} \leq O\left(\frac{\tau n^{3}\sqrt{d}}{\delta\sqrt{m}}\right)$ .

By combining (10) and the assumption on  $\omega$  (8), we obtain a bound on m.

## 5 Empirical study

Though the focus in this paper is on the theoretical analysis, we present some experiments to illustrate the theoretical claim as a proof-of-concept. Specifically, we train a feedforward fully-connected neural network and ResNets with different values of  $\tau$ , and compare their convergence behavior.

The feedforward model is composed of an input layer:  $h_0 = \phi(\mathbf{A}x)$ ; and then L feedforward layers:  $h_l = \phi(\mathbf{W}_l h_{l-1})$ , for l = 1, 2, ..., L; and finally an output layer:  $y = \mathbf{B}h_L$ . The feedforward model adopts the same initialization scheme as the ResNet model (see Section 2). The models are generated by varying the depth  $L \in \{3, 10, 30, 100, 500, 1000\}$  and the width m = 128. We conduct experiments with different settings of  $\tau$  and show how it affects the training performance. Specifically,  $\tau$  is chosen from  $\{\frac{1}{L}, \frac{1}{L^{0.25}}, \frac{1}{L^{0.25}}\}$ .

We use MNIST dataset (LeCun et al., 1998) and do the classification task. MNIST contains 60000 training examples with input dimension 784 and 10 labels. The input feature vector is normalized by subtracting feature mean and dividing the feature standard deviation. We train the model with SGD<sup>3</sup> and the size of minibatch is 256. The learning rate lr is set as lr = 0.001 without heavily tuning.

**Experiment results.** We plot the training curves of feedforward network and ResNet with varying depths and widths in Figure 1.

update something We can see that both  $\tau = \frac{1}{L}$  and  $\tau = \frac{1}{L^{0.5}}$  are able to guarantee the successful training of very deep ResNets. However, for  $\tau = \frac{1}{L^{0.25}}$ , the training loss explodes for models with

<sup>&</sup>lt;sup>3</sup> We test GD for small models and observe the same phenomenon. We use SGD for all experiments due to the expensive per-iteration cost of GD.

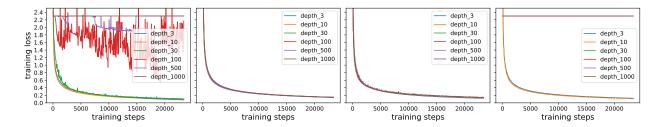


Figure 1: Training curve comparison for feedforward network and ResNet with varying  $\tau$ . The columns from left to right corresponds to feedforward network and ResNets with  $\tau = \frac{1}{L}, \frac{1}{L^{0.5}}$  and  $\frac{1}{L^{0.25}}$ , respectively.

depth 30 and more. This indicates that the bound  $\tau = \frac{1}{\sqrt{L}}$  is sharp for learning ResNet with arbitrary depth.

Moreover, Figures 1 shows that for a given width, the convergence of ResNet with  $\tau = \frac{1}{L}$  and  $\tau = \frac{1}{L^{0.5}}$  does not depend on the depth much while training feedforward network becomes much harder as the depth increases. This is well justified by the weak dependence on depth for learning deep ResNet.

#### 6 Conclusion

In this paper, we establish the convergence of gradient descent for learning over-parameterized ResNet, which is much stronger than previous work. First our convergence holds for the range of  $\tau \leq 1/\Omega(\sqrt{L}\log m)$  while that in previous work holds only for  $\tau = 1/\Omega(L\log m)$ , where  $\tau$  is a scaling factor on the residual branch. Moreover, specifically for  $\tau = 1/\Omega(L\log m)$  our convergence result does not depend on the network depth which is in sharp contrast with the polynomial dependence of network depth in previous work under the same assumption. This bridges the gap between theoretical understanding and practical advantage of ResNet structure. Conversely we show the achievable range of  $\tau$  is tight by arguing that the forward process explodes for large L if  $\tau > L^{-\frac{1}{2}}$ . This gives a full picture of learning over-parameterized ResNet. One interesting point for future research is studying the range of  $\tau$  when batch normalization (Ioffe and Szegedy, 2015) is involved. It is not clear how the proof adapts for network with batch normalization. Another direction is to bypass the semi-smooth argument and give a sharp dependence on the depth.

## Acknowledgments

The authors would like to thank Prof. Yingbin Liang for many helpful discussions and thank Zeyuan Allen-Zhu for clarifying the proof in Allen-Zhu et al. (2018b) and discussion.

#### References

#### References

- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. arXiv preprint arXiv:1811.04918, 2018a.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. arXiv preprint arXiv:1811.03962, 2018b.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. arXiv preprint arXiv:1901.08584, 2019.
- David Balduzzi, Marcus Frean, Lennox Leary, JP Lewis, Kurt Wan-Duo Ma, and Brian McWilliams. The shattered gradients problem: If resnets are the answer, then what is the question? In *International Conference on Machine Learning*, pages 342–350, 2017.
- Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. Sgd learns overparameterized networks that provably generalize on linearly separable data. arXiv preprint arXiv:1710.10174, 2017.
- Yuan Cao and Quanquan Gu. A generalization theory of gradient descent for learning overparameterized deep relu networks. arXiv preprint arXiv:1902.01384, 2019.
- Lénaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In Advances in Neural Information Processing Systems 31. 2018.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. arXiv preprint arXiv:1811.03804, 2018.
- Moritz Hardt and Tengyu Ma. Identity matters in deep learning. In *International Conference on Learning Representations*, 2016.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, pages 448–456, 2015.
- Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, pages 8168–8177, 2018.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. The role of over-parametrization in generalization of neural networks. In *International Conference on Learning Representations*, 2019.
- A Emin Orhan and Xaq Pitkow. Skip connections eliminate singularities. In *International Conference* on Learning Representations (ICLR), 2018.
- Daniel A Spielman and Shang-Hua Teng. Smoothed analysis of algorithms: Why the simplex algorithm usually takes polynomial time. *Journal of the ACM (JACM)*, 51(3):385–463, 2004.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- Andreas Veit, Michael J Wilber, and Serge Belongie. Residual networks behave like ensembles of relatively shallow networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 550–558, 2016.
- R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. Compressed Sensing, Theory and Applications, pages 210 268, 2012.
- Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5987–5995. IEEE, 2017.
- Huishuai Zhang, Wei Chen, and Tie-Yan Liu. On the local hessian in back-propagation. In Advances in Neural Information Processing Systems (NeurIPS), pages 6521–6531. 2018.
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks. arXiv preprint arXiv:1811.08888, 2018.

#### A Useful Lemmas

First we list several useful bounds on Gaussian distribution.

**Lemma 3.** Suppose  $X \sim \mathcal{N}(0, \sigma^2)$ , then

$$\mathbb{P}\{|X| \le x\} \ge 1 - \exp\left(-\frac{x^2}{2\sigma^2}\right),\tag{11}$$

$$\mathbb{P}\{|X| \le x\} \le \sqrt{\frac{2}{\pi}} \frac{x}{\sigma}.\tag{12}$$

Another bound is on the spectral norm of random matrix (Vershynin, 2012, Corollary 5.35).

**Lemma 4.** Let  $A \in \mathbb{R}^{N \times n}$ , and entries of A are independent standard Gaussian random variables. Then for every  $t \geq 0$ , with probability at least  $1 - \exp(-t^2/2)$  one has

$$s_{\max}(\mathbf{A}) \le \sqrt{N} + \sqrt{n} + t,\tag{13}$$

where  $s_{\text{max}}(\mathbf{A})$  are the largest singular value of  $\mathbf{A}$ .

## B Spectral Norm Bound at Initialization

Next we present a spectral norm bound with  $\tau \leq 1/\Omega(\sqrt{L})$  related to ResNet.

**Lemma 1.** Suppose that  $\overrightarrow{W}^{(0)}$ , A are randomly generated as in the initialization step, and  $D_0, \ldots, D_L$  are diagonal matrices such that  $\|D_l\|_2 \leq 1$  for all  $l \in [L]$  and  $D_l$  is independent from the randomness  $W_a^{(0)}$  for all a > l. Then with probability at least  $1 - L^2 \cdot \exp(-\Omega(mc^2))$  over the initialization randomness we have

$$\left\| \boldsymbol{D}_{b} \left( \boldsymbol{I} + \tau \boldsymbol{W}_{b}^{(0)} \right) \boldsymbol{D}_{b-1} \cdots \boldsymbol{D}_{a} \left( \boldsymbol{I} + \tau \boldsymbol{W}_{a}^{(0)} \right) \right\|_{2} \leq 1 + c, \tag{3}$$

where c is some small constant determined by the value of  $\tau \leq 1/\Omega(\sqrt{L})$ .

*Proof.* We first show for any vector  $h_{a-1}$  with  $||h_{a-1}|| = 1$ , we have  $||h_b|| \le 1 + c$  with high probability, where

$$h_b = D_b(I + \tau W_b^{(0)}) D_{b-1} \cdots D_a(I + \tau W_a^{(0)}) h_{a-1}.$$
(14)

Using notations  $g_l = h_{l-1} + \tau W_l h_{l-1}$  introduced in Section 2, we have  $||g_l|| \ge ||h_l||$ . Thus we have the following

$$||h_b||^2 = \frac{||h_b||^2}{||h_{b-1}||^2} \cdots \frac{||h_a||^2}{||h_{a-1}||^2} ||h_{a-1}||^2 \le \frac{||g_b||^2}{||h_{b-1}||^2} \cdots \frac{||g_a||^2}{||h_{a-1}||^2} ||h_{a-1}||^2.$$

Taking logarithm at both side, we have

$$\log \|h_b\|^2 \le \sum_{l=a}^b \log \Delta_l, \quad \text{where } \Delta_l := \frac{\|g_l\|^2}{\|h_{l-1}\|^2}.$$
 (15)

If let  $\tilde{h}_{l-1} := \frac{h_{l-1}}{\|h_{l-1}\|}$ , then we obtain that

$$\log \Delta_{l} = \log \left( 1 + 2\tau \left\langle \tilde{h}_{l-1}, \boldsymbol{W}_{l}^{(0)} \tilde{h}_{l-1} \right\rangle + \tau^{2} \| \boldsymbol{W}_{l}^{(0)} \tilde{h}_{l-1} \|^{2} \right)$$

$$\leq 2\tau \left\langle \tilde{h}_{l-1}, \boldsymbol{W}_{l}^{(0)} \tilde{h}_{l-1} \right\rangle + \tau^{2} \| \boldsymbol{W}_{l}^{(0)} \tilde{h}_{l-1} \|^{2},$$

where the inequality is because  $\log(1+x) \leq x$  for all x > -1. Let  $\xi_l := 2\tau \left\langle \tilde{h}_{l-1}, \boldsymbol{W}_l^{(0)} \tilde{h}_{l-1} \right\rangle$  and  $\zeta_l := \tau^2 \|\boldsymbol{W}_l^{(0)} \tilde{h}_{l-1}\|^2$ , then given  $h_{l-1}$  we have  $\xi_l \sim \mathcal{N}\left(0, \frac{8\tau^2}{m}\right)$ ,  $\zeta_l \sim \frac{2\tau^2}{m} \chi_m^2$ .

We see that

$$\mathbb{P}\left(\sum_{l=a}^{b} \log \Delta_{l} \ge c\right) \le \mathbb{P}\left(\sum_{l=a}^{b} \xi_{l} \ge \frac{c}{2}\right) + \mathbb{P}\left(\sum_{l=a}^{b} \zeta_{l} \ge \frac{c}{2}\right). \tag{16}$$

Next we bound terms on the right hand side one by one. For the first term we have

$$\mathbb{P}\left(\sum_{l=a}^{b} \xi_{l} \geq \frac{c}{2}\right) = \mathbb{P}\left(\exp\left(\lambda \sum_{l=a}^{b} \xi_{l}\right) \geq \exp\left(\frac{\lambda c}{2}\right)\right) \leq \mathbb{E}\left[\exp\left(\lambda \sum_{l=a}^{b} \xi_{l} - \frac{\lambda c}{2}\right)\right], \quad (17)$$

where  $\lambda$  is any positive number and the last inequality uses the Markov's inequality. Moreover,

$$\mathbb{E}\left[\exp\left(\lambda \sum_{l=a}^{b} \xi_{l}\right)\right] = \mathbb{E}\left[\exp\left(\lambda \sum_{l=a}^{b-1} \xi_{l}\right) \mathbb{E}\left[\exp\left(\lambda \xi_{b}\right)\right] \middle| \mathcal{F}_{b-1}\right]$$

$$= \exp\left(\frac{8\tau^{2}\lambda^{2}}{m}\right) \mathbb{E}\left[\exp\left(\lambda \sum_{l=a}^{b-1} \xi_{l}\right)\right]$$

$$= \dots = \exp\left(\frac{8\tau^{2}\lambda^{2}(b-a+1)}{m}\right)$$
(18)

Hence we obtain

$$\mathbb{P}\left(\sum_{l=a}^{b} \xi_{l} \ge \frac{c}{2}\right) \le \exp\left(-\frac{mc^{2}}{128\tau^{2}(b-a+1)}\right) = \exp\left(-\Omega\left(\frac{mc^{2}}{\tau^{2}(b-a+1)}\right)\right) \tag{19}$$

by choosing  $\lambda = \frac{mc}{32\tau^2l}$  and  $\tau \leq 1/\Omega(\sqrt{L})$ . Due to the symmetry of  $\sum_{l=a}^b \xi_l$ , the conclusion can be generalize to the variable  $|\sum_{l=a}^b \xi_l|$ .

Then, for the second term, we first present a concentration inequality for the general  $\chi_m^2$  distribution X (Laurent and Massart, 2000, Lemma 1)

$$\mathbb{P}\left(|X-m| \ge u\right) \le e^{-\frac{u^2}{4m}}.\tag{20}$$

Then for  $\sum_{l=a}^{b} \zeta_l$ , by applying the above inequality and Jensen's inequality, we have

$$\mathbb{P}\left(\sum_{l=a}^{b} \zeta_{l} \geq \frac{c}{2}\right) = \mathbb{E}\left[\mathbb{P}\left(\sum_{l=a}^{b} \zeta_{l} \geq \frac{c}{2} \middle| \mathcal{F}_{b-1}\right)\right]$$

$$\leq \mathbb{E}\left[\mathbb{P}\left(\left|\zeta_{b} + \sum_{l=a}^{b-1} \zeta_{l} - 2\tau^{2}(b-a+1+1)\right| \geq \frac{c}{2} - 2\tau^{2}(b-a+1) \middle| \mathcal{F}_{b-1}\right)\right]$$

$$\leq \mathbb{E}\left[\mathbb{P}\left(\left|\zeta_{b} - 2\tau^{2}\right| \geq \frac{c}{2} - (4L-2)\tau^{2} - \sum_{k=a}^{b-1} \zeta_{k} \middle| \mathcal{F}_{b-1}\right)\right]$$

$$= \mathbb{E}\left[\mathbb{P}\left(\left|\frac{m}{2\tau^{2}}\zeta_{b} - m\right| \geq \frac{m}{2\tau^{2}}\left(\frac{c}{2} - (4L-2)\tau^{2} - \sum_{k=a}^{b-1} \zeta_{k}\right) \middle| \mathcal{F}_{b-1}\right)\right]$$

$$\leq \mathbb{E}\left[\exp\left(-\frac{m}{16\tau^{4}}\left(\frac{\delta}{2} - (4L-2)\tau^{2} - \sum_{k=a}^{b-1} \zeta_{k}\right)^{2}\right)\right]$$

$$\leq \mathbb{E}\left[\exp\left(-\frac{m}{16\tau^{4}}\left(\frac{c^{2}}{4} + \Omega(L^{2})\tau^{4} - \Omega(L)\tau^{2}\right)\right)\right]$$

$$= \exp\left(-\Omega\left(\frac{mc^{2}}{\tau^{4}}\right)\right)$$
(21)

Combining equation (21) and (19), and  $\tau \leq \frac{1}{\Omega(\sqrt{L})}$ , we obtain  $||h_b|| \leq 1 + c$  with probability at least  $1 - \exp(-\Omega(mc^2))$ . Taking  $\epsilon$ -net over all m-dimensional vectors of  $h_{a-1}$ , we have with probability  $1 - \exp(-\Omega(mc^2))$  the inequality (3) holds for any a and b with  $1 \leq a \leq b < L$ . Taking a union bound over a and b, the conclusion is proved.

## C Bounded Forward/Backward Process

#### C.1 Proof at Initialization

**Lemma 2.** With probability at least  $1 - O(nL) \cdot e^{-\Omega(m)}$  over the randomness of  $\mathbf{A} \in \mathbb{R}^{m \times p}$  and  $\overrightarrow{\mathbf{W}}^{(0)} \in (\mathbb{R}^{m \times m})^L$ , we have

$$\forall i \in [n], l \in \{0, 1, \dots, L\}: \|h_{i,l}^{(0)}\| \in [1 - c, 1 + c], \tag{4}$$

where c can be arbitrarily small for  $\tau \leq 1/\Omega(\sqrt{L})$ .

*Proof.* We ignore the subscript (0) for simplicity. The upper bound of  $||h_{i,l}||$  can be easily achieved by the proof of Lemma 1. Now, we give the lower bound of  $||h_{i,l}||$ . First we have

$$||h_{i,l}|| = ||h_{i,0}|| \frac{||h_{i,1}||}{||h_{i,0}||} \cdots \frac{||h_{i,l}||}{||h_{i,l-1}||}.$$
(22)

Then we see

$$\log \|h_{i,l}\|^2 = \log \|h_{i,0}\|^2 + \sum_{a=1}^{l} \log \frac{\|h_{i,a}\|^2}{\|h_{i,a-1}\|^2}$$

$$= \log \|h_{i,0}\|^2 + \sum_{a=1}^{l} \log \left(1 + \frac{\|h_{i,a}\|^2 - \|h_{i,a-1}\|^2}{\|h_{i,a-1}\|^2}\right)$$

$$\geq \log \|h_{i,0}\|^2 + \sum_{a=1}^{l} \left(\Delta_a - \Delta_a^2\right),$$
(23)

where  $\Delta_a := \frac{\|h_{i,a}\|^2 - \|h_{i,a-1}\|^2}{\|h_{i,a-1}\|^2}$  and the last inequality uses the relation  $\log(1+x) \ge x - x^2$ . We next give a lower bound to  $\Delta_a$ . Let S be the set  $\{k : k \in [m] \text{ and } (h_{i,a-1})_k + (\mathbf{W}_a h_{i,a-1})_k > 0\}$ . We have that

$$\Delta_{a} = \frac{1}{\|h_{i,a-1}\|^{2}} \sum_{k \in S} \left[ (h_{i,a-1})_{k}^{2} + 2\tau(h_{i,a-1})_{k} (\boldsymbol{W}_{a}h_{i,a-1})_{k} + (\tau \boldsymbol{W}_{a}h_{i,a-1})_{k}^{2} \right] - \frac{1}{\|h_{i,a-1}\|^{2}} \sum_{k=1}^{m} (h_{i,a-1})_{k}^{2} 
= -\frac{1}{\|h_{i,a-1}\|^{2}} \sum_{k \notin S} (h_{i,a-1})_{k}^{2} + \frac{1}{\|h_{i,a-1}\|^{2}} \sum_{k \in S} \tau^{2} (\boldsymbol{W}_{a}h_{i,a-1})_{k}^{2} + \frac{2}{\|h_{i,a-1}\|^{2}} \sum_{k \in S} \tau(h_{i,a-1})_{k} (\boldsymbol{W}_{a}h_{i,a-1})_{k} 
\geq -\frac{1}{\|h_{i,a-1}\|^{2}} \sum_{k=1}^{m} (\tau \boldsymbol{W}_{a}h_{i,a-1})^{2} + \frac{2}{\|h_{i,a-1}\|^{2}} \tau \sum_{k=1}^{m} (h_{i,a-1})_{k} (\boldsymbol{W}_{a}h_{i,a-1})_{k} 
= -\frac{\|\tau \boldsymbol{W}_{a}h_{i,a-1}\|^{2}}{\|h_{i,a-1}\|^{2}} + \frac{2\tau \langle h_{i,a-1}, \boldsymbol{W}_{a}h_{i,a-1} \rangle}{\|h_{i,a-1}\|^{2}}, \tag{24}$$

where the inequality is due to the fact that for  $k \notin S |(h_{i,a-1})_k| < |(\tau \boldsymbol{W}_a h_{i,a-1})_k|$  and  $(h_{i,a-1})_k (\boldsymbol{W}_a h_{i,a-1})_k \le 0$ . We let  $\xi_a := \frac{2\tau \langle h_{i,a-1}, \boldsymbol{W}_a h_{i,a-1} \rangle}{\|h_{i,a-1}\|^2}$  and  $\zeta_a := \frac{\|\tau \boldsymbol{W}_a h_{i,a-1}\|^2}{\|h_{i,a-1}\|^2}$ , then  $\Delta_a \ge \xi_a - \zeta_a$ . We note that given  $h_{i,a-1}$ ,  $\xi_a \sim \mathcal{N}\left(0, \frac{8\tau^2}{m}\right)$  and  $\zeta_a \sim \frac{2\tau^2}{m}\chi_m^2$ . Due to equation (19) and (21), we have

$$\mathbb{P}\left(\left|\sum_{l=a}^{b} \xi_{l}\right| \ge \frac{c}{2}\right) \le 2 \exp\left(-\frac{mc^{2}}{128\tau^{2}(b-a+1)}\right)$$

$$\mathbb{P}\left(\sum_{l=a}^{b} \zeta_{l} \ge \frac{c}{2}\right) \le \exp\left(-\Omega\left(\frac{mc^{2}}{\tau^{4}}\right)\right)$$
(25)

Then for any c > 0, and  $\tau \leq \frac{1}{\Omega(\sqrt{L})}$ , we have

$$\mathbb{P}\left(\sum_{a=1}^{l} \Delta_{a} \leq -c\right) = \mathbb{P}\left(\sum_{a=1}^{l} \Delta_{a} \leq -c, \sum_{a=1}^{l} \xi_{a} \geq -\frac{c}{2}\right) + \mathbb{P}\left(\sum_{a=1}^{l} \Delta_{a} \leq -c, \sum_{a=1}^{l} \xi_{a} \leq -\frac{c}{2}\right) \\
\leq \mathbb{P}\left(\sum_{a=1}^{l} \zeta_{a} \geq \frac{c}{2}\right) + \mathbb{P}\left(\sum_{a=1}^{l} \xi_{a} \leq -\frac{c}{2}\right) = e^{-\Omega(c^{2}m)}.$$
(26)

We can derive a similar result that  $\mathbb{P}\left(\sum_{a=1}^{l} \Delta_a \geq c\right) \leq e^{-\Omega(c^2m)}$ . Let a=b in equation (25), we get that for a single  $\Delta_a$ ,

$$\mathbb{P}\left(|\Delta_a| \ge c\right) \le 2e^{-\Omega(Lmc^2)} \tag{27}$$

In addition, we see that

$$\mathbb{P}\left(\sum_{a=1}^{l} \Delta_a^2 \ge c\right) \le \sum_{a=1}^{l} \mathbb{P}\left(\Delta_a^2 \ge \frac{c}{l}\right) = \sum_{a=1}^{l} \mathbb{P}\left(|\Delta_a| \ge \sqrt{\frac{c}{l}}\right) \le 2le^{-\Omega(-mc)}.$$
 (28)

Thus, similar to the equation (26), we get

$$\mathbb{P}\left(\sum_{a=1}^{l} \Delta_a - \Delta_a^2 \le -c\right) \le 2Le^{-\Omega\left(-m\min\{c,c^2\}\right)},\tag{29}$$

which results in

$$\mathbb{P}\left(\log \|h_{i,l}\|^{2} \leq -c\right) \leq \mathbb{P}\left(\log \|h_{i,0}\|^{2} + \sum_{a=1}^{l} \left(\Delta_{a} - \Delta_{a}^{2}\right) \leq -c\right) \leq 2Le^{-\Omega\left(\min\{c,c^{2}\}m\right)}.$$
 (30)

Then we get the conclusion.

#### C.2 Lemmas and Proofs after Perturbation

**Lemma 5.** Suppose that  $\overrightarrow{W}^{(0)}$ , A are randomly generated as in the initialization step, and  $D_0'', \ldots, D_L''$  are diagonal matrices such that  $\|D_l''\|_2 \le 1$  for all  $l \in [L]$  and  $D_l''$  is independent from the randomness  $W_a^{(0)}$  for all a > l, and  $W_1', \ldots, W_L' \in \mathbb{R}^{m \times m}$  are perturbation matrices with  $\|W_l'\|_2 < \tau \omega$  for all  $l \in [L-1]$  and some  $\omega < 1$ . Then with probability at least  $1 - (L) \cdot \exp(-\Omega(m))$  over the initialization randomness we have

$$\|(\mathbf{I} + \tau \mathbf{W}_b^{(0)} + \tau \mathbf{W}_b') \mathbf{D}_{b-1}'' \cdots \mathbf{D}_a'' (\mathbf{I} + \tau \mathbf{W}_a^{(0)} + \tau \mathbf{W}_a')\|_2 \le O(1).$$
(31)

*Proof.* This proof is based on the result of Lemma 1. From Lemma 1, we know for any  $1 \le a \le b < L$ 

$$\|(\mathbf{I} + \tau \mathbf{W}_b^{(0)}) \mathbf{D}_{b-1}'' \cdots \mathbf{D}_a'' (\mathbf{I} + \tau \mathbf{W}_a^{(0)})\|_2 \le 1 + c.$$

Then we have

$$\|(\boldsymbol{I} + \tau \boldsymbol{W}_{b}^{(0)} + \tau \boldsymbol{W}_{b}') \boldsymbol{D}_{b-1}'' \cdots \boldsymbol{D}_{a}'' (\boldsymbol{I} + \tau \boldsymbol{W}_{a}^{(0)} + \tau \boldsymbol{W}_{a}')\|_{2}$$

$$\leq \sum_{j=0}^{b-a+1} {b-a+1 \choose j} (\tau \|\boldsymbol{W}'\|)^{j} (1+c)^{j+1} \leq (1+c) \cdot (1+(1+c)\tau^{2})^{b-a+1} \leq O(1+c),$$

due to the assumption  $\|\boldsymbol{W}_l'\| \le \tau \omega$  for  $l \in [L-1]$  and  $\omega < 1, \tau \le 1/\Omega(\sqrt{L})$ .

**Lemma 6.** Suppose for  $\omega \leq O(1)$ ,  $\tau^2 L \leq 1$ ,  $\|\mathbf{W}'_L\|_2 \leq \omega$  and  $\|\mathbf{W}'_l\|_2 \leq \tau \omega$  for  $l \in [L-1]$ . Then with probability at least  $1 - \exp(-\Omega(m\omega^{2/3}))$ , the following bounds on  $h'_{i,l}$  and  $\mathbf{D}'_{i,l}$  hold for all  $i \in [n]$  and all  $l \in [L-1]$ ,

$$\|h_{i,l}'\| \leq O(\tau^2 L \omega), \quad \|\boldsymbol{D}_{i,l}'\|_0 \leq O\left(m(\omega \tau L)^{2/3}\right), \quad \|h_{i,L}'\| \leq O(\omega), \quad \|\boldsymbol{D}_{i,L}'\|_0 \leq O\left(m\omega^{2/3}\right).$$

*Proof.* Fixing i and ignoring the subscript in i, by Claim 8.2 in Allen-Zhu et al. (2018b), for  $l \in [L-1]$ , there exists  $\mathbf{D}_{l}''$  such that  $|(\mathbf{D}_{l}'')_{k,k}| \leq 1$  and

$$h'_{l} = \mathbf{D}''_{l} \left( (\mathbf{I} + \tau \mathbf{W}_{l} + \tau \mathbf{W}'_{l}) h_{l-1} - (\mathbf{I} + \tau \mathbf{W}_{l}) h_{l-1}^{(0)} \right)$$

$$= \mathbf{D}''_{l} \left( (\mathbf{I} + \tau \mathbf{W}_{l} + \tau \mathbf{W}'_{l}) h'_{l-1} + \tau \mathbf{W}'_{l} h_{l-1}^{(0)} \right)$$

$$= \mathbf{D}''_{l} (\mathbf{I} + \tau \mathbf{W}_{l} + \tau \mathbf{W}'_{l}) \mathbf{D}''_{l-1} (\mathbf{I} + \tau \mathbf{W}_{l-1} + \tau \mathbf{W}'_{l-1}) h'_{l-2}$$

$$+ \tau \mathbf{D}''_{l} (\mathbf{I} + \tau \mathbf{W}_{l} + \tau \mathbf{W}'_{l}) \mathbf{D}''_{l-1} \mathbf{W}'_{l-1} h_{l-2}^{(0)} + \tau \mathbf{D}''_{l} \mathbf{W}'_{l} h_{l-1}^{(0)}$$

$$= \cdots$$

$$= \sum_{a=1}^{l} \tau \mathbf{D}''_{l} (\mathbf{I} + \tau \mathbf{W}_{l} + \tau \mathbf{W}'_{l}) \cdots \mathbf{D}''_{a+1} (\mathbf{I} + \tau \mathbf{W}_{a+1} + \tau \mathbf{W}'_{a+1}) \mathbf{D}''_{a} \mathbf{W}'_{a} h_{a}^{(0)}. \tag{32}$$

We claim that

$$||h_l'|| \le O(\tau^2 L\omega) \tag{33}$$

due to the fact  $\|\boldsymbol{D}_l''\|_2 \leq 1$  and the assumption  $\|\boldsymbol{W}_l'\|_2 \leq \tau \omega$  for  $l \in [L-1]$ . This implies that  $\|h_{i,l}'\|, \|g_{i,l}'\| \leq O(\tau^2 L \omega)$  for all  $l \in [L-1]$  and for all i with probability at least  $1 - O(nL) \cdot \exp(-\Omega(m))$ . One step further, we have  $\|h_L'\|, \|g_L'\| \leq O(\omega)$ .

As for the sparsity  $\|\boldsymbol{D}_l'\|_0$ , we have  $\|\boldsymbol{D}_l'\|_0 \leq O(m(\omega \tau L)^{2/3})$  for every l = [L-1] and  $\|\boldsymbol{D}_L'\|_0 \leq O(m\omega^{2/3})$ .

The argument is as follows (adapt from the Claim 5.3 in Allen-Zhu et al. (2018b)).

We first study the case for  $l \in [L-1]$ . We observe that if  $(\mathbf{D}'_l)_{j,j} \neq 0$  one must have

$$|(g_l')_j| > |(g_l^{(0)})_j|.$$

We note that  $(g_l^{(0)})_j = (h_{l-1}^{(0)} + \tau \boldsymbol{W}_l^{(0)} h_{l-1}^{(0)})_j \sim \mathcal{N}\left((h_{l-1}^{(0)})_j, \frac{2\tau^2 \|h_{l-1}^{(0)}\|^2}{m}\right)$ . Let  $\xi \leq \frac{1}{2\sqrt{m}}$  be a parameter to be chosen later. Let  $S_1 \subseteq [m]$  be a index set satisfying  $S_1 := \{j : |(g_l^{(0)})_j| \leq \xi\tau\}$ . We have  $\mathbb{P}\{|(g_l^{(0)})_j| \leq \xi\tau\} \leq O(\xi\sqrt{m})$  for each  $j \in [m]$ . By Chernoff bound, with probability at least  $1 - \exp(-\Omega(m^{3/2}\xi))$  we have

$$|S_1| \le O(\xi m^{3/2}).$$

Let  $S_2 := \{j : j \notin S_1, \text{ and } (\mathbf{D}'_l)_{j,j} \neq 0\}$ . Then for  $j \in S_2$ , we have  $|(g'_l)_j| > \xi \tau$ . As we have proved that  $||g'_l|| \leq O(\tau \omega)$ , we have

$$|S_2| \le \frac{\|g_l'\|^2}{(\xi \tau)^2} = O((\omega \tau L)^2 / \xi^2).$$

Choosing  $\xi$  to minimize  $|S_1| + |S_2|$ , we have  $\xi = (\omega \tau L)^{2/3} / \sqrt{m}$  and consequently,  $\|\boldsymbol{D}_l'\|_0 \leq O(m(\omega \tau L)^{2/3})$ . Similarly, we have  $\|\boldsymbol{D}_L'\|_0 \leq O(m\omega^{2/3})$ .

We next prove that the norm of a sparse vector after the ResNet mapping.

**Lemma 7.** If  $s \ge \Omega(d/\log m)$ , then for all  $i \in [n]$  and  $a \in [L]$  and for all s-sparse vectors  $u \in \mathbb{R}^m$  and for all  $v \in \mathbb{R}^d$ , the following bound holds with probability at least  $1 - \exp(-\Omega(s \log m))$ 

$$|v^{T} \boldsymbol{B} \boldsymbol{D}_{L} \boldsymbol{W}_{L} \boldsymbol{D}_{L-1} (\boldsymbol{I} + \tau \boldsymbol{W}_{L-1}) \cdots \boldsymbol{D}_{a} (\boldsymbol{I} + \tau \boldsymbol{W}_{a}) u| \leq O\left(\frac{\sqrt{s \log m}}{\sqrt{d}} \|u\| \|v\|\right),$$
(34)

where  $\mathbf{D}_a$  is diagonal matrix with value 0 or 1 and it is independent of  $\mathbf{W}_b$  for any  $L \geq b \geq a \geq 1$  and  $\tau = 1/\Omega(\sqrt{L})$ .

*Proof.* For any fixed vector  $u \in \mathbb{R}^m$ ,  $\|\boldsymbol{D}_{i,L}\boldsymbol{W}_L\boldsymbol{D}_{i,L-1}(\boldsymbol{I} + \tau\boldsymbol{W}_{L-1})\cdots\boldsymbol{D}_{i,a}(\boldsymbol{I} + \tau\boldsymbol{W}_a)u\| \le 1.1\|u\|$  holds with probability at least  $1 - \exp(-\Omega(m))$  (over the randomness of  $\boldsymbol{W}_l$ ,  $l \in [L]$ ).

On the above event, for a fixed vector  $v \in \mathbb{R}^d$  and any fixed  $\mathbf{W}_l$  for  $l \in [L]$ , the randomness only comes from  $\mathbf{B}$ , then  $v^T \mathbf{B} \mathbf{D}_{i,L} \mathbf{W}_L \mathbf{D}_{i,L-1} (\mathbf{I} + \tau \mathbf{W}_{L-1}) \cdots \mathbf{D}_{i,a} (\mathbf{I} + \tau \mathbf{W}_a) u$  is a Gaussian variable with mean 0 and variance no larger than  $O(\|u\| \cdot \|v\| / \sqrt{d})$ . Hence

$$\mathbb{P}\left\{|v^T \boldsymbol{B} \boldsymbol{D}_{i,L} \boldsymbol{W}_L \boldsymbol{D}_{i,L-1} (\boldsymbol{I} + \tau \boldsymbol{W}_{L-1}) \cdots \boldsymbol{D}_{i,a} (\boldsymbol{I} + \tau \boldsymbol{W}_a) u| \ge \sqrt{s \log m} \cdot \Omega(\|u\| \|v\| / \sqrt{d})\right\}$$

$$= \operatorname{erfc}(\Omega(\sqrt{s \log m})) \le \exp(-\Omega(s \log m)).$$

Take  $\epsilon$ -net over all s-sparse vectors of u and all d-dimensional vectors of v, if  $s \geq \Omega(d/\log m)$  then with probability  $1 - \exp(-\Omega(s\log m))$  the claim holds for all s-sparse vectors of u and all d-dimensional vectors of v.

## D Gradient Lower/Upper Bounds and Their Proofs

Because the gradient is pathological and data-dependent, in order to build bound on the gradient, we need to consider all possible point and all cases of data. Hence we first introduce an arbitrary loss vector and then the gradient bound can be obtained by taking a union bound.

We define the  $\mathsf{BP}_{\overrightarrow{W},i}(v,\cdot)$  operator. It back-propagates a vector v to the  $\cdot$  which could be the intermediate output  $h_l$  or the parameter  $W_l$  at the specific layer l using the forward propagation state of input i through the network with parameter  $\overrightarrow{W}$ . Specifically,

$$\begin{split} \mathsf{BP}_{\overrightarrow{\boldsymbol{W}},i}(\boldsymbol{v},\boldsymbol{h}_{l}) &:= (\boldsymbol{I} + \tau \boldsymbol{W}_{l+1})^{T} \boldsymbol{D}_{i,l+1} \cdots (\boldsymbol{I} + \tau \boldsymbol{W}_{L-1})^{T} \boldsymbol{D}_{i,L-1} \boldsymbol{W}_{L}^{T} \boldsymbol{D}_{i,L} \boldsymbol{B}^{T} \boldsymbol{v}, \\ \mathsf{BP}_{\overrightarrow{\boldsymbol{W}},i}(\boldsymbol{v},\boldsymbol{W}_{l}) &:= \tau \left( \boldsymbol{D}_{i,l} (\boldsymbol{I} + \tau \boldsymbol{W}_{l+1})^{T} \cdots (\boldsymbol{I} + \tau \boldsymbol{W}_{L-1})^{T} \boldsymbol{D}_{i,L-1} \boldsymbol{W}_{L}^{T} \boldsymbol{D}_{i,L} \boldsymbol{B}^{T} \boldsymbol{v} \right) \boldsymbol{h}_{i,l-1}^{T} \quad \forall l \in [L-1], \\ \mathsf{BP}_{\overrightarrow{\boldsymbol{W}},i}(\boldsymbol{v},\boldsymbol{W}_{L}) &:= \left( \boldsymbol{D}_{i,L} \boldsymbol{B}^{T} \boldsymbol{v} \right) \boldsymbol{h}_{i,L-1}^{T}. \end{split}$$

Moreover, we introduce

$$\mathsf{BP}_{\overrightarrow{\boldsymbol{W}}}(\overrightarrow{\boldsymbol{v}}, \boldsymbol{W}_l) := \sum_{i=1}^n \mathsf{BP}_{\overrightarrow{\boldsymbol{W}}, i}(v_i, \boldsymbol{W}_l) \quad \forall l \in [L],$$

where  $\overrightarrow{v}$  is composed of n vectors  $v_i$  for  $i \in [n]$ . If  $v_i$  is the error signal of input i, then  $\nabla_{\boldsymbol{W}_l} F_i(\overrightarrow{\boldsymbol{W}}) = \mathsf{BP}_{\overrightarrow{\boldsymbol{W}}_i}(\boldsymbol{B}h_{i,L} - y_i^*, \boldsymbol{W}_l)$ .

#### D.1 Gradient Upper Bound

**Theorem 3.** With probability at least  $1 - \exp(-\Omega(m))$  over the randomness of  $\overrightarrow{W}^{(0)}$ , A, B, it satisfies for every  $l \in [L-1]$ , every  $i \in [n]$ , and every  $\overrightarrow{W}$  with  $\|\overrightarrow{W} - \overrightarrow{W}^{(0)}\|_2 \le \omega$  for  $\omega \in [0,1]$ ,

$$\|\nabla_{\boldsymbol{W}_{l}}F(\overrightarrow{\boldsymbol{W}})\|_{F}^{2} \leq O\left(\frac{F(\overrightarrow{\boldsymbol{W}})}{d} \times \tau^{2}mn\right), \qquad \|\nabla_{\boldsymbol{W}_{L}}F(\overrightarrow{\boldsymbol{W}})\|_{F}^{2} \leq O\left(\frac{F(\overrightarrow{\boldsymbol{W}})}{d} \times mn\right). \tag{5}$$

*Proof.* For each  $i \in [n]$ , we have

$$\left\|\mathsf{BP}_{\overrightarrow{\boldsymbol{W}}}(v_i,\boldsymbol{W}_L)\right\|_F = \left\|\boldsymbol{D}_{i,L}\left(\boldsymbol{B}^Tv_i\right)h_{i,L-1}^T\right\|_F = \left\|\boldsymbol{D}_{i,L}\left(\boldsymbol{B}^Tv_i\right)\right\|\left\|h_{i,L-1}^T\right\| \leq O(\sqrt{m/d})\|v_i\|.$$

Similarly, we have for  $l \in [L-1]$ ,

$$\begin{aligned} \left\| \mathsf{BP}_{\overrightarrow{\boldsymbol{W}}}(v_i, \boldsymbol{W}_l) \right\|_F &= \tau \left( \boldsymbol{D}_{i,l} (\boldsymbol{I} + \tau \boldsymbol{W}_{l+1})^T \cdots (\boldsymbol{I} + \tau \boldsymbol{W}_{L-1})^T \boldsymbol{D}_{i,L-1} \boldsymbol{W}_L^T \boldsymbol{D}_{i,L} \boldsymbol{B}^T v_i \right) h_{i,l-1}^T \\ &\leq O(\tau \sqrt{m/d}) \|v_i\|. \end{aligned}$$

The above upper bounds hold for the initialization  $\overrightarrow{W}^{(0)}$  because of Lemma 1 and Lemma 2. They also hold for all the  $\overrightarrow{W}$  such that  $\|\overrightarrow{W} - \overrightarrow{W}^{(0)}\|_2 \le \omega$  due to Lemma 6 and Lemma 5.

Finally, taking  $\epsilon$ -net over all possible vectors  $\overrightarrow{v} = (v_1, \dots, v_n) \in (\mathbb{R}^d)^n$ , we prove that the above bounds holds for all  $\overrightarrow{v}$ . In particular, we can now plug in the choice of  $v_i = \mathbf{B}h_{i,L} - y_i^*$  and obtain the desired bounds on the true gradients.

#### D.2 Gradient Lower bound

**Theorem 5.** Let  $\omega = O\left(\frac{\delta^{3/2}}{n^{9/2}\log^3 m}\right)$ . With probability at least  $1 - \exp(-\Omega(m\omega^{2/3}))$  over the randomness of  $\overrightarrow{\boldsymbol{W}}^{(0)}, \boldsymbol{A}, \boldsymbol{B}$ , it satisfies for every  $\overrightarrow{\boldsymbol{W}}$  with  $\|\overrightarrow{\boldsymbol{W}} - \overrightarrow{\boldsymbol{W}}^{(0)}\|_2 \leq \omega$ ,

$$\|\nabla_{\boldsymbol{W}_L} F(\overrightarrow{\boldsymbol{W}})\|_F^2 \ge \Omega \left( \frac{\max_{i \in [n]} F_i(\overrightarrow{\boldsymbol{W}})}{dn/\delta} \times m \right). \tag{35}$$

This gradient lower bound on  $\|\nabla_{\boldsymbol{W}_L} F(\overrightarrow{\boldsymbol{W}})\|_F^2$  acts like the gradient dominance condition and it is the same as Allen-Zhu et al. (2018b) except that our range on  $\omega$  does not depend on the depth L.

*Proof.* The gradient lower-bound at the initialization is given in (Allen-Zhu et al., 2018b, Section 6.2) via the smoothed analysis (Spielman and Teng, 2004): with high probability the gradient is lower-bounded, although the worst case it might be 0. The proof is the same given two preconditioned results Lemma 2 and Lemma 8. We shall not repeat the proof here.

Now suppose that we have  $\|\nabla_{\boldsymbol{W}_L} F(\overrightarrow{\boldsymbol{W}}^{(0)})\|_F^2 \geq \Omega\left(\frac{\max_{i \in [n]} F_i(\overrightarrow{\boldsymbol{W}}^{(0)})}{dn/\delta} \times m\right)$ . We next bound the change of the gradient after perturbing the parameter. Recall that

$$\mathsf{BP}_{\overrightarrow{\boldsymbol{W}}^{(0)}}(\overrightarrow{\boldsymbol{v}},\boldsymbol{W}_L) - \mathsf{BP}_{\overrightarrow{\boldsymbol{W}}}(\overrightarrow{\boldsymbol{v}},\boldsymbol{W}_L) = \sum_{i=1}^n \left( (v_i^T \boldsymbol{B} \boldsymbol{D}_{i,L}^{(0)})^T (h_{i,L-1}^{(0)})^T - (v_i^T \boldsymbol{B} \boldsymbol{D}_{i,L})^T (h_{i,L-1})^T \right)$$

By Lemma 6 and Lemma 7, we know,

$$||v_i^T B D_{i,L}^{(0)} - v_i^T B D_{i,L}|| \le O(\sqrt{m\omega^{2/3}}/\sqrt{d}) \cdot ||v_i||.$$

Furthermore, we know

$$||v_i^T \mathbf{B} \mathbf{D}_{i,L}|| \le O(\sqrt{m/d}) \cdot ||v_i||.$$

By Lemma 2 and Lemma 6, we have

$$||h_{i,L-1}^{(0)}|| \le 1.01$$
 and  $||h_{i,L-1} - h_{i,L-1}^{(0)}|| \le O(\omega)$ .

Combing the above bounds together, we have

$$\|\mathsf{BP}_{\overrightarrow{\boldsymbol{W}}^{(0)}}(\overrightarrow{\boldsymbol{v}},\boldsymbol{W}_L) - \mathsf{BP}_{\overrightarrow{\boldsymbol{W}}}(\overrightarrow{\boldsymbol{v}},\boldsymbol{W}_L)\|_F^2 \leq n \|\overrightarrow{\boldsymbol{v}}\|^2 \cdot O(\sqrt{m\omega^{2/3}/d} + \omega\sqrt{m/d})^2 \leq n \|\overrightarrow{\boldsymbol{v}}\|^2 \cdot O\left(\frac{m}{d}\omega^{2/3}\right)$$

Hence the gradient lower bound still holds for  $\overrightarrow{W}$  given  $\omega < O\left(\frac{\delta^{3/2}}{n^{9/2}}\right)$ .

Finally, taking  $\epsilon$ -net over all possible vectors  $\overrightarrow{v} = (v_1, \dots, v_n) \in (\mathbb{R}^d)^n$ , we prove that the above gradient lower bound holds for all  $\overrightarrow{v}$ . In particular, we can now plug in the choice of  $v_i = \mathbf{B}h_{i,L} - y_i^*$  and it implies our desired bounds on the true gradients.

The gradient lower bound requires the following property.

**Lemma 8.** For any  $\delta$  and any pair  $(x_i, x_j)$  satisfying  $||x_i - x_j||_2 \ge \delta$ , then  $||h_{i,l} - h_{j,l}|| \ge \Omega(\delta)$  holds for all  $l \in [L]$  with probability at least  $1 - O(n^2L) \cdot \exp(-\Omega(\log^2 m))$  for  $\tau \le 1/\Omega(\sqrt{L}\log m)$  and  $m \ge \Omega(\tau^2 L^2 \delta^{-2})$ .

The proof of Lemma 8 follows that of (Allen-Zhu et al., 2018b, Appendix C.1) given the condition that  $m \ge \Omega(\tau^2 L^2 \delta^{-2})$ .

# E Semi-smoothness for $\tau \leq 1/\Omega(\sqrt{L})$

With the help of Theorem 3 and several improvements, we can obtain a tighter bound on the semi-smoothness condition of the objective function.

**Theorem 6.** Let  $\omega \in \left[\Omega\left(\left(\frac{d}{m\log m}\right)^{3/2}\right), O(1)\right]$  and  $\overrightarrow{\boldsymbol{W}}^{(0)}, \boldsymbol{A}, \boldsymbol{B}$  be at random initialization and  $\tau^2 L \leq 1$ . With probability at least  $1 - \exp(-\Omega(m\omega^{2/3}))$  over the randomness of  $\overrightarrow{\boldsymbol{W}}^{(0)}, \boldsymbol{A}, \boldsymbol{B}$ , we have for every  $\overrightarrow{\boldsymbol{W}} \in (\mathbb{R}^{m \times m})^L$  with  $\| \widecheck{\boldsymbol{W}}_L - \boldsymbol{W}_L^{(0)} \|_2 \leq \omega$  and  $\| \widecheck{\boldsymbol{W}}_l - \boldsymbol{W}_l^{(0)} \|_2 \leq \tau \omega$  for  $l \in [L-1]$ , and for every  $\overrightarrow{\boldsymbol{W}}' \in (\mathbb{R}^{m \times m})^L$  with  $\| \boldsymbol{W}_L' \|_2 \leq \omega$  and  $\| \boldsymbol{W}_l' \|_2 \leq \tau \omega$  for  $l \in [L-1]$ , we have

$$F(\overrightarrow{\overline{W}} + \overrightarrow{\overline{W}}') \leq F(\overrightarrow{\overline{W}}) + \langle \nabla F(\overrightarrow{\overline{W}}), \overrightarrow{\overline{W}}' \rangle + O(\frac{nm}{d}) \left( \| \mathbf{W}'_L \|_2 + \tau \sum_{l=1}^{L-1} \| \mathbf{W}'_l \|_2 \right)^2 + \sqrt{\frac{mn\omega^{2/3}}{d}} \cdot O\left( \| \mathbf{W}'_L \|_2 + \max\{(\tau L)^{4/3}, 1\} \sum_{l=1}^{L-1} \| \mathbf{W}'_l \|_2 \right) \sqrt{F(\overrightarrow{\overline{W}})}.$$
(36)

Before going to the proof of the theorem, we introduce a lemma.

**Lemma 9.** There exist diagonal matrices  $\mathbf{D}''_{i,l} \in \mathbb{R}^{m \times m}$  with entries in [-1,1] such that  $\forall i \in [n], \forall l \in [L-1]$ ,

$$h_{i,l} - \breve{h}_{i,l} = \sum_{a=1}^{l} (\breve{\boldsymbol{D}}_{i,l} + \boldsymbol{D}_{i,l}'') (\boldsymbol{I} + \tau \breve{\boldsymbol{W}}_{l}) \cdots (\boldsymbol{I} + \tau \breve{\boldsymbol{W}}_{a+1}) (\breve{\boldsymbol{D}}_{i,a} + \boldsymbol{D}_{i,a}'') \tau \boldsymbol{W}_{a}' h_{i,a-1},$$
(37)

and

$$h_{i,L} - \check{h}_{i,L} = (\check{\boldsymbol{D}}_{i,L} + \boldsymbol{D}''_{i,L}) \boldsymbol{W}'_{L} h_{i,L-1} + \sum_{a=1}^{L-1} (\check{\boldsymbol{D}}_{i,L} + \boldsymbol{D}''_{i,L}) \check{\boldsymbol{W}}_{L} \cdots (\boldsymbol{I} + \tau \check{\boldsymbol{W}}_{a+1}) (\check{\boldsymbol{D}}_{i,a} + \boldsymbol{D}''_{i,a}) \tau \boldsymbol{W}'_{a} h_{i,a-1}.$$
(38)

Furthermore, we then have  $\forall l \in [L-1], \|h_{i,l} - \check{h}_{i,l}\| \leq O(\tau^2 L \omega), \|\mathbf{D}''_{i,l}\|_0 \leq O(m(\omega \tau L)^{2/3}), \text{ and } \|h_{i,L} - \check{h}_{i,L}\| \leq O(\omega), \|\mathbf{D}''_{i,L}\|_0 \leq O(m\omega^{2/3}) \text{ and }$ 

$$\| \boldsymbol{B} h_{i,L} - \boldsymbol{B} \check{h}_{i,L} \| \le O(\sqrt{m/d}) \left( \| \boldsymbol{W}_L' \|_2 + \sum_{l=1}^{L-1} \tau \| \boldsymbol{W}_l' \|_2 \right)$$

hold with probability  $1 - \exp(-\Omega(m\omega^{2/3}))$  given  $\|\boldsymbol{W}_L'\|_2 \le \omega, \|\boldsymbol{W}_l'\|_2 \le \tau\omega$  for  $l \in [L-1]$  and  $\omega \le O(1), \tau^2 L \le 1$ .

*Proof.* The proof can adapt from the proof of Claim 8.2 in Allen-Zhu et al. (2018b) and the proof of Lemma 6.  $\Box$ 

Proof of Theorem 6. First of all, we know that  $loss_i := \mathbf{B} \check{h}_{i,L} - y_i^*$ 

$$\frac{1}{2} \|\boldsymbol{B}h_{i,L} - y_i^*\|^2 = \frac{1}{2} \|\ddot{loss}_i + \boldsymbol{B}(h_{i,L} - \breve{h}_{i,L})\|^2 
= \frac{1}{2} \|\ddot{loss}_i\|^2 + \ddot{loss}_i^T \boldsymbol{B}(h_{i,L} - \breve{h}_{i,L}) + \frac{1}{2} \|\boldsymbol{B}(h_{i,L} - \breve{h}_{i,L})\|^2,$$
(39)

and

$$\nabla_{\boldsymbol{W}_{l}} F(\overrightarrow{\boldsymbol{W}}) = \sum_{i=1}^{n} (loss_{i}^{T} \boldsymbol{B} \boldsymbol{D}_{i,L} \boldsymbol{W}_{L} \cdots \boldsymbol{D}_{i,l+1} (\boldsymbol{I} + \tau \boldsymbol{W}_{l}) \boldsymbol{D}_{i,l})^{T} (\tau h_{i,l-1})^{T}.$$
(40)

$$\nabla_{\boldsymbol{W}_{L}} F(\overrightarrow{\boldsymbol{W}}) = \sum_{i=1}^{n} (loss_{i}^{T} \boldsymbol{B} \boldsymbol{D}_{i,L})^{T} (h_{i,l-1})^{T}.$$

$$(41)$$

Then,

$$F(\overrightarrow{W} + \overrightarrow{W}') - F(\overrightarrow{W}) - \langle \nabla F(\overrightarrow{W}), \overrightarrow{W}' \rangle$$

$$= -\langle \nabla F(\overrightarrow{W}), \overrightarrow{W}' \rangle + \frac{1}{2} \sum_{i=1}^{n} \| \mathbf{B} h_{i,L} - y_{i}^{*} \|^{2} - \| \mathbf{B} \check{h}_{i,L} - y_{i}^{*} \|^{2}$$

$$= -\sum_{l=1}^{L} \langle \nabla \mathbf{W}_{l} F(\overrightarrow{W}), \mathbf{W}'_{l} \rangle + \sum_{i=1}^{n} l \check{oss}_{i}^{T} \mathbf{B} (h_{i,L} - \check{h}_{i,L}) + \frac{1}{2} \| \mathbf{B} (h_{i,L} - \check{h}_{i,L}) \|^{2}$$

$$\stackrel{(a)}{=} \frac{1}{2} \sum_{i=1}^{n} \| \mathbf{B} (h_{i,L} - \check{h}_{i,L}) \|^{2} + \sum_{i=1}^{n} l \check{oss}_{i}^{T} \mathbf{B} \left( (\check{\mathbf{D}}_{i,L} + \mathbf{D}''_{i,L}) \mathbf{W}'_{L} h_{i,L-1} - (\check{\mathbf{D}}_{i,L}) \mathbf{W}'_{L} \check{h}_{i,L-1} \right)$$

$$+ \sum_{i=1}^{n} \sum_{l=1}^{L-1} l \check{oss}_{i}^{T} \mathbf{B} \left( (\check{\mathbf{D}}_{i,L} + \mathbf{D}''_{i,L}) \check{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \check{\mathbf{W}}_{l+1}) (\check{\mathbf{D}}_{i,l} + \mathbf{D}''_{i,l}) \tau \mathbf{W}'_{l} h_{i,l-1} \right)$$

$$- \check{\mathbf{D}}_{i,L} \check{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \check{\mathbf{W}}_{l+1}) \check{\mathbf{D}}_{i,l} \mathbf{W}'_{l} (\tau \check{h}_{i,l-1}) \right), \tag{42}$$

where (a) is due to Lemma 9.

We next bound the RHS of (42). We first use Lemma 9 to get

$$\|\boldsymbol{B}(h_{i,L} - \check{h}_{i,L})\| \le O(\sqrt{m/d})(\|\boldsymbol{W}_L'\|_2 + \sum_{l=1}^{L-1} \tau \|\boldsymbol{W}_l'\|_2).$$
 (43)

Next we calculate that for l = L,

$$\left| \overset{\circ}{loss_{i}}^{T} \boldsymbol{B} \left( (\boldsymbol{\check{D}}_{i,L} + \boldsymbol{D}_{i,L}'') \boldsymbol{W}_{L}' h_{i,L-1} - (\boldsymbol{\check{D}}_{i,L}) \boldsymbol{W}_{L}' \check{h}_{i,L-1} \right) \right|$$

$$\leq \left| \overset{\circ}{loss_{i}}^{T} \boldsymbol{B} \left( \boldsymbol{D}_{i,L}'' \boldsymbol{W}_{L}' h_{i,L-1} \right) \right| + \left| \overset{\circ}{loss_{i}}^{T} \boldsymbol{B} \left( \boldsymbol{\check{D}}_{i,L} \boldsymbol{W}_{L}' (h_{i,L-1} - \check{h}_{i,L-1}) \right) \right|.$$

$$(44)$$

For the first term, by Lemma 7 and Lemma 9, we have

$$\left| \overset{\circ}{loss_{i}}^{T} \boldsymbol{B} \left( \boldsymbol{D}_{i,L}^{"} \boldsymbol{W}_{L}^{\prime} h_{i,L-1} \right) \right| \leq O\left( \frac{\sqrt{m\omega^{2/3}}}{\sqrt{d}} \right) \left\| \overset{\circ}{loss_{i}} \right\| \cdot \left\| \boldsymbol{W}_{L}^{\prime} h_{i,L-1} \right\|$$

$$\leq O\left( \frac{\sqrt{m\omega^{2/3}}}{\sqrt{d}} \right) \left\| \overset{\circ}{loss_{i}} \right\| \cdot \left\| \boldsymbol{W}_{L}^{\prime} \right\|_{2},$$

$$(45)$$

where the last inequality is due to  $||h_{i,L-1}|| \leq O(1)$ . For the second term, by Lemma 9 we have

$$\left| \overset{\circ}{loss}_{i}^{T} \boldsymbol{B} \left( \boldsymbol{\breve{D}}_{i,L} \boldsymbol{W}'_{L} (h_{i,L-1} - \boldsymbol{\breve{h}}_{i,L-1}) \right) \right|$$

$$\leq \| \overset{\circ}{loss}_{i} \| \cdot \left\| \boldsymbol{B} \boldsymbol{\breve{D}}_{i,L} \right\|_{2} \cdot \| \boldsymbol{W}'_{L} \|_{2} \| h_{i,L-1} - \boldsymbol{\breve{h}}_{i,L-1} \|$$

$$\leq \| \overset{\circ}{loss}_{i} \| \cdot O \left( \frac{\omega \sqrt{m}}{\sqrt{d}} \right) \cdot \| \boldsymbol{W}'_{L} \|_{2}, \tag{46}$$

where the last inequality is due to the assumption  $\|\|\boldsymbol{W}_L'\|_2\| \leq \omega$ . Similarly for  $l \in [L-1]$ , we ignore the index i for simplicity.

$$\left| \tilde{loss}^{T} \left( \mathbf{B} (\breve{\mathbf{D}}_{L} + \mathbf{D}_{L}'') \breve{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \breve{\mathbf{W}}_{l+1}) (\breve{\mathbf{D}}_{l} + \mathbf{D}_{l}'') - \mathbf{B} \breve{\mathbf{D}}_{L} \breve{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \breve{\mathbf{W}}_{l+1}) \breve{\mathbf{D}}_{l} \right) \mathbf{W}_{l}' (\tau h_{l-1}) \right| \\
= \left| \tilde{loss}^{T} \mathbf{B} \mathbf{D}_{L}'' \breve{\mathbf{W}}_{L} (\mathbf{D}_{L-1} + \mathbf{D}_{L-1}'') (\mathbf{I} + \tau \breve{\mathbf{W}}_{L-1}) \cdots (\mathbf{D}_{l} + \mathbf{D}_{l}'') (\tau \mathbf{W}_{l}' h_{l-1}) \right| \\
+ \sum_{a=l}^{L-1} \tilde{loss}^{T} \mathbf{B} \breve{\mathbf{D}}_{L} \breve{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \breve{\mathbf{W}}_{a+1}) \mathbf{D}_{a}'' (\mathbf{I} + \tau \breve{\mathbf{W}}_{a}) \cdots (\mathbf{D}_{l} + \mathbf{D}_{l}'') (\tau \mathbf{W}_{l}' h_{l-1}) \\
+ \tilde{loss}^{T} \mathbf{B} \breve{\mathbf{D}}_{L} \breve{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \breve{\mathbf{W}}_{l+1}) \breve{\mathbf{D}}_{l} \mathbf{W}_{l}' \tau (h_{l-1} - \breve{h}_{l-1}) \right| \\
+ \sum_{a=l}^{L-1} \left| \tilde{loss}^{T} \mathbf{B} \breve{\mathbf{D}}_{L} \breve{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \breve{\mathbf{W}}_{a+1}) \mathbf{D}_{a}'' (\mathbf{I} + \tau \breve{\mathbf{W}}_{a}) \cdots (\mathbf{D}_{l} + \mathbf{D}_{l}'') (\tau \mathbf{W}_{l}' h_{l-1}) \right| \\
+ \left| \tilde{loss}^{T} \mathbf{B} \breve{\mathbf{D}}_{L} \breve{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \breve{\mathbf{W}}_{a+1}) \mathbf{D}_{a}'' (\mathbf{I} + \tau \breve{\mathbf{W}}_{a}) \cdots (\mathbf{D}_{l} + \mathbf{D}_{l}'') (\tau \mathbf{W}_{l}' h_{l-1}) \right| \\
+ \left| \tilde{loss}^{T} \mathbf{B} \breve{\mathbf{D}}_{L} \breve{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \breve{\mathbf{W}}_{l+1}) \breve{\mathbf{D}}_{l} \mathbf{W}_{l}' \tau (h_{l-1} - \breve{h}_{l-1}) \right|$$
(47)

We next bound the terms in (47) one by one. For the first term, by Lemma 7 and Lemma 9, we have

$$\left| \tilde{loss}^{T} \boldsymbol{B} \boldsymbol{D}_{L}^{"} \boldsymbol{\check{W}}_{L} (\boldsymbol{D}_{L-1} + \boldsymbol{D}_{L-1}^{"}) (\boldsymbol{I} + \tau \boldsymbol{\check{W}}_{L-1}) \cdots (\boldsymbol{D}_{l} + \boldsymbol{D}_{l}^{"}) (\tau \boldsymbol{W}_{l}^{\prime} h_{l-1}) \right|$$

$$\leq O\left(\frac{\sqrt{m\omega^{2/3}}}{\sqrt{d}}\right) \left\| \tilde{loss} \right\| \cdot \left\| \boldsymbol{\check{W}}_{L} (\boldsymbol{D}_{L-1} + \boldsymbol{D}_{L-1}^{"}) (\boldsymbol{I} + \tau \boldsymbol{\check{W}}_{L-1}) \cdots (\boldsymbol{D}_{l} + \boldsymbol{D}_{l}^{"}) (\tau \boldsymbol{W}_{l}^{\prime} h_{l-1}) \right\|$$

$$\stackrel{(a)}{\leq} O\left(\frac{\sqrt{m\omega^{2/3}}}{\sqrt{d}}\right) \cdot \left\| \tilde{loss} \right\| \cdot \tau \|\boldsymbol{W}_{l}^{\prime} \|_{2},$$

$$(48)$$

where (a) is due to the fact  $\| \breve{\boldsymbol{W}}_{L}(\boldsymbol{D}_{L-1} + \boldsymbol{D}_{L-1}'')(\boldsymbol{I} + \tau \breve{\boldsymbol{W}}_{L-1}) \cdots (\boldsymbol{D}_{l} + \boldsymbol{D}_{l}'') \| = O(1)$  and  $\|h_{l-1}\| = O(1)$  holds with high probability.

We have similar bound for every summand in the second term of (47)

$$\left| \overset{\circ}{loss}^{T} B \breve{\boldsymbol{D}}_{L} \breve{\boldsymbol{W}}_{L} \cdots (\boldsymbol{I} + \tau \breve{\boldsymbol{W}}_{a+1}) \boldsymbol{D}_{a}^{"} (\boldsymbol{I} + \tau \breve{\boldsymbol{W}}_{a}) \cdots (\boldsymbol{D}_{l} + \boldsymbol{D}_{l}^{"}) (\tau \boldsymbol{W}_{l}^{'} h_{l-1}) \right|$$

$$\leq O \left( \frac{\sqrt{m(\omega \tau L)^{2/3}}}{\sqrt{d}} \right) \cdot \| \overset{\circ}{loss} \| \cdot \tau \| \boldsymbol{W}_{l}^{'} \|_{2}.$$

$$(49)$$

For the last term in (47), we have

$$\left| \overset{\circ}{loss}^{T} \mathbf{B} \check{\mathbf{D}}_{L} \check{\mathbf{W}}_{L} \cdots (\mathbf{I} + \tau \check{\mathbf{W}}_{l+1}) \check{\mathbf{D}}_{l} \mathbf{W}'_{l} \tau (h_{l-1} - \check{h}_{l-1}) \right|$$

$$\leq \| \overset{\circ}{loss} \| \cdot O\left(\sqrt{m/d}\right) \cdot \| \mathbf{W}'_{l} \|_{2} \cdot \tau \| h_{l-1} - \check{h}_{l-1} \|$$

$$\leq O(\omega \sqrt{m/d}) \cdot \| \overset{\circ}{loss} \| \cdot \tau \| \mathbf{W}'_{l} \|_{2}, \tag{50}$$

where is the last inequality is due to the bound on  $\|h_{i,l-1} - \check{h}_{i,l-1}\|_2$  in Lemma 9. Hence

$$(47) \leq O\left(L\frac{\sqrt{m(\omega\tau L)^{2/3}}}{\sqrt{d}}\right) \cdot \|\mathring{loss}\| \cdot \tau \|\mathbf{W}_{l}'\|_{2}$$

$$\leq O\left((\tau L)^{4/3} \frac{\sqrt{m\omega^{2/3}}}{\sqrt{d}}\right) \cdot \|\mathring{loss}\| \cdot \|\mathbf{W}_{l}'\|_{2}, \tag{51}$$

where the last inequality is because of  $\tau^2 L \leq 1$ .

Having all the above together and using triangle inequality, we have the result.

**Proposition 1** (Proposition 8.3 in in Allen-Zhu et al. (2018b)). Given vectors  $a, b \in \mathbb{R}^m$  and  $\mathbf{D} \in \mathbb{R}^{m \times m}$  the diagonal matrix where  $\mathbf{D}_{k,k} = \mathbf{1}_{a_k \geq 0}$ . Then, there exists a diagonal matrix  $\mathbf{D}'' \in \mathbb{R}^{m \times m}$  with

- $|D_{k,k} + D_{k,k}''| \le 1$  and  $|D_{k,k}''| \le 1$  for every  $k \in [m]$ ,
- $\mathbf{D}_{k,k}'' \neq 0$  only when  $\mathbf{1}_{a_k \geq 0} \neq \mathbf{1}_{b_k \geq 0}$ ,
- $\phi(a) \phi(b) = (D + D'')(a b).$

Proof of Lemma 9. Fixing index i and ignoring the subscript in i for simplicity, by Proposition 1, for each  $l \in [L-1]$  there exists a  $\mathbf{D}_l''$  such that  $|(\mathbf{D}_l'')_{k,k}| \leq 1$  and

$$h_{l} - \check{h}_{l} = \phi((\boldsymbol{I} + \tau \boldsymbol{W}_{l} + \tau \boldsymbol{W}_{l}')h_{l-1}) - \phi((\boldsymbol{I} + \tau \boldsymbol{W}_{l})\check{h}_{l-1})$$

$$= (\check{\boldsymbol{D}}_{l} + \boldsymbol{D}_{l}'') \left( (\boldsymbol{I} + \tau \boldsymbol{W}_{l} + \tau \boldsymbol{W}_{l}')h_{l-1} - (\boldsymbol{I} + \tau \boldsymbol{W}_{l})\check{h}_{l-1} \right)$$

$$= (\check{\boldsymbol{D}}_{l} + \boldsymbol{D}_{l}'')(\boldsymbol{I} + \tau \boldsymbol{W}_{l})(h_{l-1} - \check{h}_{l-1}) + (\check{\boldsymbol{D}}_{l} + \boldsymbol{D}_{l}'')\tau \boldsymbol{W}_{l}'h_{l-1}$$

$$= \sum_{a=1}^{l} (\check{\boldsymbol{D}}_{l} + \boldsymbol{D}_{l}'')(\boldsymbol{I} + \tau \check{\boldsymbol{W}}_{l}) \cdots (\boldsymbol{I} + \tau \check{\boldsymbol{W}}_{a+1})(\check{\boldsymbol{D}}_{a} + \boldsymbol{D}_{a}'')\tau \boldsymbol{W}_{a}'h_{a-1}.$$
(52)

For l = L, we similarly have

$$h_{L} - \check{h}_{L} = (\check{\boldsymbol{D}}_{L} + \boldsymbol{D}_{L}'')\boldsymbol{W}_{L}'h_{L-1} + \sum_{a=1}^{L-1} (\check{\boldsymbol{D}}_{L} + \boldsymbol{D}_{L}'')\check{\boldsymbol{W}}_{L} \cdots (\boldsymbol{I} + \tau \check{\boldsymbol{W}}_{a+1})(\check{\boldsymbol{D}}_{a} + \boldsymbol{D}_{a}'')\tau \boldsymbol{W}_{a}'h_{a-1}.$$
(53)

Then we have following properties.

For  $l \in [L-1]$ ,  $||h_l - \check{h}_l|| \le O(\tau^2 L \omega)$ . This is because  $||(\check{\boldsymbol{D}}_l + \boldsymbol{D}_l'')(\boldsymbol{I} + \tau \check{\boldsymbol{W}}_l) \cdots (\boldsymbol{I} + \tau \check{\boldsymbol{W}}_{a+1})(\check{\boldsymbol{D}}_a + \boldsymbol{D}_a'')|| \le 1.1$  from Lemma 5;  $||h_{a-1}|| \le O(1)$  from Lemma 2; and the assumption  $||\boldsymbol{W}_l'||_2 \le \tau \omega$  for  $l \in [L-1]$ .

For l = L,  $||h_L - \check{h}_L|| \le O\left(||\boldsymbol{W}_L'||_2 + \sum_{l=1}^{L-1} \tau ||\boldsymbol{W}_l'||_2\right) \le O(\omega)$  because of the assumptions  $||\boldsymbol{W}_L'||_2 \le \omega$ ,  $||\boldsymbol{W}_l'||_2 \le \tau \omega$  for  $l \in [L-1]$  and  $\tau^2 L \le 1$ .

For  $l \in [L]$ ,  $\|\boldsymbol{D}_l''\|_0 \leq O(m\omega^{2/3})$ . This is because  $(\boldsymbol{D}_l'')_{k,k}$  is non-zero only at coordinates k where  $(\breve{g}_l)_k$  and  $(g_l)_k$  have opposite signs, where it holds either  $(\boldsymbol{D}_l^{(0)})_{k,k} \neq (\breve{\boldsymbol{D}}_l)_{k,k}$  or  $(\boldsymbol{D}_l^{(0)})_{k,k} \neq (\boldsymbol{D}_l)_{k,k}$ . Therefore by Lemma 6, we have  $\|\boldsymbol{D}_l''\|_0 \leq O(m(\omega \tau L)^{2/3})$  if  $\|\boldsymbol{W}_l'\|_2 \leq \tau \omega$ .

#### F Proof for Main Result Theorem 1

**Theorem 1.** For the ResNet defined and initialized as in Section 2 with  $\tau \leq 1/\Omega(\sqrt{L}\log m)$ , if the network width  $m \geq \max\{L, \Omega(n^{24}\max\{(\tau L)^{14}, 1\}\delta^{-8}d\log^2 m)\}$ , then with probability at least  $1 - \exp(-\Omega(\log^2 m))$ , gradient descent with learning rate  $\eta = \Theta(\frac{d\delta}{n^4m})$  finds a point  $F(\overrightarrow{W}) \leq \varepsilon$  in  $T = \Omega(n^6\delta^{-2}\log\frac{n\log^2 m}{\varepsilon})$  iterations.

#### F.1 Convergence Result for GD

*Proof.* Using Lemma 2 we have  $||h_{i,L}^{(0)}||_2 \le 1.1$  and then using the randomness of  $\boldsymbol{B}$ , it is easy to show that  $||\boldsymbol{B}h_{i,L}^{(0)} - y_i^*||^2 \le O(\log^2 m)$  with probability at least  $1 - \exp(-\Omega(\log^2 m))$ , and therefore

$$F(\overrightarrow{\boldsymbol{W}}^{(0)}) \le O(n\log^2 m). \tag{54}$$

Assume that for every t = 0, 1, ..., T - 1, the following holds,

$$\|\boldsymbol{W}_{L}^{(t)} - \boldsymbol{W}_{L}^{(0)}\|_{F} \le \omega \stackrel{\Delta}{=} O\left(\frac{\delta^{3}}{n^{9}(\tau L)^{7}}\right)$$

$$(55)$$

$$\|\boldsymbol{W}_{l}^{(t)} - \boldsymbol{W}_{l}^{(0)}\|_{F} \le \tau \omega. \tag{56}$$

We shall prove the convergence of GD under the assumption (55) holds, so that previous statements can be applied. At the end, we shall verify that (55) is indeed satisfied.

Letting  $\nabla_t = \nabla F(\overrightarrow{W}^{(t)})$ , we calculate that

$$F(\overrightarrow{\boldsymbol{W}}^{(t+1)}) \leq F(\overrightarrow{\boldsymbol{W}}^{(t)}) - \eta \|\nabla_{t}\|_{F}^{2} + O(\eta^{2}nm/d) \|\nabla_{t}\|_{2}^{2} + \eta \sqrt{F(\overrightarrow{\boldsymbol{W}}^{(t)})} \cdot O\left(\frac{\sqrt{mn\omega^{2/3}}}{\sqrt{d}}\right) \cdot O\left(\|\nabla_{\boldsymbol{W}_{L}^{(t)}}\|_{2} + (\tau L)^{4/3} \sum_{l=1}^{L-1} \|\nabla_{\boldsymbol{W}_{l}^{(t)}}\|_{2}\right)$$

$$\leq F(\overrightarrow{\boldsymbol{W}}^{(t)}) - \eta \|\nabla_{t}\|_{F}^{2} + O\left(\frac{\eta mn(\omega)^{1/3}(\tau L)^{7/3}}{d} + \frac{\eta^{2}n^{2}m^{2}}{d^{2}}\right) \cdot F(\overrightarrow{\boldsymbol{W}}^{(t)})$$

$$\leq \left(1 - \Omega\left(\frac{\eta\delta m}{dn^{2}}\right)\right) F(\overrightarrow{\boldsymbol{W}}^{(t)}),$$

$$(57)$$

where the first inequality uses Theorem 4, the second inequality uses the gradient upper bound in Theorem 3 and the last inequality uses the gradient lower bound in Theorem 5 and the choice of  $\eta$  and the assumption on  $\omega$  (55). That is, after  $T = \Omega(\frac{dn^2}{\eta\delta m})\log\frac{n\log^2 m}{\epsilon}$  iterations  $F(\overrightarrow{\boldsymbol{W}}^{(T)}) \leq \epsilon$ .

We need to verify for each t, (55) holds. By Theorem 3,

$$\|\boldsymbol{W}_{L}^{(t)} - \boldsymbol{W}_{L}^{(0)}\|_{F} \leq \sum_{i=1}^{t-1} \|\eta \nabla_{\boldsymbol{W}_{L}} F(\overrightarrow{\boldsymbol{W}}^{(i)})\|_{F} \leq O(\eta \sqrt{nm/d}) \cdot \sum_{i=1}^{t-1} \sqrt{F(\overrightarrow{\boldsymbol{W}}^{(i)})}$$

$$\stackrel{(a)}{\leq} O(\eta \sqrt{nm/d}) \cdot O\left(\sqrt{n} \log m \cdot \frac{1}{1 - \sqrt{1 - \Omega\left(\frac{\eta \delta m}{dn^{2}}\right)}}\right)$$

$$\stackrel{(b)}{\leq} O\left(\frac{n^{3}\sqrt{d}}{\delta \sqrt{m}} \log m\right), \tag{58}$$

where (a) is due to the relation (57) and (b) is due to the fact that  $1 - \sqrt{1 - \Omega\left(\frac{\eta\delta m}{dn^2}\right)} \ge \frac{1}{2}\Omega\left(\frac{\eta\delta m}{dn^2}\right)$ . Similarly, we have for  $l \in [L-1]$ ,

$$\|\boldsymbol{W}_{l}^{(t)} - \boldsymbol{W}_{l}^{(0)}\|_{F} \leq \sum_{i=1}^{t-1} \|\eta \nabla_{\boldsymbol{W}_{l}} F(\overrightarrow{\boldsymbol{W}}^{(i)})\|_{F} \leq O(\eta \tau \sqrt{nm/d}) \cdot \sum_{i=1}^{t-1} \sqrt{F(\overrightarrow{\boldsymbol{W}}^{(i)})} \leq O\left(\frac{\tau n^{3} \sqrt{d}}{\delta \sqrt{m}} \log m\right).$$

By combining (58) and the assumption on  $\omega$  (55), we obtain a bound on m.

# G Tightness of $\tau = 1/\sqrt{L}$ and the Proof of Theorem 2

**Theorem 2.** For the ResNet defined and initialized as in Section 2, if  $\tau \geq L^{-\frac{1}{2}+c}$ , then in expectation we have

$$\mathbb{E}||h_L||^2 > L^{2c}.\tag{2}$$

*Proof.* By induction we can show for any  $k \in [m]$  and  $l \in [L-1]$ ,

$$(h_l)_k \ge \phi \left( \sum_{a=1}^l \left( \tau \boldsymbol{W}_a h_{a-1} \right)_k \right). \tag{59}$$

It is easy to verify  $(h_1)_k = \phi\left((h_0)_k + (\tau \boldsymbol{W}_1 h_0)_k\right) \ge \phi\left((\tau \boldsymbol{W}_1 h_0)_k\right)$  because of  $(h_0)_k \ge 0$ .

Then assume  $(h_l)_k \ge \phi\left(\sum_{a=1}^l (\tau \boldsymbol{W}_a h_{a-1})_k\right)$ , we show it holds for l+1.

$$(h_{l+1})_k = \phi((h_l)_k + (\tau \mathbf{W}_{l+1}h_l)_k) \ge \phi\left(\phi\left(\sum_{a=1}^l (\tau \mathbf{W}_a h_{a-1})_k\right) + (\tau \mathbf{W}_{l+1}h_l)_k\right) \ge \phi\left(\sum_{a=1}^{l+1} (\tau \mathbf{W}_a h_{a-1})_k\right),$$

where the last inequality can be shown by case study.

Next we can compute the mean and variance of  $\sum_{a=1}^{l} (\tau W_a h_{a-1})_k$  by taking iterative conditioning. We have

$$\mathbb{E}\sum_{a=1}^{l} (\tau \mathbf{W}_{a} h_{a-1})_{k} = 0, \quad \mathbb{E}\left(\sum_{a=1}^{l} (\tau \mathbf{W}_{a} h_{a-1})_{k}\right)^{2} = \frac{2\tau^{2}}{m} \sum_{a=1}^{l} \mathbb{E}\|h_{a-1}\|^{2}.$$
 (60)

Moreover,  $(\tau \mathbf{W}_a h_{a-1})_k$  are jointly Gaussian for all a with mean 0 because  $\mathbf{W}_a$ 's are drawn from independent Gaussian distributions. We use l=2 as an example to illustrate the conclusion, it can be generalized to other l. Assume that  $h_0$  is fixed. First it is easy to verify that  $(\tau \mathbf{W}_1 h_0)_k$  is Gaussian variable with mean 0 and  $(\tau \mathbf{W}_2 h_1)_k | \mathbf{W}_1$  is also Gaussian variable with mean 0. Hence  $[(\tau \mathbf{W}_1 h_0)_k, (\tau \mathbf{W}_2 h_1)_k]$  follows jointly Gaussian with mean vector [0, 0]. Thus  $(\tau \mathbf{W}_1 h_0)_k + (\tau \mathbf{W}_2 h_1)_k$  is Gaussian with mean 0. By induction, we have  $\sum_{a=1}^{l} (\tau \mathbf{W}_a h_{a-1})_k$  is Gaussian with mean 0. Then we have

$$\mathbb{E}\|h_{l}\|^{2} \geq \sum_{k=1}^{m} \mathbb{E}\left(\phi\left(\sum_{a=1}^{l} (\tau \boldsymbol{W}_{a} h_{a-1})_{k}\right)\right)^{2} = \sum_{k=1}^{m} \frac{1}{2} \mathbb{E}\left(\sum_{a=1}^{l} (\tau \boldsymbol{W}_{a} h_{a-1})_{k}\right)^{2}$$

$$= \sum_{k=1}^{m} \frac{\tau^{2} \sum_{a=1}^{l} \mathbb{E}\left[\|h_{a-1}\|^{2}\right]}{m} = \tau^{2} \sum_{a=1}^{l} \mathbb{E}\|h_{a-1}\|^{2},$$
(61)

where the first step is due to (59), the second step is due to the symmetry of Gaussian distribution and the third step is due to equation (60). Since  $(h_l)_k = \phi\left((h_{l-1})_k + (\boldsymbol{W}_l h_{l-1})_k\right)$ , we can show  $\mathbb{E}(h_l)_k^2 \geq (h_{l-1})_k^2$  given  $h_{l-1}$  by numerical integral of Gaussian variable over an interval. Hence we have  $\mathbb{E}\|h_l\|^2 \geq \mathbb{E}\|h_{l-1}\|^2 \geq \cdots \geq \mathbb{E}\|h_0\|^2 = 1$  by iteratively taking conditional expectation. Then combined with (61) and the choice of  $\tau = L^{-\frac{1}{2}+c}$ , we have  $\mathbb{E}\|h_{L-1}\|^2 \geq L^{2c}$ . Because of the random Gaussian initialization of  $\boldsymbol{W}_L$  and  $h_L = \phi(\boldsymbol{W}_L h_{L-1})$ , we have  $\mathbb{E}\|h_L\|^2 = \|h_{L-1}\|^2$ . Thus, the claim is proved.

## H More Empirical Studies

In this section we train the feedforward neural network and ResNet models on the Street View House Numbers(SVHN) dataset(Netzer et al., 2011), and compare their convergence behaviors. The model architectures is the same as Section 5. We run our experiments on both fully-connected and convolutional models. The fully-connected model zoos are generated by varying the depth  $L \in \{30, 100, 500\}$  and the width  $m \in \{128, 1024\}$ . The convolution model zoos are generated by varying the depth  $L \in \{30, 50, 100\}$  and the number of channels  $m \in \{16, 32\}$ . The width of convolution model is the number of convolution kernels of each hidden layer. We choose  $\tau \in \{\frac{1}{L}, \frac{1}{L^{0.5}}\}$  and test its influence<sup>4</sup>.

Data and hyperparameters. SVHN contains more than 70000 training examples with input dimension  $32 \times 32 \times 3$  and 10 labels. The input feature vector are normalized. We use the standard SGD optimizer. Learning rate is chosen as a function of the model width  $lr = \frac{0.1}{(m/16)}$  and minibatch size is 64. There is no pooling layer of our convolution model.

Experiments results. For a given width, we evaluate the training performances of ResNet and feedforward NN with different depths. Figure 2 and 3 show the results of fully connected models and convolutional models, respectively.  $\tau = \frac{1}{L}$  and  $\tau = \frac{1}{L^{0.5}}$  respectively. The results show that deep ResNet with a small  $\tau$  is much easier to train than feedforward NN. However, small  $\tau$  hurts the expressivity of the network, i.e., when the width is large enough (m = 1024 for fully connected models or m = 32 for convolutional models) to train a feedforward NN, ResNet with small  $\tau$  performs worse than feedforward NN.

<sup>&</sup>lt;sup>4</sup> We do not include the case of  $\tau = 1/L^{0.25}$  into comparison because the training fails to converge.

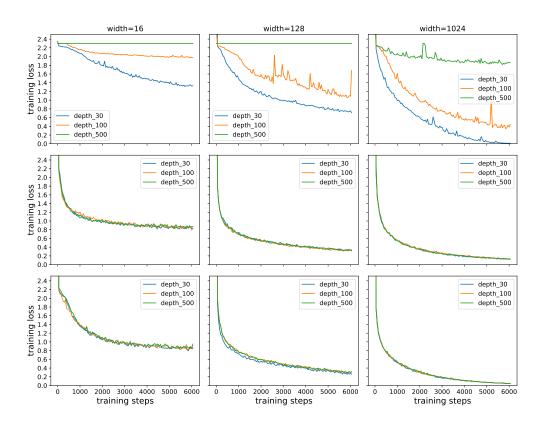


Figure 2: Training performances of fully-connected models on SVHN dataset. The rows from top to bottom are corresponding to feedforward NN, ResNet with  $\tau = \frac{1}{L}$  and  $\frac{1}{L^{0.5}}$ , respectively. The columns from left to right are corresponding to the width at each layer m=16,128 and 1024, respectively. Training loss is evaluated on 10% training random samples.

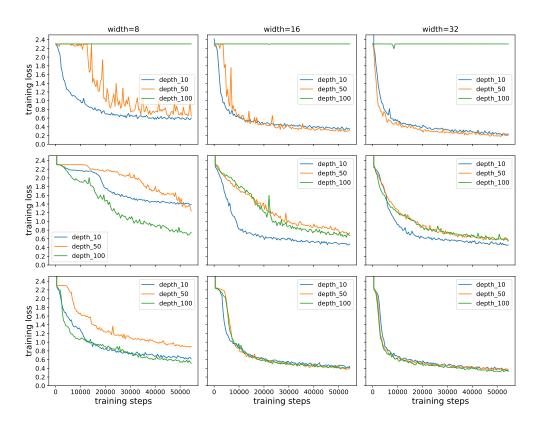


Figure 3: Training performances of convolutional models on SVHN dataset. The rows from top to bottom are corresponding to feedforward NN, ResNet with  $\tau=\frac{1}{L}$  and  $\frac{1}{L^{0.5}}$ , respectively. The columns from left to right are corresponding to the width at each layer m=8,16 and 32, respectively. Training loss is evaluated on 10% random training samples.