

1 We would like to thank the reviewers for their valuable comments. Below, we address the concerns raised.

2 **Reviewer 1.**

3 **Q. The assumption in Line 207 seems strong compared to the results in [Shamir, 2018].**

4 Please note that [S18] considers only one-block ResNets, and as mentioned in Lines 236–239, the condition in Line 207
5 is always satisfied when $L = 1$. In other words, our paper subsumes and extends the results of [S18].

6 **Q. The result in Section 3.2 seems like a trivial application of [Shamir, 2018].**

7 The example in Section 3.2 is not a result of [S18]. Rather, the example is meant to show that direct application of [S18]
8 is **not possible** for deeper ResNets. As discussed in Lines 178–182, an application of [S18] only proves $\text{err}(H_2, Y) \leq$
9 $\text{err}(H_1, Y)$, so a comparison of $\text{err}(H_2, Y)$ and $\text{err}(X, Y)$ does not follow from [S18]. Further, our example shows that
10 even $\text{err}(H_1, Y) > \text{err}(X, Y)$ is possible, showing that theoretically proving $\text{err}(H_2, Y) \leq \text{err}(X, Y)$ is challenging
11 even for $L = 2$. (We hope that the reviewer’s confusion regarding Lines 179–181 is clarified by this answer.)

12 **Q. Section 3 feels separated from the following results in the paper.**

13 The examples make the following points: (1) the advantage of ResNets; and (2) the difficulty of analyzing *deep* ResNets.
14 In particular, the second example highlights that the results in the subsequent sections are not trivial extensions of [S18].
15 We will also reflect the other comments in the revision.

16 **Reviewer 2.**

17 **Q. The paper argues that the geometric conditions [...] sufficient conditions when it should hold.**

18 Thanks for pointing this out. In fact, we do provide a sufficient architectural condition; please refer to Lines 226–235.

19 **Reviewer 3.**

20 **Q. My main concern regarding this paper is that proofs [...] future improvement based on them.**

21 Our proof is not based on a recursive analysis of each residual block. By recursive analysis, we believe that the reviewer
22 means some induction-like argument, which is not the case here. We analyze each block of the network separately using
23 partial derivatives and carefully designed perturbations, and collect the results from each block to prove the theorem.
24 We should mention that our proof technique already reveals the importance of identity skip connection in ResNets: If
25 the identity terms in residual blocks are replaced with tunable parameter matrices W_l , it is impossible to apply the
26 argument in Lines 405–409 unless strong assumptions on W_l are made. This suggests that identity skip connections
27 are indeed the key to the benign loss landscape.

28 **Q. Moreover, the second assumption made in theorem 2 [...] a rather strong assumption on the architecture.**

29 We agree that if the condition in Corollary 3 is not satisfied, the second condition may not necessarily hold. We believe
30 that weakening/removing the condition is a highly nontrivial future challenge. But as noted in Lines 240–246, we’d like
31 to stress that not only does Theorem 2 subsume previous results on shallow ResNets, but it also shows that a **chain**
32 **of multiple** skip connections, *as opposed to direct skip connections to the output*, has beneficial effects on the loss
33 landscape. It has been argued that deep ResNets have benign loss landscapes compared to fully-connected networks,
34 but a rigorous theoretical understanding is still missing; we believe that our results take a step forward in that direction.

35 **Q. The examples provided make understanding and reading the paper easier, [...] standalone results.**

36 We agree; the examples are just to note the advantage of ResNets and highlight the difficulty of analyzing *deep* ResNets.

37 **Q. Additionally, section 5 strongly relies on the near-identity assumption [...] incremental improvements.**

38 As summarized in the paper, previous theoretical results have shown that near-identity regions are expressive and have
39 certain benign optimization properties. In addition, a recent paper “*Fixup Initialization: Residual Learning Without*
40 *Normalization*” showed that initializing ResNets in near-identity regions also leads to good empirical performance. The
41 finding of that paper is rather surprising. For each residual part $\Phi_\theta^l(\cdot)$ (according to our notation), Fixup initializes the
42 last layer of $\Phi_\theta^l(\cdot)$ at zero, and initializes the other layers by using standard random schemes; and then it multiplies a
43 factor inversely proportional to depth L . This means that each $\Phi_\theta^l(\cdot)$ at initialization is zero, hence the network *does*
44 *start in the near-identity region*. Using this initialization scheme, their experiments demonstrate that ResNets can
45 be stably trained **without batch normalization**, and trained networks match the **generalization** performance of the
46 state-of-the-art models. The Fixup paper thus suggests that understanding optimization and generalization of ResNets
47 in near-identity regions is a meaningful and important problem, thus further motivating our analysis in Section 5.

48 **Q. Lastly, in line 119 it is claimed that augmenting one dimension [...] Can you please elaborate on this trick?**

49 We admit that we didn’t provide enough details about this trick. The trick requires some additional structures, e.g., the
50 residual part Φ_θ^l must have an additional output dimension that always outputs 0. Due to space limits, we refer the
51 reviewer to Remark 1 of [Shamir, 2018]; we will add more details as we revise our paper.

52 **Q. Definition of $\text{col}(\cdot)$?**

53 The notation $\text{col}(A)$ means the column space of a matrix A . We will add the definition in the next revision.