

1 We thank all reviewers for their valuable feedback. Below we address each reviewer’s questions and concerns.

2 ——— **Response to Reviewer 1** ———

3 [The theoretical results are “not complete” because only two-layer fully-connected networks are considered.] First of  
4 all, we’d like to remind the reviewer that many interesting deep learning theory papers focus on two-layer networks.  
5 There are *dozens or even hundreds* of such papers published in top conferences and journals. Just for example,  
6 [15,16,17,20,25,34,35,36,43] in our paper’s references are all theoretical studies of two-layer networks. Second, it’s  
7 already challenging to formally establish the two-layer results in our paper; we believe that the conceptual and technical  
8 contributions of the paper are interesting and insightful enough to meet the standard of publication, as acknowledged by  
9 all other reviewers. Third, we didn’t claim we have rigorous results for multi-layer nets, but only provided empirical  
10 results and partial theory to support them, and we explicitly mentioned formally proving them as a future work direction  
11 in Section 5. For these reasons, we find it incorrect and unfair to call our theoretical results “not complete” and to  
12 recommend weak reject mainly based on this. We hope the reviewer can reconsider their decision.

13 [“The experimental results are not sufficient. This work only conducts experiments on single data set, and the  
14 experimental results could not fully verify the theoretical results.”] We do not understand this comment. We provided  
15 experimental results on synthetic data, CIFAR-10, as well as MNIST in the supplementary. These are more than a  
16 “single data set.” Also, our experimental results match theoretical predictions very well, so we do not understand why  
17 they “could not fully verify the theoretical results.” Note that all other reviewers find our experiments convincing.

18 [Related work not sufficient.] We do not understand what’s the reviewer’s specific concern about related work. We have  
19 tried to provide a thorough discussion of related work, and all other reviewers think our discussion is sufficient. If the  
20 reviewer can point out exactly what they think is missing in the discussion, we are happy to incorporate it in the paper.

21 [“The part of the theoretical results is not written well, since the experimental results appear among them.”] Thank  
22 you for the feedback about writing. In fact, we separated the experiments in two sections *on purpose*: Section 3.4 is to  
23 verify the two-layer results in Section 3, and Section 4.2 is for multi-layer and convolutional nets. We thought this is  
24 the clearest way to present our results, and all other reviewers think our paper is well-written (in particular, R2 finds it  
25 “a quite enjoyable read”). Also, we did separate theory and experiments in different subsections so that they are not  
26 really muddled together. Regarding the reviewer’s comment that “theoretical results are not written well”, we do not  
27 understand the reasoning that our arrangement of experiments affects whether the theoretical results are written well.

28 ——— **Response to Reviewer 2** ———

29 [What if  $\alpha > \frac{1}{4}$ ?] If  $\alpha > \frac{1}{4}$ , the conditions are still satisfied with  $\alpha = \frac{1}{4}$ , so our result still applies, which means the  
30 network still behaves like a linear model early in training. It is indeed an interesting question to characterize what  
31 happens after the linear learning period, possibly related to some higher-order kernel as the reviewer points out.

32 ——— **Response to Reviewer 3** ———

33 [Is Assumption 3.1 approximately satisfied by real-world data? Can we transform a given dataset so that Assumption  
34 3.1 holds?] We think it’s possible that Assumption 3.1 is approximately satisfied by certain datasets, but not all. Yes, it’s  
35 possible to transform the data to move closer to this assumption. For example, for CIFAR-10 and MNIST we applied  
36 standard pre-processing to normalize each image to have zero mean and fixed norm (lines 468-469). This already  
37 enables the linear learning behavior to hold. We believe whitening the data can make the assumption better satisfied.

38 [About experiments: (1) How is the linear model calculated? (2) How long does the early phase last if we use a large  
39 learning rate?] (1) For two-layer NN experiments, the linear model is exactly calculated using Eqn. (11). For multi-layer  
40 NNs, since we use erf activation in the experiments, it’s easy to show that the corresponding linear model is  $x \mapsto cx$  for  
41 some constant  $c$  (without bias or the norm-dependent feature). We estimate  $c$  by  $c^2 \approx \frac{\lambda_{\max}(\text{NTK})}{\lambda_{\max}(XX^\top/d)}$ , since we expect  
42  $\text{NTK} \approx c^2 XX^\top/d$ . In general we can use the method sketched in Section 4.1 to compute the linear model analytically.  
43 (2) If the learning rate is 5 times larger, the number of iterations of the agreement will be 5 times smaller. Note that the  
44 correct way to think about this should be progress of learning rather than specific number of iterations. That is, the  
45 shapes of the learning curves will be the same regardless of learning rate (as long as training doesn’t diverge), i.e., the  
46 agreement will last until the linear function finishes learning.

47 ——— **Response to Reviewer 4** ———

48 [What does “we show that these common perceptions can be completely false” mean?] We simply meant that the  
49 network mimics a linear model and doesn’t use its nonlinear capacity early in training, which is a rephrasing of our  
50 main result. We will try to modify the sentence to make it clearer. Thanks for pointing out the confusion.

51 [Any intuition/reason why the networks eventually escape the linear behavior?] The network has the capacity to express  
52 complex nonlinear functions, so it should not be a surprise that it eventually becomes nonlinear. [Extension to deep  
53 networks?] We discussed extension to deep networks in Section 4. [In Section 3.4, what if ReLU was used instead of  
54 ERF?] We indeed provided experiment on ReLU in Figure 2. The corresponding linear model is given by Eqn. (11).