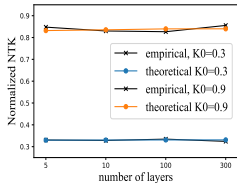
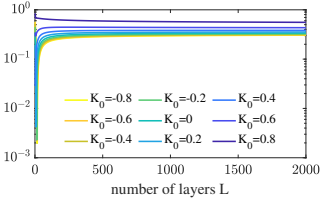


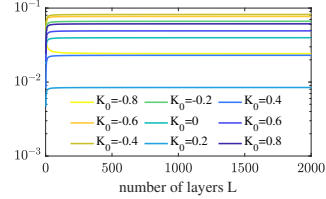
1 **Re: Reviewer #1** For Thm. 4, we randomly initialize the ResNet with width=500, scaling factor  $\gamma = 1$  and depth  
 2  $L = 5, 10, 100, 300$ , and then calculate the inner product of the Jacobians of the ResNet for two different inputs as in the  
 3 definition of NTK (line 185-line187). We repeat the procedure for 500 times and plot the means of the empirical NTKs  
 4 and the theoretical values in Fig. 1a which shows the two results match very well. For Thm. 5 and Thm. 6, Fig. 1b and  
 5 Fig. 1c show that  $\lim_{L \rightarrow \infty} |\bar{\Omega}_L(x, \tilde{x}) - 1/4| \cdot L/\log(L) \approx \text{constant}$  and  $\lim_{L \rightarrow \infty} |\bar{\Omega}_L(x, \tilde{x}) - \bar{\Omega}_1(x, \tilde{x})| \cdot L \approx \text{constant}$   
 6 with  $x^\top \tilde{x} = K_0$  chosen at 9 points.



(a) Thm. 4.  $m = 500$  and scaling  $\gamma = 1$



(b) Thm. 5.  $y$ -axis is  $|\bar{\Omega}_L(x, \tilde{x}) - 1/4| \cdot L/\log(L)$ .



(c) Thm. 6.  $y$ -axis is  $|\bar{\Omega}_L(x, \tilde{x}) - \bar{\Omega}_1(x, \tilde{x})| \cdot L$ .

7 About ResNets' using convolutional layers (Conv) and batch norms(BN) in practice: we acknowledge that Conv and BN  
 8 are indeed important and our techniques can be applied to them with some modification. We will add some discussions  
 9 and leave them for our future investigation.

10 About the scaling factor: The scaling of the bottleneck weight effectively controls the norm of the activations of ResNets,  
 11 and is commonly adopted in theoretical papers, e.g. [14] in our paper. In practice, the norm control is usually achieved  
 12 by normalization techniques (batch norm, group norm, etc.).

13 **Re: Reviewer #2** Thank you for the positive comments. About assumptions: we provide non-asymptotic results  
 14 which bound the error between the finite-width NNs and the infinite-width NNs. Specifically, for ResNets with  $\gamma = 1$ ,  
 15 our theorem ensures good error control, where the width only depends logarithmically on the depth.

16 About generalizability: considering kernel ridge regression, one can show that the generalization error is "continuous"  
 17 w.r.t. the kernel function under some integrable conditions. In this sense, the equivalence (in the view of NTK) means  
 18 that the generalization of NTKs of sufficiently deep ResNets is close to the generalization of the NTK of 1-layer ResNet.  
 19 The same applies to the poor generalization of deep FFNNs.

20 **Re: Reviewer #3.** About our motivation: our increasing depth analysis is a common practice in theoretical research.  
 21 This is because the infinite-depth behavior is very similar to the large-depth behavior of NNs. This is analogous to we  
 22 using central limit theorem in practice, but CLT essentially characterizes the limiting distribution based on infinite many  
 23 samples. Moreover, we highlight that existing results are not limited to 50 and 100 layers. For example, [10] (He et al.  
 24 2016) even show that an ultradeep ResNets of 1001 layers can still generalize very well.

25 We remind that these ultradeep ResNets are not often used in practice because of their massive sizes. If implementable,  
 26 their generalization performance can be better than existing SOTA results. As more advanced computational hardware  
 27 (faster GPU's) is developed, deeper ResNets will become more popular in practice.

28 About our experiments: some existing results on NTK-based methods aim at achieving better performance. Therefore,  
 29 their experiments use more complex structures with additional tricks, e.g. NTKs of ConvNets with Global Average  
 30 Pooling. However, the goal of our experiments is to justify the correctness of our theorems – the NTKs of deep  
 31 FFNNs do not generalize, while the NTKs of deep ResNets generalize well. Therefore, the results of our experiments  
 32 sufficiently serve the purpose of our paper.

33 About degrading kernel: the degeneration of NTKs of deep FFNNs is clearly and rigorously stated in line 237-line 250.

34 About missing references: We will discuss them in the future version, but please kindly notice that Yang's paper is  
 35 made public **after** the submission deadline of NeurIPS 2020.

36 **Re: Reviewer #4.** Thank you for the positive comments and the references. The previous literatures are significant  
 37 but have different aspects from ours. We will add more discussions in the next version.

38 About line 160: We are sorry for the misunderstanding here. Indeed, showing the equivalence of NNs trained by  
 39 gradient flow and NTK kernel predictor is non-trivial. This result is rigorously proved in [22]. Thank you a lot for  
 40 pointing out this issue. We will correct the statement in the next version.