

1 We thank all reviewers for their valuable feedback and appreciating our technical contributions. We first address some  
 2 common concerns.

3 *-More discussion on experimental results:*

4 • First, we discuss more about the relation between GNTKs and GNNs. In Section 5.2, through ablation studies, we  
 5 aimed to verify 1) whether performance of GNNs is similar to GNTKs and 2) whether various techniques used for  
 6 improving performance of GNNs are transferable to GNTKs. For example, in Figure 2, while the observation that  
 7 GNNs with more layers perform better on bioinformatics data is not new, we observed the same trend for GNTKs. In  
 8 Figure 3, we also show that jumping knowledge, which has been shown to be effective in GNNs, is also effective in  
 9 GNTKs. We will elaborate more on these connections in the final version.

10 • Second, Reviewer #2 raised a great point about computational complexity. Indeed, running GNTK is much faster  
 11 than running GNN. On IMDB-B dataset, running GIN with the default setup (official implementation of [26]) takes 19  
 12 minutes on a TITAN X GPU and running GNTK only takes 2 minutes. We will add more details.

13 *-Reproducibility:* We will open-source our code and datasets. We mistakenly chose “not applicable” for this question in  
 14 Reproducibility Response.

15 **To Reviewer #1:**

16 *-Graph convolutional layer:* Our paper does contain a graph (spatial) convolutional layer which we call the aggregation  
 17 step ( $\sum_{v \in \mathcal{N}(u) \cup \{u\}} \mathbf{h}_v^{(\ell-1)}$ ). See line 89, 96 for GNNs and line 148 for GNTKs. We will clarify this more explicitly.

18 *-Lemma 4.2/4.3, conditions on line 193:* The aim of our theoretical analysis is to characterize the function class which  
 19 can be efficiently learned by GNTK. This is a standard type of results in learning theory. The sample complexity  
 20 depends on  $\alpha$ , the coefficient of polynomials, and the norm of  $\beta$ . This kind of dependency is standard in learning  
 21 polynomials and functions that can be expressed as summation of polynomials with fast decaying coefficients. If  $\alpha$   
 22 or the norm of  $\beta$  is large, we do need more samples. See [1] and references therein. The tightness of the bound is an  
 23 open problem. In general, it is impossible to verify these conditions on real datasets, as this is a quantity related to the  
 24 underlying function which we cannot observe. Nevertheless, empirically we use simulations to show the applicability  
 of these theoretical results, as given in Fig. 1. We will include these simulation results in the paper.

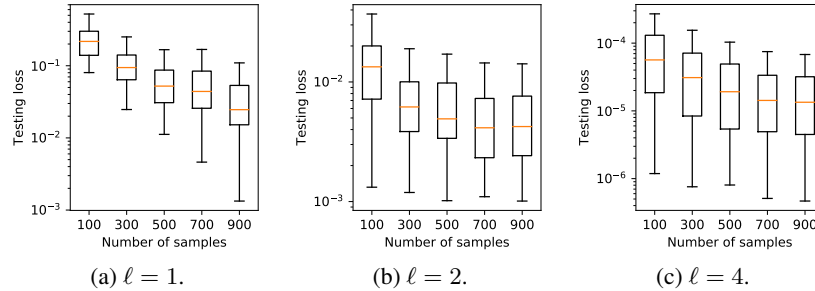


Figure 1: We perform simulations using 1000 graphs in the IMDB-B dataset. For each graph, we generate the label according to  $y = \sum_{u \in V} (\mathbf{h}_u^\top \beta)^\ell$  for  $\ell \in \{1, 2, 4\}$ , where we sample  $\beta$  from the uniform distribution over the unit ball (cf. Eq. (4) in our submission). By construction, these functions satisfy the condition in line 193. We use 100 samples for testing and use 100, 300, 500, 700 or 900 remaining samples for training. We repeat each simulation for 100 times and report the testing loss.

25 *-Theorem 4.1:* This is a classical result for kernel regression. The high-level idea is to use Rademacher complexity to  
 26 bound the difference between empirical loss and population loss, and use  $\mathbf{y}^\top \Theta^{-1} \mathbf{y} \cdot \text{tr}(\Theta)$  to bound the Rademacher  
 27 complexity. The whole proof is given in [3] and we will provide a short outline in the final version.

28 **To Reviewer #2:**

29 *-Statistical test:* We will add a statistical test in the final version. Thanks for the suggestion.

30 *-Line 118 / typos / missing conclusions:* We will fix them accordingly. Thanks for pointing out.

31 **To Reviewer #3:**

32 *-Eq. (2):* For the BLOCK operation with  $R = 2$ , we first apply the aggregation operation to gather local information,  
 33 and then apply a two-layer MLP to create non-linearity. We will clarify this in the final version.

34 *-Experimental details:* For experiments, we use the same setup as in [26]. We have provided some descriptions of the  
 35 experimental setup in Section B. We will add more detailed description in the final version.

36 *-Too many parameters:* GNTK do have some hyper-parameters. However, note that GNTK has strictly smaller number  
 37 of hyper-parameters than GNN since we do not need to tune the learning rate, momentum, weight decay, batch size and  
 38 the width of the MLP layers for GNTK. Furthermore, we found on bioinformatics datasets, we got consistently good  
 39 results by setting the number of BLOCK operations to be 10, the number of MLP layers to be 1 and  $c_u$  to be  $1/|\mathcal{N}(u)|$ .  
 40 We get 75.3% accuracy on PROTEINS, 67.9% on PTC, 83.6% on NCI1. For social network datasets, by setting the  
 41 number of BLOCK operations to be 2, the number of MLP layers to be 2 and  $c_u$  to be 1, we get 76.7% accuracy on  
 42 IMDB-B, 52.8% on IMDB-M, and 83.3% on COLLAB. We will discuss this point in the final version.