
New Insights into Graph Convolutional Networks using Neural Tangent Kernels

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Graph Convolutional Networks (GCNs) have emerged as powerful tools for learn-
2 ing on network structured data. Although empirically successful, GCNs exhibit
3 certain behaviour that has no rigorous explanation—for instance, the performance
4 of GCNs significantly degrades with increasing network depth, whereas it improves
5 marginally with depth using skip connections.

6 This paper focuses on semi-supervised learning on graphs, and explains the above
7 observations through the lens of Neural Tangent Kernels (NTKs). We derive NTKs
8 corresponding to infinitely wide GCNs (with and without skip connections). Sub-
9 sequently, we use the derived NTKs to identify that, with suitable normalisation,
10 network depth does not always drastically reduce the performance of GCNs—a fact
11 that we also validate through extensive simulation. Furthermore, we propose NTK
12 as an efficient ‘surrogate model’ for GCNs that does not suffer from performance
13 fluctuations due to hyper-parameter tuning since it is a hyper-parameter free deter-
14 ministic kernel. The efficacy of this idea is demonstrated through a comparison of
15 different skip connections for GCNs using the surrogate NTKs.

16 1 Introduction

17 Graph structured data are ubiquitous in various domains, including social network analysis, bioin-
18 formatics, communications engineering among others. In recent years, graph neural networks have
19 become an indisputable choice for various learning problems on graphs, and have been employed in
20 a wide range of applications across domains. Several variants of graph neural networks have been
21 proposed, including graph convolutional network [Kipf and Welling, 2017], graph recurrent network
22 [Scarselli et al., 2008, Li et al., 2016], graph attention network [Velickovic et al., 2018], to name a few.
23 The popularity of graph neural networks can be attributed to their ability to tackle two conceptually
24 different learning problems on graphs. In *supervised learning on graphs*, each data instance is a
25 graph and the goal is to predict a label for each graph (for example, a protein structure). In contrast,
26 *semi-supervised learning on graphs* (also called *node classification* or *graph transduction*) refers to
27 the problem of predicting the labels of nodes in a single graph. For instance, given the memberships
28 of a few individuals in a social network, the goal is to predict affiliations of others.

29 This work focuses on the latter problem of semi-supervised learning. GCNs, along with its variants
30 that locally aggregate information in the neighbourhood of each node, have proved to be superior
31 methods in practice [Defferrard et al., 2016, Kipf and Welling, 2017, Chen et al., 2018a, Wu
32 et al., 2019, Chen et al., 2020], outperforming classical, and well-studied, graph embedding based
33 approaches. Among the different variants of GCNs, we focus on the methods based on approximations
34 of spectral graph convolutions [Defferrard et al., 2016, Kipf and Welling, 2017], rather than spatial
35 graph convolutions [Hamilton et al., 2017, Xu et al., 2019]. Surprisingly, these papers suggest shallow
36 networks for the best performance, and unlike the standard neural networks that gain advantage with

37 depth, the performance of GCN has been reported to decrease for deeper nets. This appears to be
 38 due to the over smoothing effect of applying many convolutions, that is, with repeated application
 39 of the graph diffusion operator in each layer, the feature information gets averaged out to a degree
 40 where it becomes uninformative. As a solution to this, Chen et al. [2020] and Kipf and Welling
 41 [2017] proposed different formulations of skip connections in GCNs that overcome the smoothing
 42 effect and thus outperform the vanilla GCN empirically. These networks achieve state-of-the-art
 43 results by directly operating on graphs which enables effective capturing of the complex structural
 44 information as well as the features associated with the entities. However, similar to standard neural
 45 networks, tuning the hyper-parameters is particularly hard due to the highly non-convex objective
 46 function and the over-parameterised setup making it computationally intense. As a result, there is no
 47 theoretical framework that supports rigorous analysis of graph neural networks. Furthermore, the
 48 graph convolutions increase the difficulty of analysis. Motivated by this, we are interested in a more
 49 formal approach to analyze GCNs and, specifically, to understand the influence of depth.

50 Explaining the empirical evidence of deep neural networks through mathematical rigour is an active
 51 area of research. In contrast, theoretical analysis of graph neural networks has been limited in the
 52 literature. From the perspective of learning theory, generalisation error bounds have been derived for
 53 graph neural networks using complexity measures like VC Dimension and Rademacher complexity
 54 [Scarselli et al., 2018, Garg et al., 2020]. However, it is often debated whether generalisation error
 55 bounds can explain the performance of deep neural networks [Neyshabur et al., 2017]. Another line
 56 of research relies on the connection between graph convolutions and belief propagation [Dai et al.,
 57 2016] to analyse the behaviour of graph neural networks in both supervised and semi-supervised
 58 settings using cavity methods and mean field approaches [Zhou et al., 2020b, Kawamoto et al., 2019,
 59 Chen et al., 2018b]. However, the above lines of research do not completely explain the empirical
 60 trends observed in GCNs, especially with regards to the aspects analysed in our work.

61 In this paper, we explain the empirically observed trends of GCNs using the recently introduced
 62 *Neural Tangent Kernel* (NTK) [Jacot et al., 2018]. NTK was proposed to describe the behaviour and
 63 generalisation properties of randomly initialised fully connected neural networks during training by
 64 gradient descent with infinitesimally small learning rate. Jacot et al. [2018] also showed that, as the
 65 network width increases, the change in the kernel during training decreases and hence, asymptotically,
 66 one may replace an infinitely wide neural network by a deterministic kernel machine, where the
 67 kernel (NTK) is defined by the gradient of the network with respect to its parameters as

$$\Theta(x, x') = \mathbb{E}_{W \sim \mathcal{N}} \left[\left\langle \frac{\partial F(W, x)}{\partial W}, \frac{\partial F(W, x')}{\partial W} \right\rangle \right]. \quad (1)$$

68 Here $F(W, x)$ represents the output of the network at data point x and the expectation is with respect
 69 to W , that is, all the parameters of the network randomly sampled from Gaussian distribution \mathcal{N} .
 70 There has been criticism of the ‘infinite width’ assumption being too strong to model real (finite
 71 width) neural networks, and empirical results show that NTK often performs worse than the practical
 72 networks [Arora et al., 2019, Lee et al., 2019]. Nevertheless, theoretical insights on neural network
 73 training gained from NTK have proved to be valuable, particularly in showing how gradient descent
 74 can achieve good generalisation properties [Du et al., 2019a]. Subsequent works have derived
 75 NTK to analyse different neural network architectures in infinite width limit, including convolutional
 76 networks, recurrent networks among others [Arora et al., 2019, Du et al., 2018, 2019a, Alemohammad
 77 et al., 2021]. The most relevant work in the context of our discussion is the work of Du et al. [2019b]
 78 that derived NTK for graph neural networks in the supervised setting (each graph is a data instance to
 79 be classified) and empirically showed that graph NTK outperforms most graph neural networks as
 80 well as other graph kernels for the problem of graph classification.

81 **Focus of this paper and contributions.** The focus of the present paper differs from existing work
 82 on graph NTK [Du et al., 2019b] in two key aspects—we derive NTK for semi-supervised node
 83 classification and, more importantly, we use the derived NTKs to rigorously analyse corresponding
 84 GCN architectures and demonstrate the cause for surprising trends observed empirically in GCNs, as
 85 opposed to standard deep neural networks. More precisely, we make the following contributions:

86 **1.** In Section 2, we derive the NTKs for GCNs used in semi-supervised node classification [Kipf
 87 and Welling, 2017, Wu et al., 2019] in infinite width limit. In contrast to simplifying assumptions in
 88 most NTKs derivations, we allow a non-linear (sigmoid) pooling in the last layer—a natural choice in
 89 practical networks for binary classification. Using the derived NTK and through extensive simulation,

90 we show that the performance of GCN varies considerably for different hyper-parameters, but NTK
 91 captures the general trend of the best possible performance of GCN.

92 **2.** Due to the observation that NTK is a hyper-parameter free alternative to GCN that approximates
 93 the behaviour of GCNs, we suggest NTK as an efficient surrogate for GCN that could be used to
 94 identify the optimal network architecture. We demonstrate this idea in Section 3 by deriving the NTKs
 95 corresponding to GCNs with different skip connections [Chen et al., 2020, Kipf and Welling, 2017],
 96 and we make recommendation on the skip connection for improved performance through empirical
 97 studies of the NTKs. The NTK surrogate can be further used to assess the relative importance of
 98 structure and feature information in a graph dataset.

99 **3.** In Section 4, we use our NTK based analysis to investigate the popular belief that the performance
 100 of vanilla GCN degrades drastically with increasing network depth. We demonstrate that this
 101 observation is due to instabilities in the network training, which results in performance fluctuations of
 102 vanilla GCN, and that can be addressed by appropriate normalisation of the features at each level. The
 103 fluctuations can also be reduced by adding skip connections, even without appropriate normalisation.

104 **4.** In Section 5, we explain an empirical finding—unlike vanilla GCNs, the performance of NTK for
 105 certain skip connections converge with network depth. This is because the NTKs for skip connections
 106 converge with network depth, whereas this is less prominent in the case of NTK for vanilla GCNs.

107 We conclude in Section 6, and provide the NTK derivations and further experimental details in the
 108 appendix.

109 **Notation.** We represent the matrix Hadamard (entry-wise) product by \odot and the scalar product
 110 by $\langle \cdot, \cdot \rangle$. We use $M^{\odot k}$ to denote Hadamard product of matrix M with itself repeated k times. Let
 111 $\mathcal{N}(\mu, \Sigma)$ be Gaussian distribution with mean μ and co-variance Σ . For a function $\sigma(\cdot)$, we use $\dot{\sigma}(\cdot)$
 112 to represent its derivative. We use $\mathbf{1}_{n \times n}$ for the $n \times n$ matrix of ones, I_n for identity matrix of size
 113 $n \times n$, $\mathbb{E}[\cdot]$ for expectation, $\|\cdot\|_F$ denotes Frobenius norm, and $[d] = \{1, 2, \dots, d\}$.

114 2 NTK Captures the Behaviour of Vanilla GCN

115 We consider the problem of node classification in graphs in a semi-supervised setting,¹ where the
 116 labels are observed only for a subset of the nodes. We start with the formal setup and NTK derivation
 117 for the standard (vanilla) GCN proposed in Kipf and Welling [2017].

118 **Formal Setup.** Given a graph with n nodes and a set of node features $\{x_i\}_{i=1}^n \subset \mathbb{R}^f$, we may
 119 assume without loss of generality that the set of observed labels $\{y_i\}_{i=1}^m$ correspond to first m nodes.
 120 We consider a binary classification problem in this paper to simplify the NTK derivation, that is
 121 $y_i \in \{\pm 1\}$, but this could be extended to multi-class problems. The goal is to correctly predict the
 122 $n - m$ unknown labels $\{y_i\}_{i=m+1}^n$. We represent the observed labels of m nodes as $Y \in \{\pm 1\}^{m \times 1}$,
 123 and the node features as $X \in \mathbb{R}^{n \times f}$ with the assumption that entire X is available during training. We
 124 define S to be the graph diffusion operator. The analysis holds for any diffusion S , but for simulations,
 125 we consider the symmetric degree normalized diffusion $S := (D + I_n)^{-\frac{1}{2}}(A + I_n)(D + I_n)^{-\frac{1}{2}}$
 126 where A is the adjacency matrix and D is the degree matrix. We define the GCN of depth d as,

$$F_W(X, S) := \Phi \left(\sqrt{\frac{c_\sigma}{h_d}} S \dots \sigma \left(\sqrt{\frac{c_\sigma}{h_1}} S \sigma (SXW_1) W_2 \right) \dots W_{d+1} \right) \quad (2)$$

127 where $W := \{W_i \in \mathbb{R}^{h_{i-1} \times h_i}\}_{i=1}^{d+1}$ is the set of learnable weight matrices with $h_0 = f$ and $h_{d+1} = 1$,
 128 and $\Phi : \mathbb{R} \rightarrow (-1, +1)$ is re-scaled sigmoid since we consider binary node classification with labels
 129 in $\{\pm 1\}$, h_i is the size of layer $i \in [d]$ and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is the point-wise activation function. We
 130 initialise all the weights to be i.i.d $\mathcal{N}(0, 1)$ and optimise it using stochastic gradient descent. We
 131 study the limiting behavior of this network with respect to the width, that is, $h_1, \dots, h_d \rightarrow \infty$.

132 **Remark 1** (c_σ) *While this setup is similar to Kipf and Welling [2017], it is important to note that*
 133 *we additionally consider the normalisation $\sqrt{c_\sigma/h_i}$ for layer i to ensure that the input norm is*
 134 *approximately preserved. Here, c_σ is a scaling factor to normalize the input in the initialization phase*

¹More precisely, transductive setting as we assume all features are available during training at the same time.

135 and $c_\sigma = \left(\mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[(\sigma(u))^2 \right] \right)^{-1}$ from Du et al. [2019a]. We discuss the role of this normalisation
 136 in Section 4.

137 2.1 NTK for Vanilla GCN

138 We derive the NTK for vanilla GCN by first rewriting $F_W(X, S)$ as defined in (2) using the following
 139 recursive definitions:

$$g_1 := SX, \quad g_i := \sqrt{\frac{c_\sigma}{h_{i-1}}} S \sigma(f_{i-1}) \forall i \in \{2, \dots, d+1\}, \quad f_i := g_i W_i \forall i \in [d+1]$$

$$\text{Output: } F_W(X, S) := \Phi(f_{d+1}), \quad \text{where } \Phi(x) := \frac{2}{1 + \exp(-x)} - 1 \quad (3)$$

140 Using the definitions in (3), the gradient with respect to W_i can be written as

$$\frac{\partial F_W(X, S)}{\partial W_i} := g_i^T b_i \quad \text{with} \quad b_{d+1} := \dot{\Phi}(f_{d+1}), \quad b_i := \sqrt{\frac{c_\sigma}{h_i}} S^T b_{i+1} W_{i+1}^T \odot \dot{\sigma}(f_i) \quad (4)$$

141 We derive the NTK, as defined in (1), using the recursive definition of $F_W(X, S)$ in (3) and its
 142 derivative in (4). The following theorem defines the NTK between every pair of nodes, and the $n \times n$
 143 NTK matrix can be computed at once, as shown below (proof in appendix).

144 **Theorem 1 (NTK for Vanilla GCN)** For the vanilla GCN defined in (2), the NTK Θ is given by

$$\Theta = \left[\sum_{i=1}^{d+1} \Sigma_i \odot (SS^T)^{\odot(d+1-i)} \odot \left(\bigodot_{j=i}^{d+1-i} \dot{E}_j \right) \right] \odot_{f \sim \mathcal{N}(0, \Sigma_d)} \left[\dot{\Phi}(f) \dot{\Phi}(f)^T \right]. \quad (5)$$

145 Here $\Sigma_i \in \mathbb{R}^{n \times n}$ is the co-variance between nodes of the layer f_i , and is given by $\Sigma_1 := SXX^T S^T$,
 146 $\Sigma_i := SE_{i-1} S^T$ with $E_i := c_\sigma \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma_i)} [\sigma(f) \sigma(f)^T]$ and $\dot{E}_i := c_\sigma \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma_i)} [\dot{\sigma}(f) \dot{\sigma}(f)^T]$.

147 Each entry of the expected matrix in (5) can be approximately computed as follows. For $\Delta \in \mathbb{R}^{2 \times 2}$,

$$\mathbb{E}_{(p,q) \sim \mathcal{N}(0, \Delta)} \left[\dot{\Phi}(p) \dot{\Phi}(q) \right] = \frac{1}{4} - \frac{\Delta_{00} + \Delta_{11}}{16} + \frac{\Delta_{00} \Delta_{11} + 2\Delta_{01}^2}{64} + \frac{\Delta_{00}^2 + \Delta_{11}^2}{32} + \frac{\epsilon^3}{16}$$

148 for $|\epsilon| \leq \max\{\Delta_{00}, \Delta_{11}\}$.

149 **Inference using NTK.** The NTK matrix $\Theta \in \mathbb{R}^{n \times n}$ defines the pairwise kernel among all labeled
 150 and unlabeled nodes, where each entry Θ_{pq} represents the kernel between nodes (or features) x_p and
 151 x_q . For inference, consider the sub-matrix $\Theta_l \in \mathbb{R}^{m \times m}$ that consists of the kernel computed between
 152 all pairs of labeled nodes, and $\Theta_u \in \mathbb{R}^{(n-m) \times m}$ that consists of the kernel computed between
 153 all pairs of unlabeled and labeled nodes. In the case of squared loss minimisation by stochastic
 154 gradient descent with infinitesimally small learning rate $\eta \rightarrow 0$, the training dynamics resemble
 155 kernel regression [Arora et al., 2019]. Hence, the labels for unlabeled nodes Y_u can be inferred as

$$Y_u = \Theta_u \Theta_l^{-1} Y \in \mathbb{R}^{n-m} \quad (6)$$

156 which, when thresholded entry-wise at 0, yields the class prediction for unlabeled nodes.

157 The NTK derived in (5) holds for vanilla GCN with arbitrary activation function in (2). Since the
 158 focus of this work is explaining the empirical performance trends of GCNs, we focus on specific
 159 activation functions that fix the network architecture allowing the NTK to be evaluated exactly. We
 160 first consider a linear activation, that results in the SGC network [Wu et al., 2019], and derive the
 161 NTK as follows.

162 **Corollary 1 (Linear GCN)** Consider $\sigma(x) := x$ in $F_W(X, S)$, then $E_i = c_\sigma \Sigma_i$ and $\dot{E}_i = c_\sigma \mathbf{1}_{n \times n}$
 163 in Theorem 1, resulting in the following NTK

$$\Theta = c_\sigma^d \left[\sum_{i=1}^{d+1} \left(S^i X X^T (S^T)^i \right) \odot (SS^T)^{\odot(d+1-i)} \right] \odot_{f \sim \mathcal{N}(0, \Sigma_d)} \left[\dot{\Phi}(f) \dot{\Phi}(f)^T \right].$$

164 where the last expectation is approximated as in Theorem 1. The natural choice of normalisation
 165 constant c_σ is $c_\sigma = 1$ based on Remark 1.

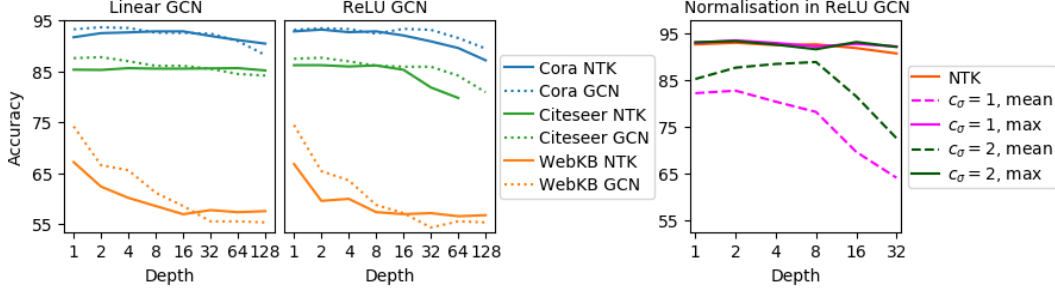


Figure 1: **(left/middle)** Performance of NTK vs GCN in linear and non-linear architectures. The performance trend of NTK matches the best performance of its corresponding GCN in both the architectures. **(right)** Impact of normalisation in ReLU GCN evaluated on Cora dataset. The correct choice of normalisation ($c_\sigma = 2$ in this case) stabilises the training of GCN even in higher depths, and enables identifying the hyper-parameters in lesser time compared to the unnormalised GCN.

166 Considering a non-linear network with ReLU activation, the NTK can be computed as shown below.

167 **Corollary 2 (ReLU GCN)** Consider $\sigma(x) := \text{ReLU}(x)$ in $F_W(X, S)$. The NTK kernel is computed
 168 as in (5), where given Σ_i at each layer, one can evaluate the entries of E_i and \dot{E}_i using a result from
 169 Bietti and Mairal [2019] as

$$\left(E_i\right)_{pq} = \frac{c_\sigma}{2} \sqrt{(\Sigma_i)_{pp} (\Sigma_i)_{qq}} \kappa_1 \left(\frac{(\Sigma_i)_{pq}}{\sqrt{(\Sigma_i)_{pp} (\Sigma_i)_{qq}}} \right) \text{ and } \left(\dot{E}_i\right)_{pq} = \frac{c_\sigma}{2} \kappa_0 \left(\frac{(\Sigma_i)_{pq}}{\sqrt{(\Sigma_i)_{pp} (\Sigma_i)_{qq}}} \right), \quad (7)$$

170 where $\kappa_0(x) := \frac{1}{\pi} (\pi - \arccos(x))$ and $\kappa_1(x) := \frac{1}{\pi} (x (\pi - \arccos(x)) + \sqrt{1 - x^2})$. Based on
 171 Remark 1, the natural choice for normalisation constant c_σ is $c_\sigma = 2$.

172 2.2 Empirical Analysis of Depth

173 Many studies have shown that the performance of vanilla GCN drastically drops with depth due to
 174 the over smoothing effect of convolutional layers [Li et al., 2018, Kipf and Welling, 2017, Chen et al.,
 175 2020]. To validate it, we empirically study the performances of GCN and its NTK counterpart. We
 176 use Tesla K80 GPU with 12GB memory from Google Colab to obtain all our experimental results.

177 **Experimental Setup.** We evaluate the performances of linear and ReLU GCN as stated
 178 in Corollary 1 and 2, respectively, and their corresponding NTKs for different depths $d =$
 179 $\{1, 2, 4, 8, 16, 32, 64, 128\}$. We fix the size of hidden layers h_i to be the same across all layers to re-
 180 duce the number of hyper-parameters. We consider a range of learning rates $\eta = \{10^{-2}, 10^{-3}, 10^{-4}\}$,
 181 different size of the hidden layers $h_i = \{16, 64, 128, 256\}$ and report the best performance among
 182 the different η and size h_i over 10,000 epochs. It is important to note that the chosen learning rates
 183 are in accordance to the theoretical analysis, that is, $\eta \rightarrow 0$. We conduct the experiments with three
 184 datasets, namely *Cora* [McCallum et al., 2000], *Citeseer* [Giles et al., 1998] and *WebKB* [Craven
 185 et al., 1998]. Since the datasets are for multi-class node classification, we combine the classes into
 186 two groups to fit our problem in focus – binary node classification. The choice of class grouping is
 187 decided by comparing the performances of different groupings and ensuring that the two groups are
 188 approximately equal sized. Appendix B includes detailed discussion on the datasets and grouping of
 189 the classes.

190 **NTK captures the performance trend of GCN.** The best performance of GCN decreases with
 191 depth in both linear and non-linear architectures, as observed in other papers. This trend in the best
 192 performance is also confirmed in NTK and thus making it a suitable method to analyse finite width
 193 GCN, despite the fact that the actual performance of the NTK is usually worse than the corresponding
 194 GCN. The left plot of Figure 1² shows the best performance of both the GCN architectures with

²NTK for Citeseer faced out-of-memory issue for depth $d = 128$ in some cases (can also be seen in Figure 2).

195 its NTK counterpart. While there is a drop in the best performance in both the GCNs and the
 196 corresponding NTKs, the drop is not as drastic as it has been reported in other papers. This is due
 197 to two factors: first, unlike the previous works that evaluated the performance for a fixed network
 198 parameterisation, we allow the size of hidden layers to be chosen as a hyper-parameter. We found that
 199 increasing the network size h_i and/or decreasing the learning rate η can reduce the performance drop
 200 with depth. For instance, in *Cora* the best performing network of depth 2 is achieved with $h_i = 16$
 201 and $\eta = 10^{-2}$, whereas, h_i has to be increased to 256 and η has to be reduced to 10^{-4} for depth
 202 128 to achieve similar performance. Second, we identify that the normalisation constant c_σ plays
 203 a crucial role in stabilising the GCN training. The right plot of Figure 1 shows the average and the
 204 best performance of unnormalised ($c_\sigma = 1$) and correctly normalised ($c_\sigma = 2$) ReLU GCN for a
 205 fixed parameterisation with $h_i = 16$ and trained with learning rate $\eta = \{10^{-2}, 10^{-3}, 10^{-4}\}$ over
 206 10,000 epochs. While the average performance of both unnormalised and normalised GCNs shows
 207 a drastic drop, correct normalisation enables the network to learn faster and achieve best results.
 208 Further detailed discussion on the role of normalisation constant c_σ is provided in Section 4.

209 3 NTK - Surrogate for GCN to Analyse Skip Connections

210 Skip connections [Chen et al., 2020, Kipf and Welling, 2017] are one way to overcome the perfor-
 211 mance degradation with depth in GCNs, but little is known about the effectiveness of different forms
 212 of available skip connections. Inspired by the observation of the previous section that the NTK is
 213 a hyper-parameter free model that captures the trends of GCNs, we propose NTK as an efficient
 214 surrogate for GCN, and we investigate different skip connections for GCN in detail in this section. We
 215 consider two formulations of skip connections with two variants each that are described in subsequent
 216 sections. To facilitate skip connections, we need to enforce constant layer size, that is, $h_i = h_{i-1}$.
 217 Therefore, we transform the input layer to H_0 of size $n \times h$ where h is the hidden layer size. This
 218 transformation is necessary as otherwise we would have to assume $h_i = f \forall i \in [d]$ and $h_i \rightarrow \infty$
 219 would not be possible. For this work, we do not consider this transformation as a learnable parameter
 220 in the network. As we consider constant layer size, the NTKs are derived considering $h \rightarrow \infty$. We
 221 first define a skip connection related to the one in Kipf and Welling [2017], where the skip connection
 222 is added to the features before convolution (we refer to it as pre-convolution or Skip-PC).

223 **Definition 1 (Skip-PC)** *In a Skip-PC (pre-convolution) network, the transformed input H_0 is added*
 224 *to the hidden layers before applying the diffusion, leading to the changes in the recursive definition of*
 225 *(3) with $g_1 := SH_0$ and*

$$226 \quad g_i := \sqrt{\frac{c_\sigma}{h}} S(\sigma(f_{i-1}) + \sigma_s(H_0)) \quad \forall i \in \{2, \dots, d+1\}, \quad f_i := g_i W_i \quad \forall i \in [d+1] \quad (8)$$

226 where $\sigma_s(\cdot)$ can be linear or ReLU accounting for two different skip connections.

227 We refer to the network with linear $\sigma_s(\cdot)$ and ReLU $\sigma_s(\cdot)$ as Linear Skip-PC and ReLU Skip-PC,
 228 respectively. The above definition deviates from Kipf and Welling [2017] in the fact that we skip to
 229 the input layer instead of the previous layer. This particular change helps in evaluating the importance
 230 of graph information in a dataset which we discuss in the following section. We also consider a skip
 231 connection similar to the one described in Chen et al. [2020].

232 **Definition 2 (Skip- α)** *Given an interpolation coefficient $\alpha \in (0, 1)$ and a function $\sigma_s(\cdot)$, a Skip- α*
 233 *network is defined such that the transformed input H_0 and the hidden layer are interpolated linearly,*
 234 *which changes the recursive definition in (3) as $g_1 := SH_0$ and*

$$235 \quad g_i := \sqrt{\frac{c_\sigma}{h}} ((1 - \alpha) S(\sigma(f_{i-1})) + \alpha \sigma_s(H_0)) \quad \forall i \in \{2, \dots, d+1\}, \quad f_i := g_i W_i \quad \forall i \in [d+1] \quad (9)$$

235 Similar to Skip-PC, $\sigma_s(\cdot)$ can be linear or ReLU accounting for two different skip connections.
 236 We refer to the network with linear $\sigma_s(\cdot)$ and ReLU $\sigma_s(\cdot)$ as Linear Skip- α and ReLU Skip- α ,
 237 respectively. Chen et al. [2020] recommends the choices for α as 0.1 or 0.2.

238 **Remark 2 (Change of the normalization factor c_σ due to Skip connections)** *Note that the nor-*
 239 *malisation constant c_σ for GCN with skip connections is not the same as defined in Remark 1 of*
 240 *vanilla GCN, since we add the transformed input to the hidden layers. Intuitively, $c_\sigma < 1$ as the norm*
 241 *of the hidden layers would increase otherwise due to the added term. We derived c_σ specifically for*
 242 *non-linear GCN with $\sigma(x) := \text{ReLU}(x)$, and it is $\simeq 0.67$. Refer to Appendix A for the proof.*

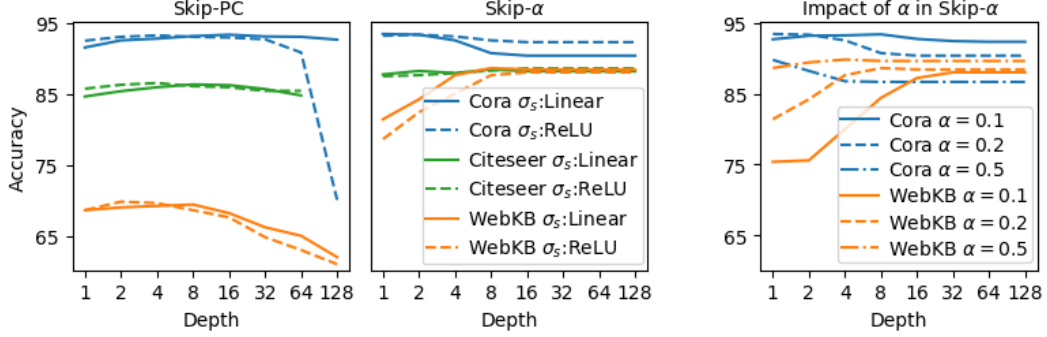


Figure 2: **(left/middle)** Performance of NTKs corresponding to the different skip connections where Skip- α is plotted for $\alpha = 0.2$. **(right)** Impact of α in Skip- α evaluated on Cora and WebKB datasets.

243 3.1 NTK for GCN with Skip Connections

244 We derive NTKs for the skip connections – Skip-PC and Skip- α . Both the NTKs maintain the form
 245 presented in Theorem 1 with the following changes to the co-variance matrices.

246 **Corollary 3 (NTK for Skip-PC)** *The NTK for an infinitely wide Skip-PC network is as presented in*
 247 *Theorem 1 where E_i is defined as in the theorem, but Σ_i is defined as*

$$\Sigma_0 := XX^T, \quad \Sigma_1 := SE_0S^T \quad \text{and} \quad \Sigma_i := SE_{i-1}S^T + \Sigma_1. \quad (10)$$

248 **Corollary 4 (NTK for Skip- α)** *The NTK for an infinitely wide Skip- α network is as presented in*
 249 *Theorem 1 where E_i is defined as in the theorem, but Σ_i is defined with $\Sigma_0 := XX^T$,*

$$\Sigma_1 := (1 - \alpha)^2 SE_0S^T + \alpha(1 - \alpha)(SE_0 + E_0S^T) + \alpha^2 E_0 \quad \text{and} \quad \Sigma_i := SE_{i-1}S^T + E_0. \quad (11)$$

250 Both Corollary 1 and 2 for linear and ReLU activations, respectively, hold for the derived NTKs
 251 corresponding to Skip-PC and Skip- α .

252 3.2 Empirical Analysis

253 Despite studies [Chen et al., 2020, Kipf and Welling, 2017] showing that having skip-connections
 254 gives a significant performance advantage, there is no clear way to choose one formulation of the skip
 255 connection over others. This practical problem can again be seen in the NTK setting as the derived
 256 NTKs have similar structure except the co-variance between the nodes, thus making it difficult to
 257 compare analytically. Therefore, we empirically study the performance of different NTKs in order
 258 to determine the preferred formulation, thereby avoiding computational intensive hyper-parameter
 259 tuning. In addition, we show that the NTK corresponding to Skip- α can be used for assessing the
 260 relevance of structure and feature information of graph in a dataset. We study the non-linear ReLU
 261 GCN with the discussed skip connections, that is, $\sigma(\cdot) := \text{ReLU}$ in (2) empirically.

262 **Experimental setup.** We evaluate the performance of NTKs corresponding to GCNs with skip
 263 connections for different depths $d = \{1, 2, 4, 8, 16, 32, 64, 128\}$ using non-linear activation $\sigma(x) :=$
 264 $\text{ReLU}(x)$ for the GCNs. The linear transformation of the input X is done by $H_0 = XT$ where T
 265 is a $f \times h$ matrix and each entry is sampled from $\mathcal{N}(0, 1)$. The interpolation coefficient α in Skip- α
 266 is chosen to be $\{0.1, 0.2, 0.5\}$. NTKs for all the formulations of skip connections discussed in the
 267 previous section are evaluated on different datasets, namely *Cora*, *Citeseer* and *WebKB*. Figure 2
 268 shows the empirical observations. Refer to Appendix B for more details.

269 We validate the expected performance advantage of GCN with skip connections over vanilla GCN,
 270 more precisely their NTK counterparts, and observe the following main findings.

271 **Non-linear σ_s and shallow net for GCNs with skip connection.** Empirical analysis reveals a
 272 distinct behavior of skip connections with $\sigma_s(\cdot)$ being linear and ReLU, which is illustrated in the
 273 left plot of Figure 2. We observe that the performance of both Skip-PC and Skip- α is not optimal at
 274 deeper depths, and hence we restrict our focus to shallow depths. In the case of shallow depths, we

275 find that using non-linear ReLU $\sigma_s(\cdot)$ in both Skip-PC and Skip- α produces the best performance.
 276 Although ReLU Skip- α initially falls short of its counterpart Linear Skip- α , it eventually outperforms
 277 or performs as good as Linear Skip- α , thus favoring ReLU $\sigma_s(\cdot)$. In addition, this experiment
 278 also validates the general practice of using shallow nets for GCNs. Consequently, we propose skip
 279 connections with ReLU $\sigma_s(\cdot)$ and using shallow nets to achieve the best performance in practice.

280 **NTK as a model to assess relevance of structure and feature information of graphs.** In the left
 281 plot of Figure 2, we notice that the performance of Skip- α on WebKB improved significantly as
 282 compared to Skip-PC and moreover, its performance continued to improve with depth, which is in
 283 contrast to other datasets. We further investigate this by analysing the interpolation coefficient α , and
 284 the corresponding results on Cora and WebKB datasets are shown in the right plot of Figure 2. Large
 285 value of α in Skip- α implies that more importance is given to feature information than the structural
 286 information of the graph. Therefore, from the figure, we infer that the structural information is not
 287 as important as the feature information for WebKB which is in contrast to Cora. Besides, NTK is
 288 a ready-to-use model without the need for hyper-parameter tuning. As a result, we propose NTK
 289 corresponding to Skip- α as a stand-alone model to determine the relative importance of structure and
 290 feature information in tasks where GCNs are employed.

291 4 Role of Normalisation in GCN

292 In Section 2, we discussed that the performance drop with depth in vanilla GCN can be reduced
 293 by varying the size of the hidden layers h_i rather than fixing the network parameterisation as done
 294 in other works. The main difference between our theoretical setup for infinite width GCN and the
 295 practical finite width GCN is the normalisation $\sqrt{c_\sigma/h_i}$. Practical networks generally ignore the
 296 normalisation factor and rely on weight initialisation and optimisation algorithm to stabilise the
 297 training. Intrigued by this, we investigate the role of normalisation applied to each layer by fixing
 298 the network parameterisation in vanilla GCN and Skip-PC empirically. Figure 3 illustrates this for
 299 different $c_\sigma = \{0.67, 1, 2\}$ and depths $d = \{8, 16, 32\}$ on *Cora* dataset. The considered architectures
 300 have non-linear activation, that is, $\sigma(x) := \text{ReLU}(x)$ and we fix the network parameterisation in both
 301 the cases. Different colors in the plot represent the epoch at which the performance is achieved. The
 302 correct choice of c_σ is 2 for ReLU in vanilla GCN (Corollary 2) and 0.67 for Skip-PC (Remark 2).

303 We make the following observation. In the case of vanilla GCN, it is clear that the best performance
 304 is achieved in almost the same number of epochs across all the depths for the correct choice of
 305 $c_\sigma = 2$. Moreover, the decrease in the performance for deeper networks is not significant. Also
 306 we need to train the network longer for $c_\sigma = 1$ to achieve similar performance of the network with
 307 correct normalisation ($c_\sigma = 2$) as we increase the depth. Thus, normalisation plays a crucial role in
 308 stabilising the training of vanilla GCN especially in higher depths. In Skip-PC, the performance of
 309 the network is not significantly affected by c_σ . This is because the residual connection ensures that
 310 the hidden layer norm is approximately equal to the input norm, and thus c_σ is not as relevant as
 311 it is in vanilla GCN case. Therefore, in practice, the absence of this normalisation in vanilla GCN
 312 explains the reported drastic degradation in performance with depth in the existing literature.

313 5 Convergence of NTK with depth

314 In Figure 2, we observe that the performance of NTKs corresponding to GCNs with skip connections
 315 does not change significantly beyond a certain depth. We investigate this behaviour of the NTK further
 316 by measuring the amount of change between NTKs of different depths. To this end, we consider the
 317 alignment between the NTKs in the eigenspace following Fowlkes et al. [2004, Section 4.2]. Formally,
 318 let Θ_i and Θ_j be the NTK of depth i and j , respectively, and $U_i^{(k)}$ and $U_j^{(k)}$ be the matrix of k leading
 319 eigenvectors of Θ_i and Θ_j , respectively, then the alignment between Θ_i and Θ_j is computed by
 320 $a = \frac{1}{k} \left\| U_i^{(k)T} U_j^{(k)} \right\|_F^2$, where $a \in [0, 1]$ with $a = 1$ indicating perfect alignment. Figure 4 shows
 321 the alignment of the NTKs for the discussed non-linear ReLU architectures ($\sigma(\cdot) := \text{ReLU}$ in (2)),
 322 evaluated on *Cora* dataset. Similar pattern is observed in other datasets as well (Appendix B).

323 **The learning happens in shallow depth.** The different alignment plots illustrate the general influ-
 324 ence of depth in GCN. We observe significant changes in the alignment between NTKs of shallow
 325 depths indicating that this is the important part where learning happens. Since the NTKs for both

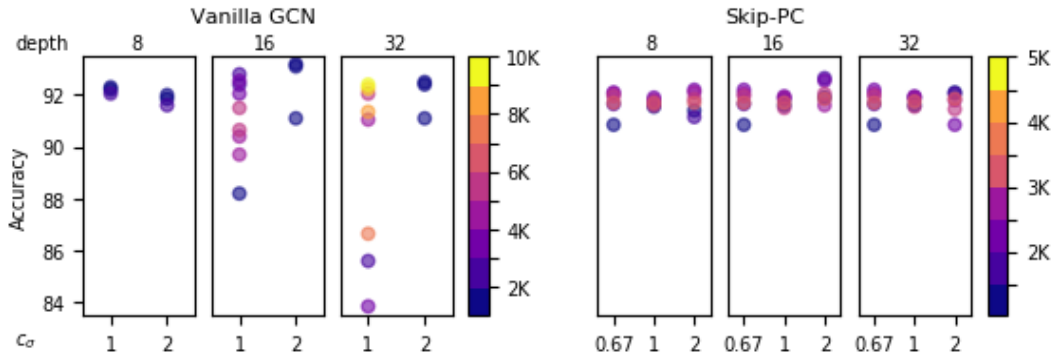


Figure 3: Role of normalisation c_σ in vanilla GCN and Skip-PC as defined in 8. The colorbar represents the number of epochs. The correct choice of c_σ stabilises the training of GCN even in higher depths in vanilla GCN.

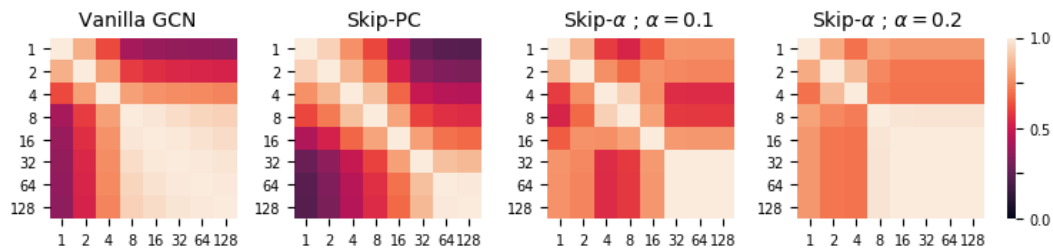


Figure 4: Convergence of NTK with depth for all the discussed ReLU architectures, evaluated on Cora dataset. The plots show perfect alignment of NTKs for higher depths in GCNs with skip connections.

326 vanilla GCN and GCN with skip connections converge with depth, it is clear that deep GCNs have no
 327 advantage or in other words, no new information is learned at deeper depths.

328 **Influence of Skip connection.** In addition, we observe that the NTKs reach almost perfect alignment
 329 with depth for GCNs with skip connection, suggesting that the networks reached saturation in
 330 learning as well. We can further distinguish the presented skip connections: overall Skip-PC has slow
 331 convergence most likely because the skip connection facilitates learning; Skip- α converges fast and
 332 as discussed in Section 3, we observe the influence of α in the learning depending on the dataset.

333 6 Conclusion

334 In this work, we derive NTKs for semi-supervised GCNs, including different formulations of skip
 335 connections. The deterministic hyper-parameter free nature of NTK makes it preferable over its
 336 neural network counterpart since it captures the behaviour of the networks very well, as demonstrated
 337 in our experiments. With the support of our empirical results and the findings from Du et al. [2019b]
 338 that the NTK for supervised GCN outperforms the neural network, we expect the NTKs for semi-
 339 supervised models to perform competitively against the respective neural networks. Nonetheless,
 340 the primary goal of our work is to use NTK to advance our understanding of GCN, particularly
 341 on the impact of depth. In addition, we suggest NTK as a surrogate to study variants of GCNs.
 342 From our surrogate analysis, we propose the NTK corresponding to the skip connection Skip- α as
 343 an efficient ready-to-use off-the-shelf model to determine the relative importance of structure and
 344 feature information in graphs, which we believe to be of great practical value. There is a possibility
 345 of expanding the usage of NTK surrogates to analyse robustness or explainability of GCNs, or other
 346 contexts that involve repeated training of networks. Another direction of research is to incorporate
 347 practical considerations of network architecture in the NTK derivation. The present paper allows
 348 sigmoid functions in the output layer, which is included through a Taylor expansion. It would be also
 349 interesting to derive NTKs considering approximations for softmax, max-pooling, dropout or batch
 350 normalisation, and use the NTKs to analyse the impact of these techniques on network performance.

351 **References**

- 352 Sina Alemohammad, Zichao Wang, Randall Balestriero, and Richard Baraniuk. The recurrent neural
353 tangent kernel. In *International Conference on Learning Representations*, 2021.
- 354 Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On
355 exact computation with an infinitely wide neural net. In *Annual Conference on Neural Information
356 Processing Systems, NeurIPS*, 2019.
- 357 Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. In *NeurIPS
358 2019-Thirty-third Conference on Neural Information Processing Systems*, volume 32, pages 12873–
359 12884, 2019.
- 360 Jie Chen, Tengfei Ma, and Cao Xiao. Fastgcn: Fast learning with graph convolutional networks via
361 importance sampling. In *International Conference on Learning Representations*, 2018a.
- 362 Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph
363 convolutional networks. In *International Conference on Machine Learning*, pages 1725–1735.
364 PMLR, 2020.
- 365 Zhengdao Chen, Lisha Li, and Joan Bruna. Supervised community detection with line graph neural
366 networks. In *International Conference on Learning Representations*, 2018b.
- 367 Mark Craven, Andrew McCallum, Dan PiPasquo, Tom Mitchell, and Dayne Freitag. Learning to
368 extract symbolic knowledge from the world wide web. Technical report, Carnegie-mellon univ
369 pittsburgh pa school of computer Science, 1998.
- 370 Hanjun Dai, Bo Dai, and Le Song. Discriminative embeddings of latent variable models for structured
371 data. In *International conference on machine learning*, pages 2702–2711. PMLR, 2016.
- 372 Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on
373 graphs with fast localized spectral filtering. In *NIPS*, 2016.
- 374 Pedro Domingos. Every model learned by gradient descent is approximately a kernel machine. *arXiv
375 preprint arXiv:2012.00152*, 2020.
- 376 Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global
377 minima of deep neural networks. In *International Conference on Machine Learning*, pages
378 1675–1685. PMLR, 2019a.
- 379 Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
380 over-parameterized neural networks. In *International Conference on Learning Representations*,
381 2018.
- 382 Simon S Du, Kangcheng Hou, Barnabás Póczos, Ruslan Salakhutdinov, Ruosong Wang, and Keyulu
383 Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Conference
384 on Neural Information Processing Systems (NeurIPS)*, 2019b.
- 385 Charless Fowlkes, Serge Belongie, Fan Chung, and Jitendra Malik. Spectral grouping using the
386 nystrom method. *IEEE transactions on pattern analysis and machine intelligence*, 26(2):214–225,
387 2004.
- 388 Vikas Garg, Stefanie Jegelka, and Tommi Jaakkola. Generalization and representational limits of
389 graph neural networks. In *International Conference on Machine Learning*, pages 3419–3430.
390 PMLR, 2020.
- 391 C Lee Giles, Kurt D Bollacker, and Steve Lawrence. Citeseer: An automatic citation indexing system.
392 In *Proceedings of the third ACM conference on Digital libraries*, pages 89–98, 1998.
- 393 William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs.
394 In *Proceedings of the 31st International Conference on Neural Information Processing Systems*,
395 pages 1025–1035, 2017.

- 396 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: convergence and
397 generalization in neural networks. In *Proceedings of the 32nd International Conference on Neural*
398 *Information Processing Systems*, pages 8580–8589, 2018.
- 399 Tatsuro Kawamoto, Masashi Tsubaki, and Tomoyuki Obuchi. Mean-field theory of graph neural
400 networks in graph partitioning. *Journal of Statistical Mechanics: Theory and Experiment*, 2019
401 (12):124007, dec 2019. doi: 10.1088/1742-5468/ab3456.
- 402 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
403 In *International Conference on Learning Representations (ICLR)*, 2017.
- 404 Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha
405 Sohl-Dickstein. Deep neural networks as gaussian processes. In *International Conference on*
406 *Learning Representations*, 2018.
- 407 Jaehoon Lee, Lechao Xiao, Samuel S Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-
408 Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models
409 under gradient descent. In *NeurIPS*, 2019.
- 410 Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks
411 for semi-supervised learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
412 volume 32, 2018.
- 413 Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural
414 networks. In *International Conference on Learning Representations (ICLR)*, 2016.
- 415 Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating the
416 construction of internet portals with machine learning. *Information Retrieval*, 3(2):127–163, 2000.
- 417 Behnam Neyshabur, Srinadh Bhojanapalli, David Mcallester, and Nati Srebro. Exploring generaliza-
418 tion in deep learning. In *Advances in Neural Information Processing Systems*, volume 30. Curran
419 Associates, Inc., 2017.
- 420 Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The
421 graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.
- 422 Franco Scarselli, Ah Chung Tsoi, and Markus Hagenbuchner. The vapnik–chervonenkis dimension
423 of graph and recursive neural networks. *Neural Networks*, 108:248 – 259, 2018.
- 424 Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
425 Bengio. Graph attention networks. *stat*, 1050:4, 2018.
- 426 Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Sim-
427 plifying graph convolutional networks. In *International conference on machine learning*, pages
428 6861–6871. PMLR, 2019.
- 429 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
430 networks? In *International Conference on Learning Representations*, 2019.
- 431 Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang,
432 Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications.
433 *AI Open*, 1:57–81, 2020a.
- 434 Pengfei Zhou, Tianyi Li, and Pan Zhang. Phase transitions and optimal algorithms for semisupervised
435 classifications on graphs: From belief propagation to graph convolution network. *Physical Review*
436 *Research*, 2(3):033325, 2020b.

437 **Checklist**

- 438 1. For all authors...
- 439 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
440 contributions and scope? [Yes]
- 441 (b) Did you describe the limitations of your work? [Yes]
- 442 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Since our
443 work is theoretical and attempts to explain the existing models.
- 444 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
445 them? [Yes]
- 446 2. If you are including theoretical results...
- 447 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 2.
448 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix.
- 449 3. If you ran experiments...
- 450 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
451 mental results (either in the supplemental material or as a URL)? [Yes] Included in the
452 supplementary material.
- 453 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
454 were chosen)? [Yes] Explained in each experiment section and in appendix.
- 455 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
456 ments multiple times)? [No] NTK is deterministic and the seed for GCN is fixed to
457 ensure that the results are reproducible.
- 458 (d) Did you include the total amount of compute and the type of resources used (e.g., type
459 of GPUs, internal cluster, or cloud provider)? [Yes] Included GPU specifications.
- 460 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 461 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 462 (b) Did you mention the license of the assets? [N/A]
- 463 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 464 (d) Did you discuss whether and how consent was obtained from people whose data you're
465 using/curating? [N/A]
- 466 (e) Did you discuss whether the data you are using/curating contains personally identifiable
467 information or offensive content? [N/A]
- 468 5. If you used crowdsourcing or conducted research with human subjects...
- 469 (a) Did you include the full text of instructions given to participants and screenshots, if
470 applicable? [N/A]
- 471 (b) Did you describe any potential participant risks, with links to Institutional Review
472 Board (IRB) approvals, if applicable? [N/A]
- 473 (c) Did you include the estimated hourly wage paid to participants and the total amount
474 spent on participant compensation? [N/A]

475 **A Proofs of NTKs for GCN and GCN with Skip Connections**

476 We provide proofs of Theorem 1 and all corollaries with additional empirical results in this section.

477 **A.1 Proof of NTK for Vanilla GCN (Theorem 1)**

478 **Co-variance between Nodes.** We will first derive the co-variance matrix of size $n \times n$ for each layer
 479 comprising of co-variance between any two nodes p and q . The co-variance between p and q in f_1
 480 and f_i are derived below. We denote p -th row of matrix M as M_p , throughout our proofs.

$$\begin{aligned} \mathbb{E} \left[(f_1)_{pk} (f_1)_{qk'} \right] &= \mathbb{E} \left[(g_1 W_1)_{pk} (g_1 W_1)_{qk'} \right] \\ &= \mathbb{E} \left[\sum_{r=1}^{h_0} (g_1)_{pr} (W_1)_{rk} \sum_{s=1}^{h_0} (g_1)_{qs} (W_1)_{sk'} \right] \stackrel{(W_1)_{xy} \sim \mathcal{N}(0,1)}{=} 0 \quad ; \text{if } r \neq s \text{ or } k \neq k' \\ \mathbb{E} \left[(f_1)_{pk} (f_1)_{qk} \right] &\stackrel{r \equiv s}{k=k'}{\mathbb{E}} \left[\sum_{r=1}^{h_0} (g_1)_{pr} (g_1)_{qr} (W_1)_{rk}^2 \right] \\ &\stackrel{(W_1)_{xy} \sim \mathcal{N}(0,1)}{=} \sum_{r=1}^{h_0} (g_1)_{pr} (g_1)_{qr} = \left\langle (g_1)_p, (g_1)_q \right\rangle \end{aligned} \quad (12)$$

$$\begin{aligned} \mathbb{E} \left[(f_i)_{pk} (f_i)_{qk} \right] &\stackrel{r \equiv s}{k=k'}{\mathbb{E}} \left[\sum_{r=1}^{h_{i-1}} (g_i)_{pr} (g_i)_{qr} (W_i)_{rk}^2 \right] \\ &\stackrel{(W_i)_{xy} \sim \mathcal{N}(0,1)}{=} \sum_{r=1}^{h_{i-1}} (g_i)_{pr} (g_i)_{qr} = \left\langle (g_i)_p, (g_i)_q \right\rangle \end{aligned} \quad (13)$$

$$(12) : \quad \left\langle (g_1)_p, (g_1)_q \right\rangle = \left\langle (SX)_p, (SX)_q \right\rangle = S_p X X^T S_q^T = (\Sigma_1)_{pq} \quad (14)$$

$$\begin{aligned} (13) : \quad \left\langle (g_i)_p, (g_i)_q \right\rangle &= \frac{c_\sigma}{h_{i-1}} \left\langle (S\sigma(f_{i-1}))_p, (S\sigma(f_{i-1}))_q \right\rangle \\ &= \frac{c_\sigma}{h_{i-1}} \sum_{k=1}^{h_{i-1}} (S\sigma(f_{i-1}))_{pk} (S\sigma(f_{i-1}))_{qk} \\ &\stackrel{h_{i-1} \rightarrow \infty}{=} c_\sigma \mathbb{E} \left[(S\sigma(f_{i-1}))_{pk} (S\sigma(f_{i-1}))_{qk} \right] \quad ; \text{law of large numbers} \\ &= c_\sigma \mathbb{E} \left[\left(\sum_{r=1}^n S_{pr} \sigma(f_{i-1})_{rk} \right) \left(\sum_{s=1}^n S_{qs} \sigma(f_{i-1})_{sk} \right) \right] \\ &= c_\sigma \mathbb{E} \left[\sum_{r=1}^n \sum_{s=1}^n S_{pr} S_{qs} \sigma(f_{i-1})_{rk} \sigma(f_{i-1})_{sk} \right] \\ &\stackrel{(a)}{=} \sum_{r=1}^n \sum_{s=1}^n S_{pr} (E_{i-1})_{rs} S_{sq}^T = S_p E_{i-1} S_q^T = (\Sigma_i)_{pq} \end{aligned} \quad (15)$$

481 (a): using $\mathbb{E}[(f_{i-1})_{rk} (f_{i-1})_{sk}] = (\Sigma_{i-1})_{rs}$ and the definition of E_{i-1} in Theorem 1.

482 **NTK for Vanilla GCN.** Let us first evaluate the tangent kernel component from W_i respective to
 483 nodes p and q . The following two results are needed to derive it.

484 **Result 1 (Inner Product of Matrices).** Let a and b be vectors of size $d_1 \times 1$ and $d_2 \times 1$, then

$$\begin{aligned} \langle ab^T, ab^T \rangle &= \text{tr} \left(ab^T (ab^T)^T \right) \\ &= \text{tr} (ab^T ba^T) = \text{tr} (a^T ab^T b) = (a^T a) \odot (b^T b) = \langle a, a \rangle \odot \langle b, b \rangle \end{aligned} \quad (16)$$

485 **Result 2** $\langle (b_r)_p, (b_r)_q \rangle$. We evaluate $\langle (b_r)_p, (b_r)_q \rangle = (b_r b_r^T)_{pq}$ which appears in the gradient.

$$\begin{aligned}
(b_r b_r^T)_{pq} &= \frac{c_\sigma}{h_r} \sum_{k=1}^{h_r} (S^T b_{r+1} W_{r+1}^T)_{pk} \dot{\sigma}(f_r)_{pk} (S^T b_{r+1} W_{r+1}^T)_{qk} \dot{\sigma}(f_r)_{qk} \\
&= \frac{c_\sigma}{h_r} \sum_{k=1}^{h_r} \sum_{i,j}^{n, h_{r+1}} S_{ip} (b_{r+1})_{ij} (W_{r+1})_{kj} \dot{\sigma}(f_r)_{pk} \dot{\sigma}(f_r)_{qk} \sum_{i',j'}^{n, h_{r+1}} S_{i'q} (b_{r+1})_{i'j'} (W_{r+1})_{kj'} \\
&= \frac{c_\sigma}{h_r} \sum_{i,j}^{n, h_{r+1}} \sum_{i',j'}^{n, h_{r+1}} (b_{r+1})_{ij} (b_{r+1})_{i'j'} S_{ip} S_{i'q} \sum_{k=1}^{h_r} (W_{r+1})_{kj} \dot{\sigma}(f_r)_{pk} \dot{\sigma}(f_r)_{qk} (W_{r+1})_{kj'} \\
&= \sum_{j,j'}^{h_{r+1}, h_{r+1}} (S^T b_{r+1})_{pj} (S^T b_{r+1})_{qj'} \frac{c_\sigma}{h_r} \sum_{k=1}^{h_r} (W_{r+1})_{kj} \dot{\sigma}(f_r)_{pk} \dot{\sigma}(f_r)_{qk} (W_{r+1})_{kj'} \\
&\stackrel{h_r \rightarrow \infty}{=} \sum_j^{h_{r+1}} (S^T b_{r+1})_{pj} (S^T b_{r+1})_{qj} c_\sigma \mathbb{E} \left[(W_{r+1}^2)_{kj} \dot{\sigma}(f_r)_{pk} \dot{\sigma}(f_r)_{qk} \right] \quad ; 0 \text{ for } j \neq j' \\
&\stackrel{(b)}{=} \langle (S^T b_{r+1})_p, (S^T b_{r+1})_q \rangle c_\sigma \mathbb{E} [\dot{\sigma}(f_r)_{pk} \dot{\sigma}(f_r)_{qk}] \\
&\stackrel{(16)}{=} (SS^T)_{pq} \langle b_{r+1}, b_{r+1} \rangle_{pq} c_\sigma \mathbb{E} [\dot{\sigma}(f_r)_{pk} \dot{\sigma}(f_r)_{qk}] \\
&= (SS^T)_{pq} \langle b_{r+1}, b_{r+1} \rangle_{pq} \left(\dot{E}_r \right)_{pq} \tag{17}
\end{aligned}$$

486 (b): $(W_{r+1})_{kj}$ is independent and $\mathbb{E} \left[(W_{r+1}^2)_{kj} \right] = 1$.

487 Now, lets derive $\left\langle \left(\frac{\partial F}{\partial W_i} \right)_p, \left(\frac{\partial F}{\partial W_i} \right)_q \right\rangle$ and $\left\langle \left(\frac{\partial F}{\partial W_1} \right)_p, \left(\frac{\partial F}{\partial W_1} \right)_q \right\rangle$ using the above results.

$$\begin{aligned}
\left\langle \left(\frac{\partial F}{\partial W_i} \right)_p, \left(\frac{\partial F}{\partial W_i} \right)_q \right\rangle &= \langle (g_i)_p^T, (b_i)_p, (g_i)_q^T, (b_i)_q \rangle \\
&\stackrel{(16)}{=} \langle (g_i)_p, (g_i)_q \rangle \odot \langle (b_i)_p, (b_i)_q \rangle \\
&\stackrel{(15),(17)}{=} (\Sigma_i)_{pq} (SS^T)_{pq} \langle b_{r+1}, b_{r+1} \rangle_{pq} \left(\dot{E}_r \right)_{pq} \\
&\stackrel{(c)}{=} (\Sigma_i)_{pq} \left((SS^T)_{pq} \right)^{d+1-i} \left(\prod_{j=i}^{d+1-i} \left(\dot{E}_j \right)_{pq} \right) \langle b_{d+1}, b_{d+1} \rangle_{pq} \\
&\stackrel{(d)}{=} (\Sigma_i)_{pq} \left((SS^T)_{pq} \right)^{d+1-i} \left(\prod_{j=i}^{d+1-i} \left(\dot{E}_j \right)_{pq} \right) \left(\dot{\Phi}(f_{d+1}) \dot{\Phi}(f_{d+1})^T \right)_{pq} \tag{18}
\end{aligned}$$

488 (c): repeated application of (17).

489 (b): definition of b_{d+1} .

490 Extending (18) to all n nodes which will result in $n \times n$ matrix,

$$\begin{aligned}
\left\langle \frac{\partial F}{\partial W_i}, \frac{\partial F}{\partial W_i} \right\rangle &= \Sigma_i \odot (SS^T)^{\odot d+1-i} \bigcirc_{j=i}^{d+1-i} \dot{E}_j \odot \dot{\Phi}(f_{d+1}) \dot{\Phi}(f_{d+1})^T \\
\mathbb{E}_{W_i} \left[\left\langle \frac{\partial F}{\partial W_i}, \frac{\partial F}{\partial W_i} \right\rangle \right] &= \Sigma_i \odot (SS^T)^{\odot d+1-i} \bigcirc_{j=i}^{d+1-i} \dot{E}_j \odot \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma_d)} \left[\dot{\Phi}(f) \dot{\Phi}(f)^T \right] \tag{19}
\end{aligned}$$

491 Finally, NTK Θ is,

$$\begin{aligned}\Theta &= \sum_{i=1}^{d+1} \mathbb{E}_{W_i} \left[\left\langle \frac{\partial F}{\partial W_i}, \frac{\partial F}{\partial W_i} \right\rangle \right] \\ &= \left[\sum_{i=1}^{d+1} \Sigma_i \odot (SS^T)^{\odot(d+1-i)} \odot \left(\bigcirc_{j=i}^{d+1-i} \dot{E}_j \right) \right] \odot_{f \sim \mathcal{N}(0, \Sigma_d)} \left[\dot{\Phi}(f) \dot{\Phi}(f)^T \right] \quad (20)\end{aligned}$$

492 We will now compute $\mathbb{E}_{f \sim \mathcal{N}(0, \Sigma_d)} \left[\dot{\Phi}(f) \dot{\Phi}(f)^T \right]$. We use Lagrange form of the remainder to
493 approximate the Taylor's expansion for the re-scaled sigmoid function $\Phi(\cdot)$ which gives better bound.

$$\begin{aligned}\Phi(x) &= \frac{2}{1 + \exp^{-x}} - 1 = \frac{x}{2} - \frac{x^3}{24} + \frac{x^5}{240} + \dots \\ \dot{\Phi}(x) &= \frac{1}{2} - \frac{x^2}{8} + \frac{x^4}{48} + \frac{x^6 \dot{\Phi}^6(\xi)}{6!} \quad ; \text{ last term is the Lagrange form of the remainder.} \quad (21)\end{aligned}$$

494 To evaluate the expectation of an entry i, j in the matrix $\dot{\Phi}(f) \dot{\Phi}(f)^T$, let us define Δ as a 2×2
495 co-variance matrix as follows, $\Delta = \begin{bmatrix} (\Sigma_{d+1})_{ii} & (\Sigma_{d+1})_{ij} \\ (\Sigma_{d+1})_{ji} & (\Sigma_{d+1})_{jj} \end{bmatrix}$

$$\begin{aligned}\mathbb{E}_{(x,y) \sim \Delta} \left[\dot{\Phi}(x) \dot{\Phi}(y) \right] &\stackrel{(21)}{=} \mathbb{E}_{(x,y) \sim \Delta} \left[\left(\frac{1}{2} - \frac{x^2}{8} + \frac{x^4}{48} + \frac{x^6 \dot{\Phi}^6(\xi)}{6!} \right) \left(\frac{1}{2} - \frac{y^2}{8} + \frac{y^4}{48} + \frac{y^6 \dot{\Phi}^6(\xi)}{6!} \right) \right] \\ &= \frac{1}{4} \mathbb{E}_{(x,y) \sim \Delta} \left[1 - \frac{x^2}{4} - \frac{y^2}{4} + \frac{x^4}{24} + \frac{y^4}{24} + \frac{x^2 y^2}{16} - \frac{x^4 y^2}{96} - \frac{x^2 y^4}{96} \right. \\ &\quad \left. + \frac{x^4 y^4}{576} + \frac{x^6 \dot{\Phi}^6(\xi)}{6!} \left(\frac{1}{2} - \frac{y^2}{8} + \frac{y^4}{48} + \frac{y^6 \dot{\Phi}^6(\xi)}{6!} \right) \right] \quad (22)\end{aligned}$$

496 **Compute** $\mathbb{E}_{x \sim \mathcal{N}(0, \lambda^2)} [x^k]$ **and** $\mathbb{E}_{(x,y) \sim \mathcal{N}(0, \Delta)} [x^i y^j]$.

$$\begin{aligned}\mathbb{E}_{x \sim \mathcal{N}(0, \lambda^2)} [x^k] &= \frac{2}{\sqrt{2\pi}\lambda} \int_0^\infty x^k \exp\left(\frac{-x^2}{2\lambda^2}\right) dx \\ &= \frac{2\lambda^k}{\sqrt{2\pi}} \int_0^\infty t^k \exp\left(\frac{-t^2}{2}\right) dt \quad ; x = \lambda t \implies dx = \lambda dt \\ &= \frac{2\lambda^k}{\sqrt{2\pi}} (k-1) \int_0^\infty t^{k-2} \exp\left(\frac{-t^2}{2}\right) dt\end{aligned}$$

$$\text{Thus, } \mathbb{E}_{x \sim \mathcal{N}(0, \lambda^2)} [x^k] = (k-1)\lambda^2 \mathbb{E}_{x \sim \mathcal{N}(0, \lambda^2)} [x^{k-2}] \quad (23)$$

497

$$\begin{aligned}\mathbb{E}_{(x,y) \sim \mathcal{N}(0, \Delta)} [x^i y^j] &= \mathbb{E}_{(x,y) \sim \mathcal{N}(0, \Delta)} [x^i (y \pm \alpha x)^j] \quad ; \alpha = \frac{\mathbb{E}[xy]}{\mathbb{E}[x^2]} \text{ then } x, y - \alpha x \text{ are independent} \\ &= \mathbb{E}_{(x,y) \sim \mathcal{N}(0, \Delta)} \left[x^i \left(\sum_{k=0}^j {}^j C_k (y - \alpha x)^k (\alpha x)^{j-k} \right) \right] \\ &= \mathbb{E}_{(x,y) \sim \mathcal{N}(0, \Delta)} \left[\sum_{k=0}^j {}^j C_k \alpha^k (y - \alpha x)^j x^{k+i} \right] \\ &\stackrel{(e)}{=} \sum_{k=0}^j {}^j C_k \alpha^k \mathbb{E}_{(x,y) \sim \mathcal{N}(0, \Delta)} [x^{k+i}] \mathbb{E}_{(x,y) \sim \mathcal{N}(0, \Delta)} [(y - \alpha x)^j] \quad (24)\end{aligned}$$

498 (e): $x, (y - \alpha x)$ are independent then $x^a, (y - \alpha x)^b$ are also independent.

499 Now, we evaluate (22) using (23) and (24) as follows.

$$\begin{aligned}
(22) &\stackrel{(23),(24)}{=} \frac{1}{4} - \frac{1}{16} (\Sigma_{2ii} + \Sigma_{2jj}) + \frac{1}{64} (\Sigma_{2ii}\Sigma_{2jj} + 2\Sigma_{2ij}^2) \\
&\quad + \frac{1}{32} (\Sigma_{2ii}^2 + \Sigma_{2jj}^2) - \frac{1}{128} (\Sigma_{2ii}^2\Sigma_{2jj} + \Sigma_{2ii}\Sigma_{2jj}^2 + 4\Sigma_{2ij}^2\Sigma_{2ii} + 4\Sigma_{2ij}^2\Sigma_{2jj}) \\
&\quad + \frac{1}{768} (3\Sigma_{2ii}^2\Sigma_{2jj}^2 + 8\Sigma_{2ij}^4 + 24\Sigma_{2ij}^2\Sigma_{2ii}\Sigma_{2jj}) \\
&\quad + \mathbb{E}_{(x,y)\sim\Delta} \left[\frac{x^6\dot{\Phi}^6(\xi)}{6!} \left(\frac{1}{2} - \frac{y^2}{8} + \frac{y^4}{48} + \frac{y^6\dot{\Phi}^6(\xi)}{6!} \right) \right] \\
&\leq \frac{1}{4} - \frac{1}{16} (\Sigma_{2ii} + \Sigma_{2jj}) + \frac{1}{64} (\Sigma_{2ii}\Sigma_{2jj} + 2\Sigma_{2ij}^2) + \frac{1}{32} (\Sigma_{2ii}^2 + \Sigma_{2jj}^2) \\
&\quad - \frac{10}{128}\epsilon^3 + \frac{35}{768}\epsilon^4 + \frac{15}{720}\epsilon^3 \quad ; |\epsilon| \leq \max\{\Delta_{00}, \Delta_{11}\}, \mathbb{E}[x^6] = 15\Delta_{00} \\
&\leq \frac{1}{4} - \frac{1}{16} (\Sigma_{2ii} + \Sigma_{2jj}) + \frac{1}{64} (\Sigma_{2ii}\Sigma_{2jj} + 2\Sigma_{2ij}^2) + \frac{1}{32} (\Sigma_{2ii}^2 + \Sigma_{2jj}^2) + \frac{1}{16}\epsilon^3 \quad (25)
\end{aligned}$$

500 where $|\epsilon| \leq \max\{\Delta_{00}, \Delta_{11}\}$.

501 We get the NTK in Theorem 1 by putting together (25) and (20).

502 **Corollary 1 (Linear GCN).** In this case, $\sigma(x) := x$ and so derivative $\dot{\sigma}(x) = 1$. Consequently, one
503 can derive $\tilde{E}_i = c_\sigma \mathbf{1}_{n \times n}$ from its definition. Therefore, we get NTK for linear GCN in Corollary 1
504 by substituting \tilde{E}_i in general NTK equation in (20).

505 **Corollary 2 (ReLU GCN).** NTK for ReLU GCN is derived by substituting (7) in general NTK
506 equation in (20) as discussed in the corollary.

507 A.2 Proof of NTK for GCN with Skip Connections (Corollary 3 and 4)

508 We derive the NTKs for GCNs with different skip connections, Skip-PC and Skip- α in this section.
509 Before we present the proofs, we note that there are typographical errors in Definitions 1 and 2, and
510 Corollary 4. The corrections in each are listed as follows,

511 1. g_1 in Definition 1 should be $g_1 := \sqrt{\frac{c_\sigma}{h}} S\sigma_s(H_0)$.

512 2. g_1 in Definition 2 should be $g_1 := \sqrt{\frac{c_\sigma}{h}} ((1 - \alpha)S\sigma_s(H_0) + \alpha\sigma_s(H_0))$.

513 3. Σ_i in Corollary 4 should be $\Sigma_i := (1 - \alpha)^2 S E_{i-1} S^T + \alpha^2 \tilde{E}_0$ where $\tilde{E}_0 =$
514 $\mathbb{E}_{f \sim \mathcal{N}(0, \Sigma_0)} [\sigma_s(f)\sigma_s(f)^T]$. We replace E_0 with \tilde{E}_0 in both Corollary 3 and 4 to be clear.

515 We clarify that the errors are only typographical and it did not carry forward to the experiments,
516 thus leaving the empirical results and discussions unaffected. The above mentioned errors will be
517 corrected in the final version of the paper. We derive the NTKs for Skip-PC and Skip- α using these
518 definitions.

519 We observe that the definitions of $g_i \forall i \in [1, d + 1]$ are different for GCN with skip connections from
520 the vanilla GCN. Despite the difference, the definition of gradient with respect to W_i in (4) does not
521 change as g_i in the gradient accounts for the change and moreover, there is no new learnable parameter
522 since the input transformation $H_0 = XT$ where T_{ij} is sampled from $\mathcal{N}(0, 1)$ is not learnable in our
523 setting. Given the fact that the gradient definition holds for GCN with skip connection, the NTK will
524 retain the form from NTK for vanilla GCN as evident from the above derivation. The change in g_i
525 will only affect the co-variance between nodes. Hence, we will derive the co-variance matrix for the
526 discussed skip connections, Skip-PC and Skip- α in the following sections.

527 **Skip-PC: Co-variance between nodes.** The co-variance between nodes p and q in f_1 and f_i are
528 derived below.

$$\begin{aligned}
\mathbb{E} \left[(f_1)_{pk} (f_1)_{qk} \right] &= \left\langle (g_1)_p, (g_1)_q \right\rangle \\
&= \frac{c_\sigma}{h} \left\langle (S\sigma_s(H_0))_p, (S\sigma_s(H_0))_q \right\rangle \\
&= \frac{c_\sigma}{h} \sum_{k=1}^h (S\sigma_s(H_0))_{pk} (S\sigma_s(H_0))_{qk} \\
&\stackrel{h \rightarrow \infty}{=} c_\sigma \mathbb{E} \left[(S\sigma_s(H_0))_{pk} (S\sigma_s(H_0))_{qk} \right] \quad ; \text{law of large numbers} \\
&= S_p \tilde{E}_0 S_q^T \quad ; \tilde{E}_0 = c_\sigma \mathbb{E}_{f \sim \mathcal{N}(0, XX^T)} [\sigma_s(f) \sigma_s(f)^T] \\
&= (\Sigma_1)_{pq} \tag{26}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E} \left[(f_i)_{pk} (f_i)_{qk} \right] &= \left\langle (g_i)_p, (g_i)_q \right\rangle \\
&= \frac{c_\sigma}{h} \left\langle (S(\sigma(f_{i-1}) + \sigma_s(H_0)))_p, (S(\sigma(f_{i-1}) + \sigma_s(H_0)))_q \right\rangle \\
&= \frac{c_\sigma}{h} \sum_{k=1}^h (S\sigma(f_{i-1}) + S\sigma_s(H_0))_{pk} (S\sigma(f_{i-1}) + S\sigma_s(H_0))_{qk} \\
&\stackrel{h \rightarrow \infty}{=} c_\sigma \mathbb{E} \left[(S\sigma(f_{i-1}) + S\sigma_s(H_0))_{pk} (S\sigma(f_{i-1}) + S\sigma_s(H_0))_{qk} \right] \quad ; \text{law of large numbers} \\
&= c_\sigma \left[\mathbb{E} \left[(S\sigma(f_{i-1}))_{pk} (S\sigma(f_{i-1}))_{qk} \right] + \mathbb{E} \left[(S\sigma(f_{i-1}))_{pk} (S\sigma_s(H_0))_{qk} \right] \right. \\
&\quad \left. + \mathbb{E} \left[(S\sigma_s(H_0))_{pk} (S\sigma(f_{i-1}))_{qk} \right] + \mathbb{E} \left[(S\sigma_s(H_0))_{pk} (S\sigma_s(H_0))_{qk} \right] \right] \\
&= S_p E_{i-1} S_q^T + c_\sigma \mathbb{E} \left[(S\sigma(f_{i-1}))_{pk} (S\sigma_s(XW_0))_{qk} \right] \\
&\quad + c_\sigma \mathbb{E} \left[(S\sigma_s(XW_0))_{pk} (S\sigma(f_{i-1}))_{qk} \right] \\
&\quad + c_\sigma \mathbb{E} \left[\sum_{r=1}^n \sum_{s=1}^n S_{pr} S_{qs} \sigma_s(XW_0)_{rk} \sigma_s(XW_0)_{sk} \right] \\
&\stackrel{(f)}{=} S_p E_{i-1} S_q^T + c_\sigma S_p \mathbb{E} [\sigma_s(XW_0)_{rk} \sigma_s(XW_0)_{sk}] S_q^T \\
&= S_p E_{i-1} S_q^T + S_p \tilde{E}_0 S_q^T \\
&= S_p E_{i-1} S_q^T + (\Sigma_i)_{pq} \\
&= (\Sigma_i)_{pq} \tag{27}
\end{aligned}$$

529 (f): $\mathbb{E} \left[(S\sigma(f_{i-1}))_{pk} (S\sigma_s(XW_0))_{qk} \right]$ and $\mathbb{E} \left[(S\sigma_s(XW_0))_{pk} (S\sigma(f_{i-1}))_{qk} \right]$ evaluate to 0
530 by conditioning on W_0 first and rewriting the expectation based on this conditioning.
531 The terms within expectation are independent when conditioned on W_0 , and hence it is
532 $\mathbb{E}_{W_0} \left[\mathbb{E}_{\Sigma_{i-1}|W_0} \left[(S\sigma(f_{i-1}))_{pk} | W_0 \right] \mathbb{E}_{\Sigma_{i-1}|W_0} \left[(S\sigma_s(XW_0))_{qk} | W_0 \right] \right]$ by taking h in W_0 going to in-
533 finity first. Here, $\mathbb{E}_{\Sigma_{i-1}|W_0} \left[(S\sigma_s(XW_0))_{qk} | W_0 \right] = 0$.

534 We get the co-variance matrix for all pairs of nodes $\Sigma_1 = S\tilde{E}_0 S^T$ and $\Sigma_i = S E_{i-1} S^T + \Sigma_1$ from
535 (26) and (27).

536 **Skip- α : Co-variance between nodes.** Let p and q be two nodes and the co-variance between p and
537 q in f_1 and f_i are derived below.

$$\begin{aligned}
\mathbb{E} \left[(f_1)_{pk} (f_1)_{qk} \right] &= \left\langle (g_1)_p, (g_1)_q \right\rangle \\
&= \frac{c_\sigma}{h} \sum_{k=1}^h ((1-\alpha)S\sigma_s(H_0) + \alpha\sigma_s(H_0))_{pk} ((1-\alpha)S\sigma_s(H_0) + \alpha\sigma_s(H_0))_{qk} \\
&\stackrel{h \rightarrow \infty}{=} c_\sigma \mathbb{E} \left[((1-\alpha)S\sigma_s(H_0) + \alpha\sigma_s(H_0))_{pk} ((1-\alpha)S\sigma_s(H_0) + \alpha\sigma_s(H_0))_{qk} \right] \\
&= c_\sigma \left[(1-\alpha)^2 \mathbb{E} \left[(S\sigma_s(H_0))_{pk} (S\sigma_s(H_0))_{qk} \right] \right. \\
&\quad \left. + (1-\alpha)\alpha \left(\mathbb{E} \left[(S\sigma_s(H_0))_{pk} (\sigma_s(H_0))_{qk} \right] + \mathbb{E} \left[(S\sigma_s(H_0))_{qk} (\sigma_s(H_0))_{pk} \right] \right) \right. \\
&\quad \left. + \alpha^2 \mathbb{E} \left[(\sigma_s(H_0))_{pk} (\sigma_s(H_0))_{qk} \right] \right] \\
&= (1-\alpha)^2 S_p \tilde{E}_0 S_{.q}^T + (1-\alpha)\alpha \left(S_p \cdot (\tilde{E}_0)_{.q} + (\tilde{E}_0)_p \cdot S_{.q}^T \right) + \alpha^2 (\tilde{E}_0)_{pq} \\
&= (\Sigma_1)_{pq} \tag{28}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E} \left[(f_i)_{pk} (f_i)_{qk} \right] &= \left\langle (g_i)_p, (g_i)_q \right\rangle \\
&= \frac{c_\sigma}{h} \sum_{k=1}^h ((1-\alpha)S\sigma(f_{i-1}) + \alpha\sigma_s(H_0))_{pk} ((1-\alpha)S\sigma(f_{i-1}) + \alpha\sigma_s(H_0))_{qk} \\
&\stackrel{h \rightarrow \infty}{=} c_\sigma \mathbb{E} \left[((1-\alpha)S\sigma(f_{i-1}) + \alpha\sigma_s(H_0))_{pk} ((1-\alpha)S\sigma(f_{i-1}) + \alpha\sigma_s(H_0))_{qk} \right] \\
&= c_\sigma \left[(1-\alpha)^2 \mathbb{E} \left[(S\sigma(f_{i-1}))_{pk} (S\sigma(f_{i-1}))_{qk} \right] + \alpha^2 \mathbb{E} \left[(\sigma_s(H_0))_{pk} (\sigma_s(H_0))_{qk} \right] \right. \\
&\quad \left. + (1-\alpha)\alpha \left(\mathbb{E} \left[(S\sigma(f_{i-1}))_{pk} (\sigma_s(H_0))_{qk} \right] + \mathbb{E} \left[(\sigma_s(H_0))_{pk} (S\sigma(f_{i-1}))_{qk} \right] \right) \right] \\
&\stackrel{(g)}{=} (1-\alpha)^2 S_p E_{i-1} S_{.q}^T + \alpha^2 (\tilde{E}_0)_{pq} = (\Sigma_i)_{pq} \tag{29}
\end{aligned}$$

538 (g): same argument as (f) in derivation of Σ_i in Skip-PC.

539 We get the co-variance matrix for all pairs of nodes $\Sigma_1 = (1-\alpha)^2 S \tilde{E}_0 S^T + \alpha(1-$
540 $\alpha) (S \tilde{E}_0 + \tilde{E}_0 S^T) + \alpha^2 \tilde{E}_0$ and $\Sigma_i = (1-\alpha)^2 S E_{i-1} S^T + \alpha^2 \tilde{E}_0$ from (28) and (29).

541 A.3 Normalisation constant c_σ (Remark 1 and 2).

542 We derive the normalisation constant c_σ loosely, as the purpose of c_σ is to preserve the input norm
543 approximately. We focus on general form of a network with skip connection (not GCN in particular),
544 where the output vector of size h from any hidden layer l with weight matrix $W \in \mathbb{R}^{h \times h}$ and
545 transformed input vector X_0 of size h can be written as $g_l := \sqrt{\frac{c_\sigma}{h}} (\sigma(Wg_{l-1}) + X_0) \in \mathbb{R}^{h \times 1}$. The
546 role of the normalisation constant c_σ is to maintain $\|g_l\|_2 \simeq \|X_0\|_2$ and is derived as follows.

$$\begin{aligned}
\|X_0\|_2^2 &= \|g_l\|_2^2 = \frac{c_\sigma}{h} \sum_{k=1}^h (\sigma(Wg_{l-1}) + X_0)_k^2 \\
\|X_0\|_2^2 &= c_\sigma \mathbb{E} \left[(\sigma(Wg_{l-1})_k)^2 + (X_0)_k^2 + 2\sigma(Wg_{l-1})_k (X_0)_k \right] \quad ; h \rightarrow \infty \\
\|X_0\|_2^2 &= c_\sigma \mathbb{E}_{u \sim \mathcal{N}(0, \|X_0\|_2^2)} \left[(\sigma(u))^2 \right] + \|X_0\|_2^2 \quad ; \mathbb{E} [\sigma(Wg_{l-1})_k (X_0)_k] = 0 \\
c_\sigma &= \left(\mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[(\sigma(u))^2 \right] + 1 \right)^{-1} \quad ; \text{normalised } X_0 \tag{30}
\end{aligned}$$

547 We use this c_σ for GCN with skip connection in our work and it evaluates to $2/3$ for $\sigma(x) := \text{ReLU}(x)$
548 in GCN as stated in Remark 2. The evident change for a network without skip connection is to
549 not add X_0 in $g_l := \sqrt{\frac{c_\sigma}{h}} \sigma(Wg_{l-1})$ and consequently by following the proof, we get $c_\sigma =$
550 $\left(\mathbb{E}_{u \sim \mathcal{N}(0,1)} [(\sigma(u))^2] \right)^{-1}$ as mentioned in Remark 1.

551 B Additional Experimental Results

552 B.1 Datasets for binary node classification

553 Since the considered datasets *Cora*, *Citeseer* and *WebKB* are for multi-class node classification, we
554 converted the datasets to have binary class by grouping the classes into two sets. Table 1 shows the
555 label grouping for each dataset and total number of nodes with the grouped labels respectively. The
556 classes in all the datasets are well balanced and sensible to learn for binary classification problem
557 which is proved from the performance of a simple graph neural network like linear vanilla GCN. The
558 train-test split for each dataset is 708 and 2000 nodes for Cora, 312 and 2000 for Citeseer, and 377
559 and 500 for WebKB for all the experiments.

	Cora		Citeseer		WebKB	
	Class Groups	#nodes	Class Groups	#nodes	Class Groups	#nodes
Class 1	Neural_Networks Theory Probabilistic_Methods	1595	Agents AI ML	1435	student	415
Class 2	Case_Based Rule_Learning Reinforcement Genetic_Algorithms	1103	DB IR HCI	1877	faculty staff course project	462
Total		2708		3312		877

Table 1: Class grouping in datasets for binary node classification.

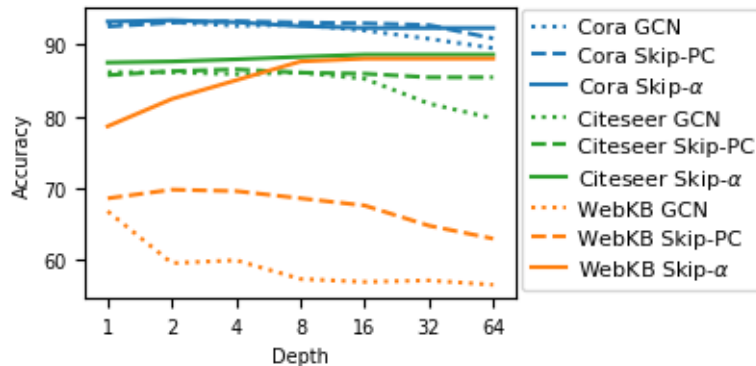


Figure 5: Performance validation of vanilla GCN, Skip-PC and Skip- α with $\sigma(\cdot) := \text{ReLU}$, $\sigma_s(\cdot) := \text{ReLU}$ and $\alpha = 0.2$ using the respective NTKs.

560 **B.2 Vanilla GCN vs GCN with Skip Connections**

561 We established ReLU GCN is preferred over linear in Section 2 and ReLU for the input transformation
 562 in Section 3. Hence, we focus on $\sigma(\cdot) := \text{ReLU}$ and $\sigma_s(\cdot) := \text{ReLU}$ with $\alpha = 0.2$ for Skip- α to
 563 validate the performance of vanilla GCN and GCN with skip connections, Skip-PC and Skip- α . We
 564 use the respective NTKs to validate the performance. Figure 5 shows that GCN with skip connection
 565 outperforms vanilla GCN even in deeper depths, and Skip- α gives better performance than Skip-PC
 566 with depth.

567 **Note.** In Figure 2, the performance of Skip-PC with ReLU $\sigma_s(\cdot)$ evaluated on Citeseer when depth
 568 = 32 is different from what is plotted due to some numerical precision error. We evaluated the
 569 performance at depth 30, 31 and used it to plot.

570 **B.3 Convergence of NTK with depth - Cora, Citeseer, WebKB**

571 We presented the convergence of NTK with depth for ReLU GCN with and without skip connections
 572 evaluated on Cora dataset in Figure 4. Here, we present the convergence plot for Linear GCN
 573 evaluated on Cora and all discussed linear and ReLU networks evaluated Citeseer and WebKB. The
 574 observation is similar to the discussion in Section 5. Figures 6, 7 and 8 show the convergence plots
 575 for linear GCN evaluated on Cora, ReLU and linear GCNs with and without skip connections for
 576 Citeseer and WebKB, respectively.

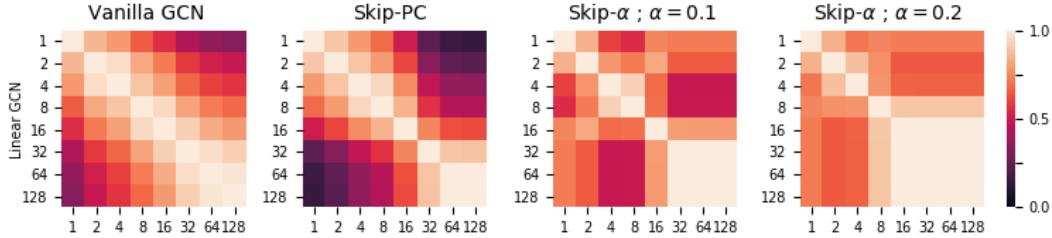


Figure 6: Convergence of NTK with depth for all the discussed linear architectures evaluated on Cora dataset.

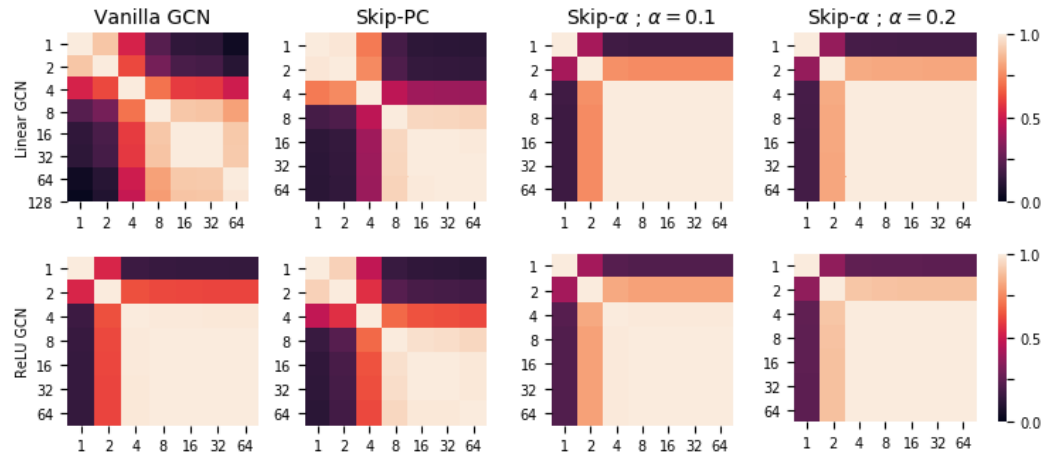


Figure 7: Convergence of NTK with depth for all the discussed linear and ReLU architectures evaluated on Citeseer dataset.

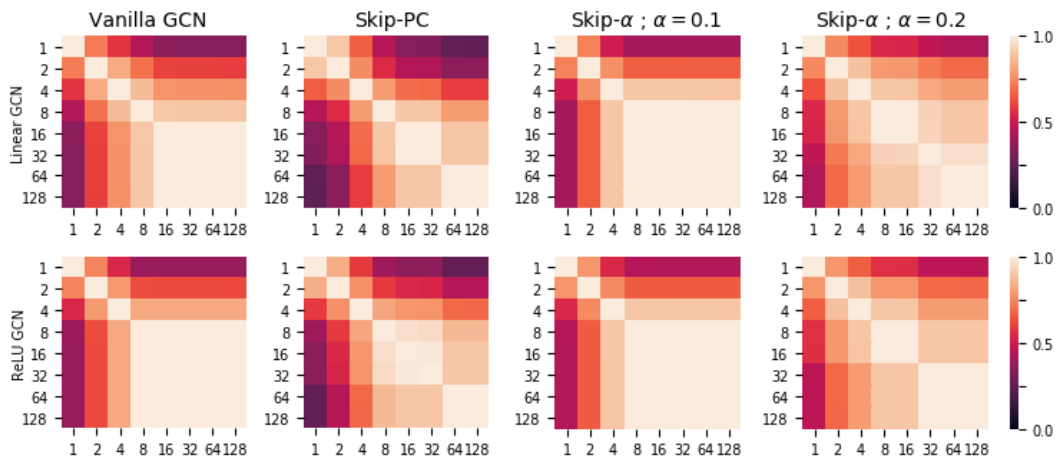


Figure 8: Convergence of NTK with depth for all the discussed linear and ReLU architectures evaluated on WebKB dataset.