# Deeper Insights into Deep Graph Convolutional Networks: Stability and Generalization

Guangrui Yang, Ming Li\*, Han Feng\*, Xiaosheng Zhuang

#### **Abstract**

Graph convolutional networks (GCNs) have emerged as powerful models for graph learning tasks, exhibiting promising performance in various domains. While their empirical success is evident, there is a growing need to understand their essential ability from a theoretical perspective. Existing theoretical research has primarily focused on the analysis of single-layer GCNs, while a comprehensive theoretical exploration of the stability and generalization of deep GCNs remains limited. In this paper, we bridge this gap by delving into the stability and generalization properties of deep GCNs, aiming to provide valuable insights by characterizing rigorously the associated upper bounds. Our theoretical results reveal that the stability and generalization of deep GCNs are influenced by certain key factors, such as the maximum absolute eigenvalue of the graph filter operators and the depth of the network. Our theoretical studies contribute to a deeper understanding of the stability and generalization properties of deep GCNs, potentially paving the way for developing more reliable and well-performing models.

#### **Index Terms**

Graph convolutional n	etworks (GCNs); Generalization ga	p; Deep GCNs; Uniform stability.
		<b>•</b>

# 1 Introduction

RAPH-structured data is pervasive across diverse domains, including knowledge graphs, traffic networks, and social networks to name a few [1], [2]. Several pioneering works [3], [4] introduced the initial concept of graph neural networks (GNNs), incorporating recurrent mechanisms and necessitating neural network parameters to define contraction mappings. Concurrently, Micheli [5] introduced the neural network for graphs, commonly referred to as NN4G, over a comparable timeframe. It is worth noting that the NN4G diverges from recurrent mechanisms and instead employs a feed-forward architecture, exhibiting similarities to contemporary GNNs. In recent years, (contemporary) GNNs have gained significant attention as an effective methodology for modeling graph data [6]–[11]. To obtain a comprehensive understanding of GNNs and deep learning for graphs, we refer the readers to relevant survey papers for an extensive overview [12]–[15].

Among the various GNN variants, one of the most powerful and frequently used GNNs is graph convolutional networks (GCNs). A widely accepted perspective posits that GCNs can be regarded as an extension or generalization of traditional spatial filters, which are commonly employed in Euclidean data analysis, to the realm of non-Euclidean data. Due to its success on non-Euclidean data, GCN has attracted widespread attention on its theoretical exploration. Recent works on GCNs includes understanding over-smoothing [16]–[19], interpretability and explainability [20]–[24], expressiveness [25]–[27], and generalization [28]–[41]. In this paper, we specifically address the generalization of GCNs to provide a bound on their generalization gap.

Investigating the generalization of GCNs is essential in understanding its underlying working principles and capabilities from a theoretical perspective. However, the theoretical establishment in this area is still in its infancy. In recent work [36], Verma and Zhang provided a novel technique based on algorithmic stability to investigate the generalization capability of single-layer GCNs in semi-supervised learning tasks. Their results indicate that the stability of a single-layer GCN trained with the stochastic gradient descent (SGD) algorithm is dependent on the largest absolute eigenvalue of graph filter operators. This finding highlights the crucial role of graph filters in determining the generalization capability of single-layer GCNs, providing guidance for designing effective graph filters for these networks. On the other hand, a number of prior

1

This work was supported in part by the National Natural Science Foundation of China (No. U21A20473, No. 62536006, No. 62172370). M. Li also acknowledged the support from the "Pioneer" and "Leading Goose" R&D Program of Zhejiang (No. 2024C03262). G. Yang acknowledged the support from the Opening Project of Guangdong Province Key Laboratory of Computational Science at the Sun Yat-sen University (No. 2024008). H. Feng was supported in part by the Research Grants Council of Hong Kong (Project no. CityU 11303821 ,and CityU 11315522). X. Zhuang was supported in part by the Research Grants Council of Hong Kong (Project no. CityU 11309122, CityU 11301224, and CityU 11300825).

Guangrui Yang is with the Department of Mathematics, College of Mathematics and Informatics, South China Agricultural University, Guangzhou, China. (e-mail: yanggrui@mail2.sysu.edu.cn).

Ming Li is with Zhejiang Key Laboratory of Intelligent Education Technology and Application, Zhejiang Normal University, Jinhua, China (e-mail: mingli@zjnu.edu.cn).

Han Feng and Xiaosheng Zhuang are with Department of Mathematics, City University of Hong Kong, Hong Kong, China (e-mail: hanfeng@cityu.edu.hk; xzhuang7@cityu.edu.hk).

<sup>\*</sup>Corresponding authors

studies have shown that deep GCNs possess greater expressive power than their single-layer counterparts. Consequently, it is essential to extend the generalization results of single-layer GCNs to their multi-layer counterparts. This will help us understand the effect of factors (e.g., graph filters, number of layers) on the generalization capability of deep GCNs.

In this paper, we investigate the generalization properties of deep GCNs. Building on the stability framework of [36], we analyze the uniform stability of deep GCNs in semi-supervised learning, while developing a more refined theoretical treatment suited to deep architectures. Our analysis reveals a strong connection between the generalization gap of deep GCNs and the characteristics of the graph filter, particularly the number of layers. In particular, we show that when the maximum absolute eigenvalue (or the largest singular value) of the graph filter operator remains invariant with respect to graph size, the generalization gap diminishes asymptotically at a rate of  $O(1/\sqrt{m})$  as the training sample size m grows. This result explains why normalized graph filters generally outperform non-normalized ones in deep GCNs. Furthermore, our findings indicate that increasing depth can enlarge the generalization gap and consequently degrade performance, thereby offering theoretical guidance for selecting an appropriate number of layers when designing deep GCNs. We then empirically validate our theoretical results through experiments on three benchmark datasets: Cora, Citeseer, and Pubmed, demonstrating strong consistency between theory and practice. In addition, we further discuss how our theoretical framework extends to advanced architectures, including GCNII [42] and Graph Transformer [43], thereby highlighting its broader applicability and its potential to inspire future theoretical studies on more complex GNN variants.

The key contributions of our paper are as follows:

- We establish the uniform stability of deep GCNs trained with SGD, thereby extending the earlier results on single-layer GCNs presented in [36].
- We provide a rigorous upper bound for the generalization gap of deep GCNs and highlight the key factors that govern their generalization ability. Moreover, we further discuss how our theoretical framework extends naturally to advanced GNN architectures, including GCNII and Graph Transformer models.
- We conduct empirical studies on three benchmark datasets for node classification, which strongly validate our theoretical findings regarding the influence of graph filters, as well as the depth and width of deep GCNs.

The remainder of this paper is organized as follows. In Section 2, an overview of prior studies on the generalization of GCNs (or generic GNNs) is presented, along with a comparative analysis highlighting the similarities and distinctions between our work and previous research. Section 3 offers an exposition of the essential concepts. The primary findings of this paper are given in Section 4. Experimental studies designed to validate our theoretical findings are presented in Section 5. In Section 6, we discuss how our findings extend to advanced GNN architectures, including GCNII and Graph Transformer models. Section 7 concludes the paper with additional remarks. The detailed proofs of our theoretical results are deferred to the **Appendix** section.

# 2 RELATED WORK

Theoretical studies on the generalization capability of GCNs mainly employ three methodologies: Vapnik–Chervonenkis (VC) dimension [30], [34], Rademacher complexity [31]–[35], and algorithmic stability [36], [37], [44], [45]. Other approaches include PAC-Bayesian theory [38], [39], neural tangent kernels (NTKs) [40], [41], algorithm alignment [46], [47], and methods from statistical physics and random matrix theory [48]. For a broader perspective, we refer readers to the recent survey [49], which provides a comprehensive overview of generalization theory for message-passing GNNs.

VC-Dimension and Rademacher Complexity. Scarselli et al. [30] study the generalization capability of GNNs by deriving upper bounds on the growth order of their VC-dimension. While VC-dimension is a classical tool for establishing learning bounds, it does not capture the structure of the underlying graph. Similarly, [34] provides VC-dimension–based error bounds for GNNs, but the results are trivial and fail to reflect the benefits of degree normalization. To address graph-specific effects, Esser et al. [34] analyze upper bounds using transductive Rademacher complexity (TRC), highlighting how graph convolutions and network architectures influence generalization. Tang et al. [35] establish high-probability generalization bounds for popular GNNs via TRC-based analysis of transductive SGD. However, their bounds scale with the parameter dimension, limiting tightness for large models.

Algorithmic Stability. Beyond capacity-based measures, algorithmic stability serves as an important framework for understanding GNN generalization. Building on the work of Hardt et al. [50], Verma and Zhang [36] show that one-layer GCNs exhibit uniform stability and provide generalization bounds that scale with the largest absolute eigenvalue of the graph filter operator. Extending this line, Liu et al. [44] analyze the stability of single-layer GCNs trained with an SGD-proximal algorithm under  $\ell_p$ -regularization, yielding a more refined theoretical understanding. These studies, however, remain restricted to single-layer architectures. Cong et al. [51] examine GNNs under uniform transductive stability, showing that deeper models improve stability and reduce generalization error, whereas our work adopts a different stability formulation. Ng and Yip [37] investigate stability and generalization in two-layer GCNs under an eigen-domain formulation, relying on spectral graph convolution [52]. Because this formulation requires computationally expensive eigendecomposition of the graph Laplacian, it does not scale to large node-classification tasks. Within this methodological line, the closest studies to ours are [36] and [37], but our analysis focuses on deep GCNs without assuming a spectral-based formulation.

Other Methodologies. Alternative perspectives on GNN generalization also exist. The pioneering work of [38] introduces PAC-Bayesian analysis for GCNs and message-passing neural networks, later extended in [39] to provide tighter bounds linked to the graph diffusion matrix. The NTK framework introduced by [40] enables analysis of infinitely wide GNNs trained by gradient descent, with [41] extending this framework to multi-layer settings. However, NTK-based analyses typically focus on graph classification rather than the more challenging transductive node-classification setting. Additional work explores distinct theoretical frameworks, including topology-sampling techniques [53], analysis on large random graphs [54], and NTK-based loss landscape analysis of wide GCNs [55]. For further perspectives, we refer readers to the survey [56], which synthesizes emerging theoretical approaches to characterizing GNN capabilities.

#### 3 Preliminaries and Notations

In this section, we describe the problem setup considered in this paper and review fundamental concepts of uniform stability for training algorithms, which form the basis of our subsequent analysis. For clarity, we first summarize the main symbols used in this paper in the table below.

TABLE 1 Frequently used notations.

Notation	Description
$g(\mathbf{L})$	graph filter operator used in the considered deep GCNs
$C_g$	the 2-norm of $g(\mathbf{L})$ , i.e., $C_g := \ g(\mathbf{L})\ _2$
$C_{\mathbf{X}}$	Frobenius norm of the input feature $\mathbf{X}$ , i.e., $C_{\mathbf{X}} := \ \mathbf{X}\ _F$
K	number of hidden layers of the considered deep GCNs
$\alpha_{\sigma}, v_{\sigma}$	parameters w.r.t the continuity of activation function $\sigma(\cdot)$
$ abla \sigma$	the derivative of activation function $\sigma(\cdot)$
$\alpha_\ell, v_\ell$	parameters w.r.t the continuity of the loss function $\ell(\cdot,\cdot)$
M	the upper bound of loss function $\ell(\cdot,\cdot)$
$\mathcal{A}_{\mathcal{S}}$	the learning algorithm for deep GCNs trained on dataset ${\cal S}$
m	the number of samples in the trained dataset ${\cal S}$
$\eta$	the learning rate of $\mathcal{A}_{\mathcal{S}}$
T	number of iterations for training $\mathcal{A}_{\mathcal{S}}$ using SGD
$\mu_m$	the uniform stability of a learning algorithm $\mathcal{A}_{\mathcal{S}}$
$\boldsymbol{\delta}_{\mathbf{x}}$	the indicator vector with respect to node $\mathbf{x}$
$oldsymbol{\delta}_i$	the indicator vector with respect to index $i$
$\mathbf{X}^{(k)}$	the output feature matrix of the $k$ -th layer
$\triangle \mathbf{X}^{(k)}$	the variation of $\mathbf{X}^{(k)}$ in two GCNs
$\mathbf{W}^{(k)}$	the parameter matrix specific to the $k$ -th layer
B	upper bound for 2-norm of $\{\mathbf{W}^{(1)},\ldots,\mathbf{W}^{(K)},\mathbf{w}\}$
$\triangle \mathbf{W}^{(k)}$	the variation of $\mathbf{W}^{(k)}$ in two GCNs
$\triangle \theta$	$\triangle \theta = \{ \triangle \mathbf{W}^{(1)}, \dots, \triangle \mathbf{W}^{(K)}, \triangle \mathbf{w} \}$
$\mathbf{W}_t^{(k)}$	the learnt $\mathbf{W}^{(k)}$ trained after $t$ iterations
$ riangle \mathbf{W}_t^{(k)}$	the variation of $\mathbf{W}_t^{(k)}$ of two GCNs trained after $t$ iterations
$\triangle  heta_t$	$\triangle \theta_t = \{ \triangle \mathbf{W}_t^{(1)}, \dots, \triangle \mathbf{W}_t^{(K)}, \triangle \mathbf{w}_t \}$

# 3.1 Deep Graph Convolutional Networks

Let  $\mathcal{G}=(\mathcal{V},\mathcal{E},\mathbf{A})$  denote an undirected graph with a node set  $\mathcal{V}$  of size N, an edge set  $\mathcal{E}$  and the adjacency matrix  $\mathbf{A}\in\mathbb{R}^{N\times N}$ . As usual,  $\mathbf{L}:=\mathbf{D}-\mathbf{A}$  is denoted as its conventional graph Laplacian, where  $\mathbf{D}\in\mathbb{R}^{N\times N}$  signifies the degree diagonal matrix. Furthermore,  $g(\mathbf{L})\in\mathbb{R}^{N\times N}$  represents a graph filter and is defined as a function of  $\mathbf{L}$  (or its normalized versions). We denote by  $C_g=\|g(\mathbf{L})\|_2$  the maximum absolute eigenvalue of a symmetric filter  $g(\mathbf{L})$  or the maximum singular value of an asymmetric  $g(\mathbf{L})$ .

We denote by  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^{\top} \in \mathbb{R}^{N \times d_0}$  the input features ( $d_0$  stands for input dimension) and  $\mathbf{x}_j \in \mathbb{R}^{d_0}$  the node feature of node j, while  $C_{\mathbf{X}} = \|\mathbf{X}\|_F$  represents the Frobenius norm of  $\mathbf{X}$ . For the input feature  $\mathbf{X}$ , a deep GCN with  $g(\mathbf{L})$  updates the representation as follows:

$$\mathbf{X}^{(k)} = \sigma(g(\mathbf{L})\mathbf{X}^{(k-1)}\mathbf{W}^{(k)}), \quad k = 1, 2, \dots, K,$$

where  $\mathbf{X}^{(k)} \in \mathbb{R}^{N \times d_k}$  is the output feature matrix of the k-th layer with  $\mathbf{X}^{(0)} = \mathbf{X}$ , the matrix  $\mathbf{W}^{(k)} \in \mathbb{R}^{d_{k-1} \times d_k}$  represents the trained parameter matrix specific to the k-th layer. The function  $\sigma(\cdot)$  denotes a nonlinear activation function applied within the GCN model. For simplicity, we set a final output in a single dimension, that is, the final output label of N nodes is given by

$$\mathbf{y} = \sigma(g(\mathbf{L})\mathbf{X}^{(K)}\mathbf{w}),\tag{1}$$

where  $\mathbf{y} \in \mathbb{R}^N$  and  $\mathbf{w} \in \mathbb{R}^{d_K}$ .

As defined above, the deep GCN (1) with learnable parameters

$$\theta = \{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \dots, \mathbf{W}^{(K)}, \mathbf{w}\}\$$

is a K + 1 layers GCN with K hidden layers and a final output layer, and in the case of K = 0, it degenerates into the single-layer GCN studied in [36].

# 3.2 The SGD Algorithm

We denote by  $\mathcal{D}$  the unknown joint distribution of input features and output labels. Let

$$\mathcal{S} := \left\{ (\mathbf{x}_j, y_j) \right\}_{j=1}^m$$

be the training set i.i.d sampled from  $\mathcal{D}$  and  $\mathcal{A}_{\mathcal{S}}$  be a learning algorithm for a deep GCN trained on  $\mathcal{S}$ . For a deep GCN model (1) with parameters  $\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(K)}, \mathbf{w}\}$ , denote  $\mathcal{A}_{\mathcal{S}}(\mathbf{x}) = f(\mathbf{x}|\theta_{\mathcal{S}}) = \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})\mathbf{X}^{(K)}\mathbf{w})$  as the output of node  $\mathbf{x}$ , where  $\theta_{\mathcal{S}}$  is the corresponding learned parameter and  $\boldsymbol{\delta}_{\mathbf{x}}$  is the indicator vector with respect to node  $\mathbf{x}$ . For a loss function  $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}_+$ , the generalization error or risk  $R(\mathcal{A}_{\mathcal{S}})$  is defined by

$$R(\mathcal{A}_{\mathcal{S}}) := \mathbb{E}_{\mathbf{z}} \Big[ \ell(f(\mathbf{x}|\theta_S), y) \Big],$$

where the expectation is taken over  $\mathbf{z} = (\mathbf{x}, y) \sim \mathcal{D}$ , and the empirical error or risk  $R_{emp}(\mathcal{A}_{\mathcal{S}})$  is

$$R_{emp}(\mathcal{A}_{\mathcal{S}}) := \frac{1}{m} \sum_{j=1}^{m} \ell(f(\mathbf{x}_j | \theta_S), y_j).$$

When considering a randomized algorithm  $A_S$ ,

$$\epsilon_{gen}(\mathcal{A}_{\mathcal{S}}) := \mathbb{E}_{\mathcal{A}} \Big[ R(\mathcal{A}_{\mathcal{S}}) - R_{emp}(\mathcal{A}_{\mathcal{S}}) \Big]$$
 (2)

gives the generalization gap between the generalization error and the empirical error, where the expectation  $\mathbb{E}_{\mathcal{A}}$  corresponds to the inherent randomness of  $\mathcal{A}_{\mathcal{S}}$ .

In this paper,  $\mathcal{A}_{\mathcal{S}}$  is considered to be the algorithm given by the SGD algorithm. Following the approach employed in [36], our analysis focuses solely on the randomness inherent in  $\mathcal{A}_{\mathcal{S}}$  arising from the SGD algorithm, while disregarding the stochasticity introduced by parameter initialization. The SGD algorithm for a deep GCN (1) aims to optimize its empirical error on a dataset  $\mathcal{S}$  by updating parameters iteratively. For  $t \in \mathbb{N}$  and considering the parameters  $\theta_{t-1}$  obtained after t-1 iterations, the t-th iteration of SGD involves randomly drawing a sample  $(\mathbf{x}_t, y_t)$  from the dataset  $\mathcal{S}$ . Subsequently, parameters  $\theta$  are iteratively updated as follows:

$$\theta_t = \theta_{t-1} - \eta \nabla_{\theta} \ell(f(\mathbf{x}_t | \theta_{t-1}), y_t), \tag{3}$$

with the learning rate  $\eta > 0$ .

#### 3.3 Uniform Stability

For the sake of estimating the generalization gap  $\epsilon_{gen}(A_S)$  of  $A_S$ , we invoke the notion of uniform stability of  $A_S$  as adopted in [36], [57].

Let

$$S^{\setminus i} = \{(\mathbf{x}_j, y_j)\}_{j=1}^{i-1} \cup \{(\mathbf{x}_j, y_j)\}_{j=i+1}^m$$

be the dataset obtained by removing the i-th data point in S, and

$$S^{i} = \left\{ (\mathbf{x}_{j}, y_{j}) \right\}_{j=1}^{i-1} \cup \left\{ (\mathbf{x}'_{i}, y'_{i}) \right\} \cup \left\{ (\mathbf{x}_{j}, y_{j}) \right\}_{j=i+1}^{m}$$

the dataset obtained by replacing the i-th data point in S. Then, the formal definition of uniform stability of a randomized algorithm  $\mathcal{A}_{S}$  is given in the following.

**Definition 1 (Uniform Stability** [36]). A randomized algorithm  $A_S = f(\mathbf{x}|\theta_S)$  is considered to be  $\mu_m$ -uniformly stable in relation to a loss function  $\ell$  when it fulfills the following condition:

$$\sup_{S,\sigma} \left| \mathbb{E}_{\mathcal{A}} [\ell(\hat{y}, y)] - \mathbb{E}_{\mathcal{A}} [\ell(\hat{y}', y)] \right| \le \mu_m, \tag{4}$$

where  $\mathbf{z} = (\mathbf{x}, y) \sim \mathcal{D}$ ,  $\hat{y} = f(\mathbf{x}|\theta_S)$  and  $\hat{y}' = f(\mathbf{x}|\theta_{S\setminus i})$ .

As shown in Definition 1,  $\mu_m$  indicates a bound on how much the variation of the training set S can influence the output of  $A_S$ . It further implies the following property:

$$\sup_{\mathcal{S}, \mathbf{z}} \left| \mathbb{E}_{\mathcal{A}} \left[ \ell(\hat{y}, y) \right] - \mathbb{E}_{\mathcal{A}} \left[ \ell(\hat{y}', y) \right] \right| \le 2\mu_m, \tag{5}$$

where  $\mathbf{z} = (\mathbf{x}, y) \sim \mathcal{D}$ ,  $\hat{y} = f(\mathbf{x}|\theta_S)$  and  $\hat{y}' = f(\mathbf{x}|\theta_{S^i})$ .

Moreover, it is shown that the uniform stability of a learning algorithm  $\mathcal{A}_{\mathcal{S}}$  can yield the following upper bound on the generalization gap  $\epsilon_{qen}(\mathcal{A}_{\mathcal{S}})$ .

*Lemma 1 (Stability Guarantees* [36]). Suppose that a randomized algorithm  $\mathcal{A}_{\mathcal{S}}$  is  $\mu_m$ -uniformly stable with a bounded loss function  $\ell$ . Then, with a probability of at least  $1 - \delta$ , considering the random draw of  $\mathcal{S}, \mathbf{z}$  with  $\delta \in (0,1)$ , the following inequality holds for the expected value of the generalization gap:

$$\epsilon_{gen}(\mathcal{A}_{\mathcal{S}}) \le 2\mu_m + \left(4m\mu_m + M\right)\sqrt{\frac{\log\frac{1}{\delta}}{2m}},$$

where M is an upper bound of the loss function  $\ell$ , i.e.,  $0 \le \ell(\cdot, \cdot) \le M$ .

# 4 MAIN RESULTS

This section presents an established upper bound on the generalization gap  $\epsilon_{gen}(\mathcal{A}_{\mathcal{S}})$  as defined in (2) for deep GCNs trained using the SGD algorithm. Notably, this generalization bound, derived from a meticulous analysis of the comprehensive back-propagation algorithm, demonstrates the enhanced insight gained through the utilization of SGD.

# 4.1 Assumptions

First, we make some assumptions about the considered deep GCN model (1), which are necessary to derive our results. **Assumption 1.** The activation function  $\sigma : \mathbb{R} \to \mathbb{R}$  is assumed to satisfy the following:

1)  $\alpha_{\sigma}$ -Lipschitz:

$$|\sigma(x) - \sigma(y)| \le \alpha_{\sigma}|x - y|, \ \forall \ x, y \in \mathbb{R}.$$

2)  $\nu_{\sigma}$ -smooth:

$$|\nabla \sigma(x) - \nabla \sigma(y)| \le \nu_{\sigma} |x - y|, \ \forall \ x, y \in \mathbb{R}.$$

3)  $\sigma(0) = 0$ .

With these assumptions, the derivative of  $\sigma$ , denoted by  $\nabla \sigma$ , is bounded, i.e.,  $|\nabla \sigma(\cdot)| \leq \alpha_{\sigma}$ , and  $||\sigma(\mathbf{X})||_F \leq \alpha_{\sigma} ||\mathbf{X}||_F$  holds for any matrix  $\mathbf{X}$ . It can be easily verified that activation functions such as ELU and tanh satisfy the above assumptions.

**Assumption 2.** Let  $\hat{y}$  and y be the predicted and true labels, respectively. We denote the loss function  $\ell : [y_{\min}, y_{\max}] \times [y_{\min}, y_{\max}] \to \mathbb{R}$  by  $\ell(\hat{y}, y)$ . Similar to [37], we adopt the following assumptions for  $\ell$ .

- 1) The loss function  $\ell$  exhibits continuity with respect to the variables  $(\hat{y}, y)$  and possesses continuous differentiability with respect to  $\hat{y}$ .
- 2) The loss function  $\ell$  satisfies  $\alpha_{\ell}$ -Lipschitz with respect to  $\hat{y}$ :

$$|\ell(\hat{y}, y) - \ell(\hat{y}', y)| \le \alpha_{\ell} |\hat{y} - \hat{y}'|, \ \forall \ \hat{y}, \hat{y}', y \in [y_{\min}, y_{\max}].$$

3) The loss function  $\ell$  meets  $\nu_{\ell}$ -smooth with respect to  $\hat{y}$ :

$$\left| \frac{\partial \ell}{\partial \hat{y}}(\hat{y}, y) - \frac{\partial \ell}{\partial \hat{y}}(\hat{y}', y) \right| \le \nu_{\ell} |\hat{y} - \hat{y}'|, \ \forall \ \hat{y}, \hat{y}', y \in [y_{\min}, y_{\max}].$$

With these assumptions,  $|\frac{\partial \ell}{\partial \hat{y}}(\hat{y},y)| \leq \alpha_{\ell}$ , and  $\ell$  is bounded, i.e.,  $0 \leq \ell(\hat{y},y) \leq M$ .

**Assumption 3.** The learned parameters  $\{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(K)}, \mathbf{w}\}$  during the training procedure with limited iterations satisfies

$$\max \left\{ \|\mathbf{W}^{(1)}\|_{2}, \dots, \|\mathbf{W}^{(K)}\|_{2}, \|\mathbf{w}\|_{2} \right\} \le B.$$

## 4.2 Generalization Gap

This section presents the main results of this paper. Under the assumptions made in Section 4.1, the bound on the generalization gap of deep GCNs is provided in the following theorem.

**Theorem 1** (Generalization gap for deep GCNs). Consider the deep GCN model, defined in equation (1), which comprises K hidden layers and utilizes  $g(\mathbf{L})$  as the graph filter operator. The model is trained on  $\mathcal S$  using SGD for T iterations. Under Assumptions 1, 2 and 3 stated in Section 4.1, the following expected generalization gap is valid with a probability of at least  $1 - \delta$ , where  $\delta \in (0, 1)$ :

$$\epsilon_{gen}(\mathcal{A}_{\mathcal{S}}) \leq \frac{1}{\sqrt{m}} \left\{ O\left( \left( (K+1)\eta \kappa_1 + \eta \kappa_2 \right)^T \right) + M \sqrt{\frac{\log \frac{1}{\delta}}{2}} \right\},$$
(6)

where

$$\kappa_1 := (v_\ell \alpha_\sigma^2 + \alpha_\ell \nu_\sigma) (B\alpha_\sigma C_g)^{2K} C_g^2 C_{\mathbf{X}}^2 + \alpha_\ell (B\alpha_\sigma C_g)^{K-1} \alpha_\sigma^2 C_g^2 C_{\mathbf{X}}, \tag{7}$$

and

$$\kappa_2 := \nu_\sigma \left( B \alpha_\sigma C_g \right)^K C_g^2 C_\mathbf{X}^2 \left( \sum_{j=0}^{K-1} (j+1) (B \alpha_\sigma C_g)^j \right). \tag{8}$$

A fundamental correlation between the generalization gap and the parameters governing deep GCNs is induced by Theorem 1. This correlation implies that the uniform stability of deep GCNs, trained using the SGD algorithm, exhibits an increase with the number of samples when the upper bound approaches zero as the sample size m tends to infinity. Specifically, it is observed that if the value of  $C_g$  (presenting the largest absolute eigenvalue of a symmetric  $g(\mathbf{L})$  or the maximum singular value of an asymmetric  $g(\mathbf{L})$  remains unaffected by the size N, a generalization gap decaying at the order of  $O(1/\sqrt{m})$  is obtained. To compare with the result in [36], let us discuss at length the role of  $g(\mathbf{L})$  and the hidden layer number *K* on the generalization gap.

According to (7) and (8),  $\kappa_1 = O\left(C_g^{2K+2}\right)$  and  $\kappa_2 = O\left(C_g^{2K+1}\right)$ . Therefore, the bound on the generalization gap of deep GCNs in Theorem 1 is

$$\epsilon_{gen}(\mathcal{A}_{\mathcal{S}}) \le \frac{1}{\sqrt{m}} \left( O\left(C_g^{2T(K+1)}\right) + M\sqrt{\frac{\log \frac{1}{\delta}}{2}} \right).$$
(9)

When K = 0, the GCN model (1) degenerates into the single-layer GCN model considered in [36]. At this point, according to (9), we have

$$\epsilon_{gen}(\mathcal{A}_{\mathcal{S}}) \le \frac{1}{\sqrt{m}} \left( O\left(C_g^{2T}\right) + M\sqrt{\frac{\log \frac{1}{\delta}}{2}} \right),$$
(10)

which is the same as the result of [36].

**Remarks.** Based on (9), we present certain observations regarding the impact of filter  $g(\mathbf{L})$  and the hidden layer number K on the generalization capacity of deep GCNs in (1).

- **Normalized vs. Unnormalized Graph Filters:** We examine the three most commonly utilized filters: 1)  $g_1(\mathbf{L}) =$  $\mathbf{A} + \mathbf{I}$ , 2)  $g_2(\mathbf{L}) = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} + \mathbf{I}$ , and 3)  $g_3(\mathbf{L}) = \mathbf{D}^{-1} \mathbf{A} + \mathbf{I}$ . For the unnormalized filter  $g_1$ , its maximum absolute eigenvalue is bounded by O(N). Consequently, as the value of m approaches the magnitude to N, the upper bound indicated by (9) tends towards  $O(N^p)$  for some p > 0, leading to an impractical upper bound when N become infinitely large. On the contrary, for two normalized filters  $g_2$  and  $g_3$ , their largest absolute eigenvalues are bounded and independent of graph size N. Therefore, both filters yield a diminishing generalization gap at a rate of  $O(\frac{1}{\sqrt{m}})$  as m goes to infinity. This discovery underscores the superior performance of normalized filters over unnormalized counterparts in deep GCNs. This observation is consistent with the findings in [36], [37].
- **Low-pass vs. High-pass Graph Filters:** Our theoretical results are not restricted to the choice of  $q(\mathbf{L})$  as either a low-pass or a high-pass filter. To illustrate, consider two exponential filters with symmetric L: i) a low-pass filter  $g_{\text{low}}(\lambda) = e^{-b\lambda^2}$  and ii) a high-pass filter  $g_{\text{high}}(\lambda) = 1 - e^{-a\lambda^2}$ , where a, b > 0. In this setting, it is straightforward to verify that

$$||g_{\text{high}}(\mathbf{L})||_2 < ||g_{\text{low}}(\mathbf{L})||_2 = 1.$$

Consequently, both filters lead to a vanishing generalization gap at the rate of  $O\left(\frac{1}{\sqrt{m}}\right)$  as  $m \to \infty$ . The Role of Parameter K: It is evident that, when the values of  $C_g$  and T are fixed, the upper bound (9) exhibits an exponential dependence on parameter K. This observation implies that a larger value K leads to an increase in the upper bound of the generalization gap, thereby offering valuable insights for the architectural design of deep GCNs. This finding diverges from the ones presented in [36], [37], as these studies do not account for generic deep GCNs and overlook the significance of the parameter *K*.

Furthermore, based on Theorem 1, we give a brief analysis of the impact of  $d_k$  (width of the k-th layer) on the generalization. Actually, the impact of  $d_k$  on the generalization is reflected in its impact on B. More specifically, let us consider the case where parameters  $\{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(K)}, \mathbf{w}\}$  belong to the set  $\mathcal{X}_{\xi}$ , where

$$\mathcal{X}_{\xi} := \{ \mathbf{W} : \| \mathbf{W} \|_{\infty} \le \xi \},$$

i.e.,  $\mathcal{X}_{\xi}$  is the collection of all matrices whose elements' absolute values are all less than  $\xi$ . At this point, for  $\mathbf{W}^{(k)} \in \mathbb{R}^{d_{k-1} \times d_k}$ , we have

$$\sup_{\mathbf{W}^{(k)} \in \mathcal{X}_{\varepsilon}} \|\mathbf{W}^{(k)}\|_{2} \leq \sup_{\mathbf{W}^{(k)} \in \mathcal{X}_{\varepsilon}} \|\mathbf{W}^{(k)}\|_{F} \leq \xi \sqrt{d_{k-1}d_{k}}.$$

Therefore, a larger  $d_k$  (i.e., width of the k-th layer) results in a larger upper bound of  $\|\mathbf{W}^{(k)}\|_2$ , which implies that a larger  $d_k$  results in a larger B (see Assumption 3 in Section 4.1). Finally, Theorem 1 indicates that a larger B leads to a larger bound on the generalization gap, thus we conclude that a larger  $d_k$  leads to a larger bound on the generalization gap. To justify this argument, we add some experimental studies in Section 5. The empirical results are consistent with our analysis.

TABLE 2	
Comparison of the generalization gap estimated based on uniform stability	/.

Reference	Model Architecture	Estimated Upper Bound of the Generalization Gap
[36]	shallow	$\frac{1}{\sqrt{m}} \left( O\left( (1 + \eta \upsilon_{\ell} \upsilon_{\sigma} C_g^2)^T \right) + M \sqrt{\frac{\log \frac{1}{\delta}}{2}} \right)$
[37]	shallow	$\frac{1}{\sqrt{m}} \left( O\left( \eta \alpha_{\ell} \alpha_{\sigma} c_{2,T} \sum_{t=0}^{T-1} c_{6,t} \prod_{s=t+1}^{T-1} (1 + \eta c_{5,s}) \right) + M \sqrt{\frac{\log \frac{1}{\delta}}{2}} \right)$
[44]	shallow	$\frac{1}{\sqrt{m}} \left\{ O\left(C_g^2 \eta C_{p,\lambda} \sum_{t=1}^T (C_{p,\lambda} (1 + (\alpha_\sigma^2 + \alpha_\ell) \eta C_g^2))^{t-1}\right) + M\sqrt{\frac{\log \frac{1}{\delta}}{2}} \right\}$
Ours	deep	$\frac{1}{\sqrt{m}} \left\{ O\left( \left( (K+1)\eta \kappa_1 + \eta \kappa_2 \right)^T \right) + M\sqrt{\frac{\log \frac{1}{\delta}}{2}} \right\}$

Note: m is the number of samples in the trained dataset; M is the upper bound of loss function  $\ell(\cdot,\cdot)$ ;  $\eta>0$  is the learning rate;  $\delta\in(0,1)$ ; T is the number of iterations for training  $A_{\mathcal{S}}$  using SGD;  $C_g$  represents the 2-norm of filter  $g(\mathbf{L})$ ;  $\alpha_{\sigma}$ ,  $v_{\sigma}$  are two parameters w.r.t the continuity of activation function  $\sigma(\cdot)$ ;  $\alpha_{\ell}$ ,  $v_{\ell}$  are two parameters w.r.t the continuity of the loss function  $\ell(\cdot,\cdot)$ .  $c_{2,t}$ ,  $c_{6,t}$ ,  $c_{5,t} > 0$  ( $t=0,1,\ldots,T$ ) represent some specific parameters defined in [37].  $C_{p,\lambda} = \frac{28}{p(p-1)\lambda_t} (B_{\ell}/\lambda)^{(3-p)/p}$ , where  $B_{\ell} > 0$  is a parameter related to loss function  $\ell(\cdot,\cdot)$ ,  $1 , <math>\lambda > 0$  is the regularization parameter and  $\lambda_t > 0$  is another regularization parameter dependent on  $\lambda$  and t, as detailed in [44]. K is number of hidden layers of the considered deep GCNs;  $\kappa_1$  and  $\kappa_2$  are two parameters as defined in (7) and (8).

Table 2 offers a concise summary of various upper bounds on the generalization gap, derived through the application of uniform stability. From Table 2, we can see that all the works derive a generalization gap decaying at the order of  $O(1/\sqrt{m})$ . However, compared to the other three works which only consider shallow GCNs, our work explores the case of deep GCNs. We should point out that the generalization of single-layer GCNs into deep GCNs is not trivial. To derive the results for deep GCNs, we tackle two significant challenges that arise specifically in the context of deep GCNs, which are unique to deep GCNs and are non-existent in single-layer models. The first challenge is the derivation of the gradient of the final output with respect to the learnable parameters across multiple layers, which requires determining how the gradient of the overall error of a GCN is shared among neurons in different hidden layers. In particular, in Appendix A, we provide a recursive formula to compute the related gradients. The second challenge is the evaluation of gradient variations between GCNs trained on different datasets. In the single layer case, since the input feature is the same, the variation of the related gradient is only dependent on the variations of learnable parameters. While, in the case of deep GCNs, the variation of the related gradients is also dependent on the variations of the gradients of the final output with respect to the hidden layer outputs. Please see Lemma 7 and its proof for details (see Appendix C).

# 4.3 Stability Upper Bound

In this subsection, we establish the uniform stability of SGD for deep GCNs, which is the key to further proving Theorem 1.

**Theorem 2 (Uniform stability of deep GCNs).** Consider the deep GCNs defined by equation (1), which are trained on a dataset S using the SGD algorithm for a total of T iterations and denoted as  $A_S$ . Assume that Assumptions 1, 2 and 3 stated in Section 4.1 are satisfied. Then,  $A_S$  is  $\mu_m$ -uniformly stable, with  $\mu_m$  satisfying the following condition:

$$\mu_m \le \frac{C}{m} \sum_{t=1}^{T} \left( 1 + (K+1)\eta \kappa_1 + \eta \kappa_2 \right)^{t-1},$$
(11)

where

$$C:=(K+1)\eta\alpha_\ell^2(B\alpha_\sigma C_g)^{2K}\alpha_\sigma^2C_g^2C_{\mathbf{X}}^2,$$

 $\kappa_1$  and  $\kappa_2$  are defined by (7) and (8), respectively.

With a straightforward calculation, one can see that

$$\mu_m \le \frac{1}{m} O\bigg( \Big( (K+1)\eta \kappa_1 + \eta \kappa_2 \Big)^T \bigg),$$

which decays at the rate of  $\frac{1}{m}$  as m tends to infinity. Together with Lemma 1, it yields the result of Theorem 1. **Proof Sketch for Theorem 2**. We prove Theorem 2 in the following two steps.

- **Step 1:** We begin by bounding the stability of deep GCNs with respect to perturbations in the learned parameters caused by changes in the training set. The result is given in Lemma 2.
- Step 2: Next, we provide a bound for the perturbation of the learned parameters. The result is presented in Theorem

Consider  $\mathcal{A}_{\mathcal{S}}$ , a set of deepGCNs defined by (1), trained on the dataset  $\mathcal{S}$  using SGD for T iterations. Let  $\theta_t = \{\mathbf{W}_t^{(1)}, \dots, \mathbf{W}_t^{(K)}, \mathbf{w}_t\}$  and  $\theta_t' = \{\mathbf{W}_t^{(1)'}, \dots, \mathbf{W}_t^{(K)'}, \mathbf{w}_t'\}$  (with  $\theta_0 = \theta_0'$ ) denote the parameters of two GCNs trained on  $\mathcal{S}$  and  $\mathcal{S}^i$  after t iterations, respectively. We set  $\Delta \mathbf{w}_t = \mathbf{w}_t - \mathbf{w}_t'$  and  $\Delta \mathbf{W}_t^{(k)} = \mathbf{W}_t^{(k)} - \mathbf{W}_t^{(k)'}$  to be the perturbation of learning parameters and define

$$\|\triangle\theta_t\|_* = \|\triangle\mathbf{w}_t\|_2 + \sum_{k=1}^K \|\triangle\mathbf{W}_t^{(k)}\|_2.$$
 (12)

In the following lemma, it is shown that the stability of  $A_S$  can be bounded by  $\|\triangle\theta_T\|_*$ .

**Lemma 2.** Let  $\theta_t$  and  $\theta_t'$  be the learnt parameters of two GCNs trained on S and  $S^i$  using SGD in the t-th iteration with  $\theta_0 = \theta_0'$ , and  $\Delta \theta_t := \theta_t - \theta_t'$ . Suppose that all the assumptions made in Section 4.1 hold. Then, after T iterations, we have that for any  $\mathbf{z} = (\mathbf{x}, y)$  taken from D,

$$\left| \mathbb{E}_{\mathcal{A}} \left[ \ell(\hat{y}, y) \right] - \mathbb{E}_{\mathcal{A}} \left[ \ell(\hat{y}', y) \right] \right| \le \alpha_{\ell} B^{K} \alpha_{\sigma}^{K+1} C_{g}^{K+1} C_{\mathbf{X}} \cdot \mathbb{E}_{\mathcal{A}} \left[ \| \triangle \theta_{T} \|_{*} \right], \tag{13}$$

where  $\hat{y} = f(\mathbf{x}|\theta_T)$  and  $\hat{y}' = f(\mathbf{x}|\theta_T')$ .

We provide the proof of Lemma 2 in Appendix B.

Combining (5) and (13), the stability of  $A_S$  has a bound

$$\mu_m \le \frac{\alpha_\ell B^K \alpha_\sigma^{K+1} C_g^{K+1} C_{\mathbf{X}}}{2} \sup_{\mathcal{S}} \left\{ \mathbb{E}_{\mathcal{A}} \left[ \| \triangle \theta_T \|_* \right] \right\}. \tag{14}$$

So, to estimate the uniform stability of  $\mathcal{A}_{\mathcal{S}}$ , we need to bound  $\mathbb{E}_{\mathcal{A}}[\|\triangle\theta_T\|_*]$ . Now, let us recall (3) for parameter updating, for training on  $\mathcal{S}$ ,

$$\mathbf{w}_{t} = \mathbf{w}_{t-1} - \eta \nabla_{\mathbf{w}} \ell(f(\mathbf{x}_{t} | \theta_{t-1}), y_{t}),$$

$$\mathbf{W}_{t}^{(k)} = \mathbf{W}_{t-1}^{(k)} - \eta \nabla_{\mathbf{W}^{(k)}} \ell(f(\mathbf{x}_{t} | \theta_{t-1}), y_{t}),$$

 $k = 1, 2, \dots, K$ , and for training on  $S^i$ 

$$\begin{aligned} \mathbf{w}_t' &= \mathbf{w}_{t-1}' - \eta \nabla_{\mathbf{w}} \ell(f(\mathbf{x}_t'|\theta_{t-1}'), y_t'), \\ \mathbf{W}_t^{(k)'} &= \mathbf{W}_{t-1}^{(k)'} - \eta \nabla_{\mathbf{W}^{(k)}} \ell(f(\mathbf{x}_t'|\theta_{t-1}'), y_t'), \end{aligned}$$

k = 1, 2, ..., K, where  $(\mathbf{x}_t, y_t) \in \mathcal{S}$  and  $(\mathbf{x}_t', y_t') \in \mathcal{S}^i$  are the samples drawn at the t-th SGD iteration. Therefore,  $\triangle \theta_t = \{\triangle \mathbf{W}_t^{(1)}, \ldots, \triangle \mathbf{W}_t^{(K)}, \triangle \mathbf{w}_t\}$  has the following iterations:

$$\Delta \mathbf{w}_t = \Delta \mathbf{w}_{t-1} - \eta \Big( \nabla_{\mathbf{w}} \ell(f(\mathbf{x}_t | \theta_{t-1}), y_t) - \nabla_{\mathbf{w}} \ell(f(\mathbf{x}_t' | \theta_{t-1}'), y_t') \Big),$$

and for k = 1, 2, ..., K,

$$\Delta \mathbf{W}_{t}^{(k)} = \Delta \mathbf{W}_{t-1}^{(k)} - \eta \Big( \nabla_{\mathbf{W}^{(k)}} \ell(f(\mathbf{x}_{t}|\theta_{t-1}), y_{t}) - \nabla_{\mathbf{W}^{(k)}} \ell(f(\mathbf{x}_{t}'|\theta_{t-1}'), y_{t}') \Big),$$

with  $\|\triangle \theta_0\|_* = 0$ .

So, we need to bound

$$\nabla_{\mathbf{w}} \ell(f(\mathbf{x}_t | \theta_{t-1}), y_t) - \nabla_{\mathbf{w}} \ell(f(\mathbf{x}_t' | \theta_{t-1}'), y_t')$$

and

$$\nabla_{\mathbf{W}^{(k)}} \ell(f(\mathbf{x}_t | \theta_{t-1}), y_t) - \nabla_{\mathbf{W}^{(k)}} \ell(f(\mathbf{x}_t' | \theta_{t-1}'), y_t')$$

to obtain a bound of  $\|\triangle\theta_t\|_*$ . There are two scenarios to consider: i) At step t, SGD picks a sample  $\mathbf{z}_t = (\mathbf{x}_t, \mathbf{y}_t)$  which is identical in  $\mathcal{S}$  and  $\mathcal{S}^i$ , and occurs with probability (m-1)/m; and ii) At step t, SGD picks the only samples that  $\mathcal{S}$  and  $\mathcal{S}^i$  differ,  $\mathbf{z}_t = (\mathbf{x}_t, \mathbf{y}_t)$  and  $\mathbf{z}_t' = (\mathbf{x}_t', \mathbf{y}_t')$  which occurs with probability 1/m. We provide the results in the following Lemma 3 and Lemma 4.

**Lemma 3.** Consider two GCNs with parameters  $\theta_t$  and  $\theta'_t$ , respectively. Then, the following holds for any sample  $\mathbf{z}_t = (\mathbf{x}_t, y_t)$ :

$$\|\nabla_{\mathbf{w}}\ell(f(\mathbf{x}_t|\theta_{t-1}), y_t) - \nabla_{\mathbf{w}}\ell(f(\mathbf{x}_t|\theta'_{t-1}), y_t)\|_F \le \kappa_1 \|\triangle\theta_{t-1}\|_*, \tag{15}$$

and for k = 1, 2, ..., K,

$$\|\nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_{t}|\theta_{t-1}), y_{t}) - \nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_{t}|\theta'_{t-1}), y_{t})\|_{F} \le (\kappa_{1} + \rho_{k})\|\Delta\theta_{t-1}\|_{*}, \tag{16}$$

where  $\kappa_1$  and  $\rho_k$  are defined by (7) and (A.12).

**Lemma 4.** Consider two GCNs with parameters  $\theta_t$  and  $\theta_t'$ , respectively. Then, the following holds for any two samples  $\mathbf{z}_t = (\mathbf{x}_t, y_t)$  and  $\mathbf{z}_t' = (\mathbf{x}_t', y_t')$ :

$$\|\nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_t|\theta_{t-1}), y_t) - \nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_t'|\theta_{t-1}', y_t'))\|_F \le 2\alpha_\ell B^K \alpha_\sigma^{K+1} C_g^{K+1} C_{\mathbf{X}}, \tag{17}$$

for k = 1, 2, ..., K + 1. Note that  $\mathbf{W}^{(K+1)} = \mathbf{w}$ .

The proofs of Lemma 3 and Lemma 4 are given in Appendix C. We now provide a bound for  $\mathbb{E}_{\mathcal{A}}[\|\Delta\theta_T\|_*]$ .

**Theorem 3.** Let  $\theta_t$  and  $\theta_t'$  be the learnt parameters of two GCNs trained on  $\mathcal{S}$  and  $\mathcal{S}^i$  using SGD in the t-th iteration with  $\theta_0 = \theta_0'$ . The assumptions made in Section 4.1 hold. Then, after T iterations,  $\Delta \theta_T$  satisfies

$$\mathbb{E}_{\mathcal{A}}\left[\|\triangle\theta_T\|_*\right] \le c \sum_{t=1}^T \left(1 + (K+1)\eta\kappa_1 + \eta\kappa_2\right)^{t-1},\tag{18}$$

where  $c:=\frac{2(K+1)\eta\alpha_{\ell}B^{K}\alpha_{\sigma}^{K+1}C_{g}^{K+1}C_{\mathbf{X}}}{m}$ , and  $\kappa_{1}$  and  $\kappa_{2}$  are defined by (7) and (8), respectively.

The proof of Theorem 3, using Lemma 3 and Lemma 4, is provided in Appendix D. Combining (14) and Theorem 3, we obtain that the uniform stability  $\mu_m$  of  $A_S$  has a bound as

$$\mu_m \le \alpha_{\ell} B^K \alpha_{\sigma}^{K+1} C_g^{K+1} C_{\mathbf{X}} \sup_{\mathcal{S}} \left\{ \mathbb{E}_A \left[ \| \triangle \theta_T \|_* \right] \right\}$$
$$\le \frac{C}{m} \sum_{t=1}^{T} \left( 1 + (K+1) \eta \kappa_1 + \eta \kappa_2 \right)^{t-1},$$

which completes the proof of Theorem 2.

# **5 EXPERIMENTS**

In this section, we conduct some empirical studies using three benchmark datasets commonly utilized for the node classification task, namely Cora, Citeseer, and Pubmed [58], [59]. Table 3 summarizes the basic statistics of these datasets.

TABLE 3

Statistics of the three benchmark datasets.

	Cora	Citeseer	Pubmed
# Nodes	2,708	3,327	19,717
# Edges	5,429	4,732	44,338
# Features	1,433	3,703	500
# Classes	7	6	3
Label Rate	0.052	0.036	0.003

In our experiments, we follow the standard transductive learning problem formulation and the training/test setting used in [60]. To rigorously test our theoretical insights, our experiments aim to answer the following key questions:

- Q1: How does the design of graph filters (i.e., g(L)) influence the generalization gap?
- Q2: How does the generalization gap change with the number of hidden layers (i.e., K)?
- Q3: How does the width (i.e., the number of hidden units: d) affect the generalization gap?

To address each question, we empirically estimate the generalization gap by calculating the absolute difference in loss between training and test samples. We adopt the official TensorFlow implementation (https://github.com/tkipf/gcn) for GCN [60] and the Adam optimizer with default settings. The number of iterations is fixed to T=200 for all the simulations. Results and Discussion for Q1. We analyze two types of graph filters in our study: 1) the normalized graph filter, defined as  $g(\mathbf{L}) = \tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1/2}$  with  $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$  and  $\tilde{\mathbf{D}}_{ii} = \sum_j \tilde{\mathbf{A}}_{ij}$  (which was first employed in the vanilla GCN [60] and has subsequently become widely used in follow-up works on GCNs), and 2) the random walk filter,  $g(\mathbf{L}) = \mathbf{D}^{-1}\mathbf{A} + \mathbf{I}$ . To fit our theoretical finding, we compare the performance of two 5-layer GCN models (with width d=32 for each layer), each employing one of these filters. Table 4 presents the numerical records of  $R_{emp}(\mathcal{A}_{\mathcal{S}})$ ,  $R(\mathcal{A}_{\mathcal{S}})$ ,  $\epsilon_{gen}(\mathcal{A}_{\mathcal{S}})$ ,  $C_g$  for both filters. The results indicate clearly that the 5-layer GCN with the normalized graph filter exhibits a smaller generalization gap compared to the one with the random walk filter. Furthermore, Fig. 1 illustrates the performance of each filter across different datasets over iterations, demonstrating the superior performance of the normalized graph filter. Overall, the empirical findings in Table 4 and Fig. 1 align well with our theoretical finding regarding the impact of  $C_g$  on the generalization gap.

Results and Discussion for Q2. In this experimental study, we try different settings of K, i.e., the number of hidden layers. Specifically, for  $K = \{1, 2, 3, 4, 5\}$ , we compare the performance of two K-layer GCNs (with width d = 32 for each layer): one employing the normalized graph filter  $g(\mathbf{L}) = \tilde{\mathbf{D}}^{-1/2} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-1/2}$ , and one using the random walk filter  $g(\mathbf{L}) = \mathbf{D}^{-1} \mathbf{A} + \mathbf{I}$ . Fig. 2 shows the performance comparison results for each K. It demonstrates clearly that, consistent with the aforementioned results for Q1, GCN with a normalized graph filter (with smaller  $C_g$ ) consistently exhibits

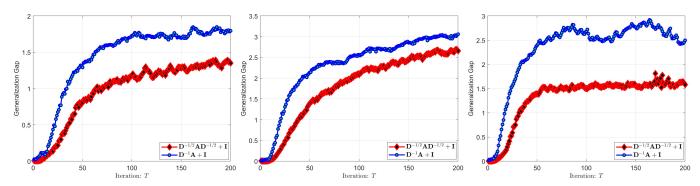


Fig. 1. Comparison of trends in the generalization gap: Cora (left), Citeseer (middle), Pubmed (right).

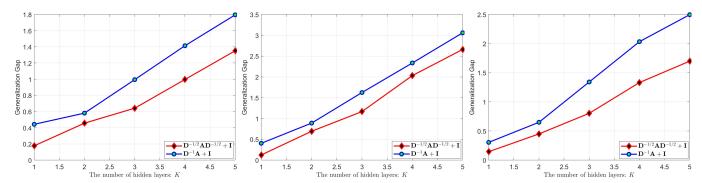


Fig. 2. Comparison of the generalization gap with different settings of network depth K: Cora (left), Citeseer (middle), Pubmed (right).

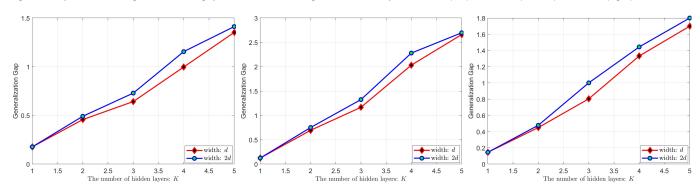


Fig. 3. Comparison of the generalization gap with different settings of network width d: Cora (left), Citeseer (middle), Pubmed (right).

TABLE 4
The generalization gap with different graph filter for three datasets.

Dataset	Graph filter $g(\mathbf{L})$	$R_{emp}(\mathcal{A}_{\mathcal{S}})$	$R(\mathcal{A}_{\mathcal{S}})$	$\epsilon_{gen}(\mathcal{A_S})$	$C_g$
Cora	$\mathbf{ ilde{D}}^{-1/2}\mathbf{ ilde{A}}\mathbf{ ilde{D}}^{-1/2}$ $\mathbf{D}^{-1}\mathbf{A}+\mathbf{I}$	1.488 1.914	0.136 0.118	<b>1.352</b> 1.796	1 4.746
Citeseer	$\mathbf{ ilde{D}}^{-1/2}\mathbf{ ilde{A}}\mathbf{ ilde{D}}^{-1/2} \\ \mathbf{D}^{-1}\mathbf{A} + \mathbf{I}$	2.896 3.206	0.235 0.145	<b>2.661</b> 3.061	1 4.690
Pubmed	$\mathbf{ ilde{D}}^{-1/2}\mathbf{ ilde{A}}\mathbf{ ilde{D}}^{-1/2}$ $\mathbf{D}^{-1}\mathbf{A}+\mathbf{I}$	1.594 2.534	0.023 0.037	<b>1.571</b> 2.497	1 7.131

smaller generalization gaps compared to those with the random walk filter. Also, it is observed that the generalization gap becomes larger as K increases, further validating our theoretical assertions regarding the influence of K on the model's generalization gap.

**Results and Discussion for Q3.** To empirically investigate the impact of width d (i.e., the number of hidden units) on the generalization gap, we conduct additional experiments using a 5-layer GCN equipped with a normalized graph filter. The experiments specifically involve a comparison between a 5-layer GCN configured with a width of 2d for each layer and the previously studied model with d width (d=32), as illustrated in Fig. 3. This setup allows for a direct comparison under varying network configurations, providing insights into how changes in the number of hidden units influence the generalization gap. As demonstrated in Fig. 3, across all the datasets examined, a d-width GCN consistently exhibits smaller generalization gaps compared to one with a 2d-width. This observation is in harmony with our theoretical explanation

presented after Theorem 1, that is, the factor B (i.e., the upper bound of 2-norm of the parameters  $\{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(K)}, \mathbf{w}\}$ ) directly influences factors  $\kappa_1$  and  $\kappa_2$  in the upper bound of the generalization gap.

## 6 THEORETICAL IMPLICATIONS

Our work establishes a theoretical framework for analyzing the generalization gap of traditional deep GCNs, which further provides insights into extending the analysis to other classes of graph neural networks, including Graph Transformers. As illustrative examples, we briefly discuss how the theoretical proof methodology developed in our framework can be applied to GCNII and Graph Transformer, which are representative models of more advanced GNNs, thereby demonstrating the broader applicability of our theoretical framework.

#### 6.1 Extension to GCNII

With input features  $\mathbf{X}^{(0)} = \mathbf{X} \in \mathbb{R}^{N \times d}$ , GCNII defines its k-th layer as

$$\mathbf{X}^{(k)} = \sigma \left( \left( (1 - a_k) g(\mathbf{L}) \mathbf{X}^{(k-1)} + a_k \mathbf{X}^{(0)} \right) \cdot \left( (1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)} \right) \right),$$

for  $k=1,2,\ldots,K$ , where  $a_k,b_k\in(0,1)$  are two hyperparameters,  $\mathbf{X}^{(k)}$  is the output feature matrix of the k-th layer,  $\mathbf{W}^{(k)}$  is the trained parameter matrix specific to the k-th layer, graph filter  $g(\mathbf{L})=\tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1/2}$ , and  $\mathbf{I}_d$  is the  $d\times d$  identity matrix. The output for node  $\mathbf{x}$  is

$$f(\mathbf{x}|\theta) = \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} \left( (1 - a_{K+1}) g(\mathbf{L}) \mathbf{X}^{(K)} + a_{K+1} \mathbf{X}^{(0)} \right) \mathbf{w} \right),$$

where  $\theta = \{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \dots, \mathbf{W}^{(K)}, \mathbf{w}\}$  (all trainable parameters, with  $\mathbf{w} \in \mathbb{R}^d$  the output layer parameter);  $\delta_{\mathbf{x}} \in \mathbb{R}^N$  is the indicator vector for node  $\mathbf{x}$ ;  $a_{K+1} \in (0,1)$  is a hyperparameter for the output layer residual connection. Let  $\theta_t$  and  $\theta_t'$  be the learned parameters of two GCNs trained on  $\mathcal{S}$  and  $\mathcal{S}^i$  using SGD in the t-th iteration with  $\theta_0 = \theta_0'$ , and  $\Delta \theta_t := \theta_t - \theta_t'$ . For each layer k, the perturbation of layer outputs  $\|\Delta \mathbf{X}^{(k)}\|_F = \|\mathbf{X}^{(k)} - \mathbf{X}^{(k)'}\|_F$  satisfies the recursive bound:

$$\|\triangle \mathbf{X}^{(k)}\|_{F} \le c_{1}^{(k)} \|\triangle \mathbf{X}^{(k-1)}\|_{F} + c_{2}^{(k)} \|\triangle \mathbf{W}^{(k)}\|_{2},\tag{19}$$

where  $c_1^{(k)} = (1-a_k)(1-b_k+b_kB)\alpha_\sigma C_g$  and  $c_2^{(k)} = \alpha_\sigma b_k \left((1-a_k)C_gB_{\mathbf{X}}^{(k-1)} + a_kC_{\mathbf{X}}\right)$  with  $B_{\mathbf{X}}^{(k-1)}$  the bound of  $\|\mathbf{X}^{(k-1)}\|_F$  (see (A.22) in the Appendix E). The first term on the right side of the iterative formula captures propagation of perturbations from the previous layer, while the second term captures perturbation from  $\mathbf{W}^{(k)}$ .

By induction, it yields that

$$\|\triangle \mathbf{X}^{(k)}\|_F \le e^{(k)} (\sum_{j=1}^k \|\triangle \mathbf{W}^{(k)}\|_2),$$
 (20)

where  $e^{(k)} = \max\{c_1^{(k)}e^{(k-1)},c_2^{(k)}\}$  with  $e^{(0)} = 0$ . We provide the proof of (19) and (20) in Appendix E. Then, combining layer-wise bounds and using the Lipschitz property of  $\sigma$ , one can have the output perturbation  $|f(x|\theta) - f(x|\theta')|$  bounded by the total parameter perturbation  $\|\Delta\theta\|_* = \sum\limits_{j=1}^K \|\mathbf{W}^{(j)} - \mathbf{W}^{(j)'}\|_2 + \|\mathbf{w} - \mathbf{w}'\|_2$  (see Appendix E for technical details) as

$$|f(\mathbf{x}|\theta) - f(\mathbf{x}|\theta')| \le \alpha_{\sigma} \cdot \varrho ||\Delta \theta||_{*}, \tag{21}$$

where  $\varrho = \max \left\{ (1 - a_{K+1}) B C_g \cdot e^{(K)}, (1 - a_{K+1}) C_g B_{\mathbf{X}}^{(K)} + a_{K+1} C_{\mathbf{X}} \right\}$ . Then,

$$\left| \mathbb{E}_{\mathcal{A}} [\ell(\hat{y}, y)] - \mathbb{E}_{\mathcal{A}} [\ell(\hat{y}', y)] \right| = \left| \mathbb{E}_{\mathcal{A}} [\ell(f(\mathbf{x}|\theta_T), y) - \ell(f(\mathbf{x}|\theta_T'), y)] \right| \le \alpha_{\ell} \mathbb{E}_{\mathcal{A}} [|f(\mathbf{x}|\theta_T) - f(\mathbf{x}|\theta_T')|] \le \varrho \alpha_{\ell} \cdot \mathbb{E}_{\mathcal{A}} [||\Delta \theta_T||_*].$$

This implies that the stability of  $A_S$  for GCNII has a bound

$$\mu_m \leq \frac{\varrho \alpha_\ell}{2} \sup_{S} \Big\{ \mathbb{E}_{\mathcal{A}}[\|\triangle \theta_T\|_*] \Big\}.$$

Note that when  $a_k = 0, b_k = 1$  for all k, GCNII degenerates into the traditional GCN, we have  $\varrho = B^K \alpha_\sigma^K C_g^{K+1} C_{\mathbf{X}}$ , and thus

$$\mu_m \le \frac{\alpha_{\ell} B^K \alpha_{\sigma}^K C_g^{K+1} C_{\mathbf{X}}}{2} \sup_{\mathcal{S}} \left\{ \mathbb{E}_{\mathcal{A}}[\| \triangle \theta_T \|_*] \right\},\,$$

which is consistent with (14).

To further bound  $\|\triangle\theta_T\|_*$ , the crucial step is to bound the perturbation of the gradient of  $f(\mathbf{x}|\theta)$  with respect to the parameters  $\theta = \{\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_K, \mathbf{w}\}$  and obtain the result similar to Lemma 7 in Appendix C, which can be achieved by following the technique in our paper. Here, we provide the result for  $\|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta')\|_F$ :

$$\|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta')\|_F \le \left(\nu_{\sigma} \varrho \cdot \left((1 - a_{K+1})C_g B_{\mathbf{X}}^{(K)} + a_{K+1}C_{\mathbf{X}}\right) + \alpha_{\sigma} \cdot (1 - a_{K+1})C_g e^{(K)}\right) \cdot \|\triangle\theta\|_*, \tag{22}$$

where  $\varrho = \max\left\{(1-a_{K+1})BC_g\cdot e^{(K)}, (1-a_{K+1})C_gB_{\mathbf{X}}^{(K)} + a_{K+1}C_{\mathbf{X}}\right\}$ . Note that when  $a_k=0, b_k=1$  for all k, GCNII degenerates into the traditional GCN, we have  $\varrho=B^K\alpha_\sigma^KC_g^{K+1}C_{\mathbf{X}}, B_{\mathbf{X}}^{(K)}=B^K\alpha_\sigma^KC_g^KC_{\mathbf{X}}$  and  $e^{(K)}=B^{K-1}\alpha_\sigma^KC_g^KC_{\mathbf{X}}$ . At this point,

$$\left\| \nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta') \right\|_{F} \le \left( \upsilon_{\sigma} B^{2K} \alpha_{\sigma}^{2K} C_{g}^{2K+2} C_{\mathbf{X}}^{2} + B^{K-1} \alpha_{\sigma}^{K+1} C_{g}^{K+1} C_{\mathbf{X}} \right) \|\triangle \theta\|_{*},$$

which is consistent with (A.10) in Appendix C. For the bound of  $\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_F$ , we refer the readers to the proof process of (A.27) in Appendix E.

Finally, these structured analysis results can lead to the results corresponding Lemma 3 and Lemma 4, and thus enable bounding the stability of GCNII.

# 6.2 Extension to Graph Transformer

To extend our theoretical framework to more complex models like Graph Transformer, the key is to bound the generalization gap of Graph Transformer by quantifying how perturbations in the training set (e.g., removing or replacing a node) propagate to changes in model outputs. Graph Transformer introduce new learnable parameters: query  $(\mathbf{W}_Q)$ , key  $(\mathbf{W}_K)$ , and value  $(\mathbf{W}_V)$  projection matrices, alongside attention scalers and feed-forward layers, for which a self-attention layer is defined [43] as

$$F\left(\mathbf{x}_{n}\right) = \mathbf{a}^{\top} \operatorname{Relu}\left(\mathbf{W}_{O} \sum_{i \in \mathcal{T}^{n}} \mathbf{W}_{V} \mathbf{x}_{i} \cdot \operatorname{softmax}_{n}\left(\left(\mathbf{W}_{K} \mathbf{x}_{i}\right)^{\top} \mathbf{W}_{Q} \mathbf{x}_{n}\right)\right),$$

where  $\mathbf{x}_i$  denotes features of node i,  $\mathcal{T}^n$  is the set of nodes for the aggregation computing of node n, and softmax $_n(h(i,n)) = \exp(h(i,n))/\sum_{j\in\mathcal{T}^n}\exp(h(j,n))$ . Despite their architectural complexity (e.g., self-attention mechanisms, query/key/value projections), gradient decomposition still remains to be conducted via the product rule and chain rule, accounting for the propagation of attention-weight variations to the final output. Besides, a Lipschitz-type inequality for softmax may be critically needed, for which we claim that for  $\mathbf{z}=(z_1,z_2,\ldots,z_p)$ ,  $\mathbf{z}'=(z_1',z_2',\ldots,z_p')$  with  $\|\mathbf{z}-\mathbf{z}'\|_{\infty}\leq 1$ ,

$$\|\operatorname{softmax}(\mathbf{z}) - \operatorname{softmax}(\mathbf{z}')\|_1 \le 2e\|\mathbf{z} - \mathbf{z}'\|_{\infty}.$$
 (23)

Actually, the proof is not hard to set up by straight forward boundedness and the mean value theorem of exponential functions (see the technical details in Appendix F).

For trainable parameters  $\mathbf{W}_Q$ ,  $\mathbf{W}_K$ ,  $\mathbf{W}_V$ , set the attention output is:

$$F(\mathbf{x}_n) = \mathbf{a}^{\top} \text{ReLu} \Big( \mathbf{W}_O \sum_{i \in \mathcal{T}^n} \mathbf{W}_V \mathbf{x}_i \cdot \text{Attn}(\mathbf{x}_n)_i \Big),$$

where  $S_{i,n} = (\mathbf{W}_K \mathbf{x}_i)^T (\mathbf{W}_Q \mathbf{x}_n)$  is the scaled dot-product score,  $A_{i,n} = \operatorname{softmax}_n(S_{i,n})$  are attention weights, and  $\operatorname{Attn}(\mathbf{x}_n) = \sum_{i \in \mathcal{T}^n} \mathbf{W}_V \mathbf{x}_i \cdot A_{i,n}$  the attention output. Then the gradient decomposition with respect to  $\mathbf{W}_K$  is given by

$$\nabla_{\mathbf{W}_{K}}F(\mathbf{x}_{n}) = \underbrace{\nabla_{\mathrm{ReLU}(\mathbf{Z})}F(\mathbf{x}_{n})}_{\text{(I)}} \cdot \underbrace{\nabla_{\mathbf{Z}}\mathrm{ReLU}(\mathbf{Z})}_{\text{(2)}} \cdot \underbrace{\nabla_{\mathrm{Attn}(\mathbf{x}_{n})}\mathbf{Z}}_{\text{(3)}} \cdot \underbrace{\nabla_{\mathbf{A}}\mathrm{Attn}(\mathbf{x}_{n})}_{\text{(4)}} \cdot \underbrace{\nabla_{\mathbf{S}}\mathbf{A}}_{\text{(5)}} \cdot \underbrace{\nabla_{\mathbf{W}_{K}}S}_{\text{(6)}}$$

where  $\mathbf{Z} = \mathbf{W}_O \cdot \operatorname{Attn}(\mathbf{x}_n)$ ,  $\mathbf{A} = \{A_{i,n}\}$ , and  $\mathbf{S} = \{S_{i,n}\}$ . Then calculating each item gives that

$$\nabla_{\mathbf{W}_K} F(\mathbf{x}_n) = \mathbf{a}^\top \mathbb{I}_{\geq 0} (\mathbf{W}_O \cdot \operatorname{Attn}(\mathbf{x}_n)) \cdot \mathbf{W}_Q \cdot \mathbf{W}_V \cdot \left( \sum_{i \in \mathcal{T}^n} \mathbf{A}_{i,n} (\mathbf{x}_i - \bar{\mathbf{x}}_n) \mathbf{x}_i^\top \right) \cdot (\mathbf{W}_Q \mathbf{x}_n)^\top.$$

By leveraging the Lipschitz continuity of the gradient with respect to its trainable parameters, it can lead to bounding the gradient perturbation in terms of the total parameter perturbation  $\|\Delta\theta\|_* = \|\mathbf{W}_K - \mathbf{W}_K'\|_2 + \|\mathbf{W}_V - \mathbf{W}_V'\|_2 + \|\mathbf{W}_O - \mathbf{W}_O'\|_2 + \|\mathbf{W}_Q - \mathbf{W}_Q'\|_2 + \|\mathbf{a} - \mathbf{a}'\|_2$  by

$$\|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta')\|_2 \le 2eK_{\max} B^3 C_{\mathbf{X}}^3 \|\Delta \theta\|_*, \tag{24}$$

where  $K_{\max} \geq |\mathcal{T}^n|$  is the maximum neighborhood size, B is the upper bound of weight matrices (technical details in Appendix F). It mirrors the Lemma 7 in our approach for deep GCNs, where we recursively decomposed gradients across layers (see Lemma 7). For Graph Transformer, similar recursive relations can be derived for attention layers, with additional terms capturing interactions between  $\mathbf{W}_Q\mathbf{X}, \mathbf{W}_K\mathbf{X}, \mathbf{W}_V\mathbf{X}$ . For GCNs, we bounded gradient variations using norms of graph filters and layer parameters (e.g.,  $\|g(\mathbf{L})\|_2$ ,  $\|\mathbf{W}^{(k)}\|_2$ ). For Graph Transformer, this will be extended to: singular values of  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$  (analogous to  $C_g$  in GCNs), as they control the "strength" of feature projections and Lipschitz constants of softmax and feed-forward activations (replacing  $\alpha_\sigma$  for GCN activations, and leads to an analogous to Theorem 2 for deep GCNs.

#### 7 CONCLUSION AND FURTHER REMARKS

This paper explores the generalization of deep GCNs by providing an upper bound on their generalization gap. Our generalization bound is obtained based on the algorithmic stability of deep GCNs trained by the SGD algorithm. Our analysis demonstrates that the algorithmic stability of deep GCNs is contingent upon two factors: the largest absolute eigenvalue (or maximum singular value) of graph filter operators and the number of layers utilized. In particular, if the aforementioned eigenvalue (or singular value) remains invariant regardless of changes in the graph size, deep GCNs exhibit robust uniform stability, resulting in an enhanced generalization capability. Additionally, our results suggest that a greater number of layers can increase the generalization gap and subsequently degrade the performance of deep GCNs. This provides guidance for designing well-performing deep GCNs with a proper number of layers [61]. Most importantly, the result of single-layer GCNs in [36] can be regarded as a special case of our results in deep GCNs without hidden layers.

While our study is primarily focused on exploring the fundamental principles of generalizability and stability in the context of a simple deep GCN model framework, the theoretical insights obtained here can also offer preliminary perspectives on several research topics that have drawn increasing attention in the graph neural network community. These include, among others, the over-smoothing problem in deep architectures [62], [63], the design of models tailored for heterophilic graphs [64], [65], and the emerging topic of graph out-of-distribution (OOD) generalization [66], [67]. Our theoretical study can provide potential hints toward these directions, but more fine-grained and comprehensive work is still needed to fully address them. Below, we elaborate on these aspects in turn, aiming to clarify their conceptual connections with our work, outline possible directions for extending our theoretical framework, and highlight three open and challenging questions that can serve as seeds for future exploration.

How can the impact of over-smoothing in deep GCNs be mitigated? We first note that, given a trivial deep GCN model characterized by over-smoothed node embeddings (which typically result in significant training errors), our theoretical upper bound still holds — that is, for a given graph filter, an increase in layers could potentially increase this upper bound in a probabilistic sense. This also motivates the exploration of advanced deep GCN models that incorporate mechanisms to counteract over-smoothing, such as the skip connection technique used in GCNII [42] and its follow-up works. As detailed in Section 5, our theoretical results can in fact be extended to the setting of GCNII, thereby providing analytical support for architectures that integrate skip connections. In both theory and practice, reducing the maximum absolute eigenvalue of graph filter operators is achievable through the strategic implementation of skip connections across layers, which can potentially reduce the generalization gap. From this perspective, our findings may inspire further studies into sophisticated deep GCN architectures designed to mitigate over-smoothing, offering a promising direction for both theoretical and practical advancements.

What is the role of heterophily in GCN generalization? It is also valuable to consider extending our theoretical analysis to models specifically designed for heterophilic graphs, where nodes often connect to neighbors with dissimilar labels. This would require incorporating the homophily/heterophily ratio of the input graph signal into the upper bound estimation, thereby capturing how graph signal characteristics influence generalization. Although our empirical study here considers two types of low-pass filters on homophilic benchmark datasets (Cora, Citeseer, Pubmed), our theoretical framework is not restricted to low-pass scenarios alone. As remarked in Section 4.2, the analysis framework is in principle applicable to a broader range of filtering schemes; however, the derivations in our proofs do not explicitly examine the impact of specific quantities such as the homophily/heterophily ratio, leaving this as an open aspect for further refinement. To ensure a consistent and fair empirical evaluation, as demonstrated in [36], we adopt homophilic datasets that are standard in prior stability and generalization analyses of GCNs. For analyses involving high-pass filters, it would be appropriate to engage with heterophilic benchmark datasets (e.g., Texas, Wisconsin, Cornell). Relevant to this discussion is the recent work [48], which employs analytical tools from statistical physics and random matrix theory to precisely characterize generalization in simple GCNs on the contextual stochastic block model (CSBM). Such studies, although based on specific graph signal assumptions, could inspire refinements to our theoretical framework by jointly considering graph signal characteristics (homophily/heterophily) and model complexities (filter types, depth, and width).

Can insights from in-distribution generalization inform OOD generalization? Beyond the above considerations, another relevant line of research that has recently attracted considerable attention is graph out-of-distribution (OOD) generalization [66], [67]. It is worth clarifying that the problem setting and theoretical assumptions in OOD generalization are distinct from those in the in-distribution generalization framework considered in this work. In-distribution generalization focuses on scenarios where both training and test data are drawn from the same underlying distribution, enabling rigorous analysis under well-defined stochastic assumptions, such as those adopted in our stability-based framework. In contrast, OOD generalization addresses cases involving distribution shifts, which often require additional modeling principles (e.g., invariance to spurious correlations, causal structure modeling, or domain adaptation techniques) and seek performance guarantees that hold across domains. Despite these differences, the two areas can be mutually beneficial: in-distribution analyses, such as our characterization of bias-variance trade-offs and the influence of spectral properties of graph filters on generalization, may offer insights for developing more OOD-robust architectures; conversely, OOD-oriented approaches, such as invariant risk minimization or causal subgraph intervention, may inspire new regularization schemes or architectural components that also enhance in-distribution performance. Related to this discussion, the authors in [68] analyze a one-layer GCN trained on the CSBM via logistic regression, providing theoretical insights into improved linear separability and out-of-distribution generalization in semi-supervised node classification. Extending the current stability-based framework

to accommodate mild forms of distribution shift thus presents an appealing research direction that could bridge these two lines of work and advance the understanding of generalization in graph neural networks.

Taken together, these discussions highlight that our theoretical framework, while developed under a specific in-distribution setting, has the potential to be extended and adapted to address a broader range of challenges in graph learning.

Building on the above open questions, which outline core challenges for future exploration, it is also important to consider more concrete research directions and methodological extensions. For example, the theoretical analysis presented in this study could be extended to encompass other commonly used learning algorithms in graph neural networks, moving beyond the scope of SGD. Our theoretical results may also inform the exploration of strategies to enhance the generalization capability of deep graph neural networks, such as investigating the efficacy of regularization techniques, conducting advanced network architecture searches, or developing adaptive graph filters. In addition, establishing the potential connection between model stability, generalization, and the issues of over-smoothing and over-squashing represents another promising avenue. Understanding these interrelationships could contribute to the development of novel techniques and algorithms that address these challenges, thereby complementing the broader problem-oriented directions discussed above and improving the overall effectiveness of deep graph neural networks in dealing with more complex tasks.

#### **ACKNOWLEDGMENT**

The authors also wish to thank Dr. Yi Wang (City University of Hong Kong, Hong Kong SAR, China) and Dr. Xianchen Zhou (National University of Defense Technology, China) for their insightful discussions and dedicated assistance with the experimental studies.

# **APPENDIX: PRELIMINARIES**

The proofs of our main results are given in this section. We first make some statements about the notations used in the paper.  $\mathbf{W}^{\top}$  denotes the transpose of a matrix  $\mathbf{W}$ ; the (i,j)-entry of  $\mathbf{W}$  is denoted as  $\mathbf{W}_{ij}$ ; however when contributing to avoid confusion, the alternative notation  $\mathbf{W}(i,j)$  will be used.  $\|\cdot\|_2$  denotes the 2-norm of a matrix or vector and  $\|\cdot\|_F$  denotes the Frobenius norm.  $\delta_i$  denotes the unit pulse signal at node i that all elements are 0 except the i-th one, which is 1. Let  $f: \mathbb{R}^{m \times n} \to \mathbb{R}$  be a real-valued function of variable  $\mathbf{W} \in \mathbb{R}^{m \times n}$ . Then, the gradient of f with respect to  $\mathbf{W}$  is denoted as

 $\nabla_{\mathbf{W}} f = \frac{\partial f}{\partial \mathbf{W}} = \left(\frac{\partial f}{\partial \mathbf{W}_{ij}}\right) \in \mathbb{R}^{m \times n}.$ 

To make it easier to understand the derivation of our results, we first provide the following inequalities, which will be used frequently in the derivation.

For any matrix  $A_1$ ,  $A_2$ ,  $A'_1$  and  $A'_2$ , we have:

•  $\|\mathbf{A}_1\mathbf{A}_2\|_F \leq \|\mathbf{A}_1\|_2 \|\mathbf{A}_2\|_F$ . To prove this, let  $\mathbf{A}_1 = \mathbf{U}\Sigma\mathbf{V}^{\top}$  be the SVD of  $\mathbf{A}_1$ , where  $\mathbf{U}$  and  $\mathbf{V}$  are both orthogonal matrix. Then,

$$\|\mathbf{A}_1\mathbf{A}_2\|_F = \|\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^{\top}\mathbf{A}_2\|_F = \|\boldsymbol{\Sigma}\mathbf{V}^{\top}\mathbf{A}_2\|_F \leq \|\boldsymbol{\Sigma}\|_2\|\mathbf{V}^{\top}\mathbf{A}_2\|_F = \|\mathbf{A}_1\|_2\|\mathbf{A}_2\|_F.$$

Similarly, we also have  $\|\mathbf{A}_1\mathbf{A}_2\|_F \leq \|\mathbf{A}_1\|_F \|\mathbf{A}_2\|_2$ .

•  $\|\mathbf{A}_1\mathbf{A}_2 - \mathbf{A}_1'\mathbf{A}_2'\|_F \le \|\mathbf{A}_1 - \mathbf{A}_1'\|_F \|\mathbf{A}_2\|_2 + \|\mathbf{A}_1'\|_F \|\mathbf{A}_2 - \mathbf{A}_2'\|_2$ . To show this, note that

$$\begin{aligned} \|\mathbf{A}_{1}\mathbf{A}_{2} - \mathbf{A}_{1}'\mathbf{A}_{2}'\|_{F} &= \|(\mathbf{A}_{1} - \mathbf{A}_{1}')\mathbf{A}_{2} + \mathbf{A}_{1}'(\mathbf{A}_{2} - \mathbf{A}_{2}')\|_{F} \\ &\leq \|(\mathbf{A}_{1} - \mathbf{A}_{1}')\mathbf{A}_{2}\|_{F} + \|\mathbf{A}_{1}'(\mathbf{A}_{2} - \mathbf{A}_{2}')\|_{F}. \end{aligned}$$

Then, the proof is complete using the first inequality  $\|\mathbf{A}_1\mathbf{A}_2\|_F \leq \|\mathbf{A}_1\|_F \|\mathbf{A}_2\|_2$ ,

•  $\|\mathbf{A}_1 \odot \mathbf{A}_2\|_F \le \alpha \|\mathbf{A}_1\|_F \le \|\mathbf{A}_1\|_F \|\mathbf{A}_2\|_F$ , where  $\alpha$  is the maximum absolute value of the entries of  $\mathbf{A}_2$ . Note that  $\alpha \|\mathbf{A}_1\|_F \le \|\mathbf{A}_1\|_F \|\mathbf{A}_2\|_F$  holds true because  $\alpha \le \|\mathbf{A}_2\|_F$ . Furthermore,

$$\|\mathbf{A}_1 \odot \mathbf{A}_2\|_F = \sqrt{\sum_{ij} \left(\mathbf{A}_1(i,j)\mathbf{A}_2(i,j)\right)^2}$$

$$\leq \sqrt{\sum_{ij} \left(\alpha \mathbf{A}_1(i,j)\right)^2} \leq \alpha \sqrt{\sum_{ij} \left(\mathbf{A}_1(i,j)\right)^2} = \alpha \|\mathbf{A}_1\|_F.$$

# **APPENDIX A: GRADIENT COMPUTATION FOR SGD**

To work with the SGD algorithm, we provide a recursive formula for the gradient of the final output  $f(\mathbf{x}|\theta)$  at node  $\mathbf{x}$  in the GCNs model (1) with respect to the learnable parameters.

• For the final layer,

$$\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) = \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} \right) \left[ \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \right]^{\top}, \tag{A.1}$$

• For the hidden layer  $k = 1, 2, \dots, K$ ,

$$\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) = \left[ g(\mathbf{L}) \mathbf{X}^{(k-1)} \right]^{\top} \left( \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} \right), \tag{A.2}$$

where  $\mathbf{R}^{(k)} := \nabla \sigma(g(\mathbf{L})\mathbf{X}^{(k-1)}\mathbf{W}^{(k)})$  and

$$\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k-1)}} = g(\mathbf{L})^{\top} \left( \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} \right) \left[ \mathbf{W}^{(k)} \right]^{\top}, \tag{A.3}$$

with

$$\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} = \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} \right) \left[ \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \right]^{\top} \mathbf{w}^{\top}. \tag{A.4}$$

The notation  $\odot$  represents the Hadamard product of two matrices. (A.1) and (A.4) are easy to verify, while (A.2) and (A.3) are not. In the following, a detailed procedure is provided to derive (A.2) and (A.3).

First, since  $\mathbf{X}_{ij}^{(k)} = \sigma(\boldsymbol{\delta}_i^{\top} g(\mathbf{L}) \mathbf{X}^{(k-1)} \mathbf{W}^{(k)} \boldsymbol{\delta}_j)$ ,

$$\begin{split} \frac{\partial \mathbf{X}_{ij}^{(k)}}{\partial \mathbf{W}^{(k)}} &= \frac{\partial \sigma \left( \boldsymbol{\delta}_{i}^{\top} g(\mathbf{L}) \mathbf{X}^{(k-1)} \mathbf{W}^{(k)} \boldsymbol{\delta}_{j} \right)}{\partial \mathbf{W}^{(k)}} \\ &= \nabla \sigma \left( \boldsymbol{\delta}_{i}^{\top} g(\mathbf{L}) \mathbf{X}^{(k-1)} \mathbf{W}^{(k)} \boldsymbol{\delta}_{j} \right) \frac{\partial \left\{ \boldsymbol{\delta}_{i}^{\top} g(\mathbf{L}) \mathbf{X}^{(k-1)} \mathbf{W}^{(k)} \boldsymbol{\delta}_{j} \right\}}{\partial \mathbf{W}^{(k)}} \\ &= \nabla \sigma \left( \boldsymbol{\delta}_{i}^{\top} g(\mathbf{L}) \mathbf{X}^{(k-1)} \mathbf{W}^{(k)} \boldsymbol{\delta}_{j} \right) \left[ g(\mathbf{L}) \mathbf{X}^{(k-1)} \right]^{\top} \boldsymbol{\delta}_{i} \boldsymbol{\delta}_{i}^{\top}, \end{split}$$

and

$$\frac{\partial \mathbf{X}_{ij}^{(k)}}{\partial \mathbf{X}^{(k-1)}} = \frac{\partial \sigma(\boldsymbol{\delta}_i^{\top} g(\mathbf{L}) \mathbf{X}^{(k-1)} \mathbf{W}^{(k)} \boldsymbol{\delta}_j)}{\partial \mathbf{X}^{(k-1)}} = \nabla \sigma(\boldsymbol{\delta}_i^{\top} g(\mathbf{L}) \mathbf{X}^{(k-1)} \mathbf{W}^{(k)} \boldsymbol{\delta}_j) g(\mathbf{L})^{\top} \boldsymbol{\delta}_i \boldsymbol{\delta}_j^{\top} [\mathbf{W}^{(k)}]^{\top}.$$

Let  $\mathbf{R}^{(k)} = \nabla \sigma (g(\mathbf{L})\mathbf{X}^{(k-1)}\mathbf{W}^{(k)})$ . Then,

$$\begin{split} \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{W}^{(k)}} &= \sum_{i,j} \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}_{ij}^{(k)}} \cdot \frac{\partial \mathbf{X}_{ij}^{(k)}}{\partial \mathbf{W}^{(k)}} = \sum_{i,j} \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} (i,j) \cdot \frac{\partial \mathbf{X}_{ij}^{(k)}}{\partial \mathbf{W}^{(k)}} \\ &= \sum_{i,j} \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} (i,j) \cdot \mathbf{R}^{(k)} (i,j) \left[ g(\mathbf{L}) \mathbf{X}^{(k-1)} \right]^{\top} \boldsymbol{\delta}_{i} \boldsymbol{\delta}_{j}^{\top} \\ &= \left[ g(\mathbf{L}) \mathbf{X}^{(k-1)} \right]^{\top} \sum_{i,j} \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} (i,j) \cdot \mathbf{R}^{(k)} (i,j) \boldsymbol{\delta}_{i} \boldsymbol{\delta}_{j}^{\top} \\ &= \left[ g(\mathbf{L}) \mathbf{X}^{(k-1)} \right]^{\top} \left( \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} \right), \end{split}$$

and

$$\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k-1)}} = \sum_{i,j} \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}_{ij}^{(k)}} \cdot \frac{\partial \mathbf{X}_{ij}^{(k)}}{\partial \mathbf{X}^{(k-1)}}$$

$$= g(\mathbf{L})^{\top} \Big( \sum_{i,j} \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} (i,j) \cdot \mathbf{R}^{(k)} (i,j) \boldsymbol{\delta}_{i} \boldsymbol{\delta}_{j}^{\top} \Big) [\mathbf{W}^{(k)}]^{\top}$$

$$= g(\mathbf{L})^{\top} \Big( \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} \Big) [\mathbf{W}^{(k)}]^{\top}.$$

This completes the derivation of (A.2) and (A.3).

Based on the above recursive formula, we prove the following lemma recursively.

*Lemma 5.* Let the assumptions made in Section 4.1 hold. Then, we have the following results for the GCNs model (1) during the training procedure.

• Hidden layer output  $\mathbf{X}^{(k)}(k=1,2...,K)$  satisfies

$$\|\mathbf{X}^{(k)}\|_F \le B^k \alpha_\sigma^k C_g^k C_{\mathbf{X}}.\tag{A.5}$$

• The gradient of f with respect to  $\mathbf{X}^{(k)}$  (k = 1, 2, ..., K) satisfies

$$\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F \le B^{K+1-k} \alpha_{\sigma}^{K+1-k} C_g^{K+1-k}. \tag{A.6}$$

• The gradient of f with respect to  $\mathbf{W}^{(k)}$   $(k = 1, \dots, K+1)$  satisfies

$$\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta)\|_{F} \le B^{K} \alpha_{\sigma}^{K+1} C_{q}^{K+1} C_{\mathbf{X}}, \tag{A.7}$$

where  $\mathbf{W}^{(K+1)} := \mathbf{w}$ .

*Proof*. Now, we give a complete proof for Lemma 5.

• Firstly, for k = 1, 2, ..., K, since  $\|\sigma(\mathbf{Z})\|_F \le \alpha_{\sigma} \|\mathbf{Z}\|_F$  holds for any matrix  $\mathbf{Z}$ , we have

$$\|\mathbf{X}^{(k)}\|_F = \|\sigma(g(\mathbf{L})\mathbf{X}^{(k-1)}\mathbf{W}^{(k)})\|_F \le \alpha_\sigma \|g(\mathbf{L})\mathbf{X}^{(k-1)}\mathbf{W}^{(k)}\|_F.$$

Then, by applying the inequality  $\|\mathbf{A}_1\mathbf{A}_2\|_F \leq \|\mathbf{A}_1\|_2 \|\mathbf{A}_2\|_F$  twice, we obtain  $\|\mathbf{X}^{(k)}\|_F \leq B\alpha_\sigma C_g \|\mathbf{X}^{(k-1)}\|_F$ . Note that  $\|\mathbf{X}^{(1)}\|_F \leq B\alpha_\sigma C_g \|\mathbf{X}^{(0)}\|_F = B\alpha_\sigma C_g C_{\mathbf{X}}$ , it further yields that

$$\|\mathbf{X}^{(k)}\|_F \le B^k \alpha_\sigma^k C_q^k C_{\mathbf{X}}, \quad k = 1, 2, \dots, K,$$

which completes the proof of (A.5).

• To show (A.6), note that for  $k=1,2,\ldots,K-1$ , by applying  $\|\mathbf{A}_1\mathbf{A}_2\|_F \le \|\mathbf{A}_1\|_2 \|\mathbf{A}_2\|_F$  twice, we obtain

$$\left\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\right\|_{F} = \left\|g(\mathbf{L})^{\top} \left(\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)}\right) \left[\mathbf{W}^{(k+1)}\right]^{\top}\right\|_{F} \leq \|g(\mathbf{L})\|_{2} \left\|\left(\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)}\right)\right\|_{F} \|\mathbf{W}^{(k+1)}\|_{2}.$$

Since  $C_g = \|g(\mathbf{L})\|_2$ ,  $\|\mathbf{W}^{(k+1)}\|_2 \leq B$  and the absolute value of the elements in  $\mathbf{R}^{(k+1)}$  is less than  $\alpha_{\sigma}$ , we further have  $\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F \leq B\alpha_{\sigma}C_g\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k+1)}}\|_F$ . Meanwhile, since  $|\nabla\sigma(\mathbf{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})\mathbf{X}^{(K)}\mathbf{w})| \leq \alpha_{\sigma}$ ,

$$\left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} \right\|_F = \left\| \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} \right) \left[ \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \right]^{\top} \mathbf{w} \right\|_F \le B \alpha_{\sigma} C_g.$$

Therefore, for  $k = 1, 2, \dots, K$ ,

$$\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F \le B^{K+1-k} \alpha_{\sigma}^{K+1-k} C_g^{K+1-k}.$$

This completes the proof of (A.6).

• Now, let's prove (A.7). Firstly, note that  $|\nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w})| \leq \alpha_{\sigma}$ , so

$$\left\|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta)\right\|_F = \left\|\nabla\sigma\left(\boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L})\mathbf{X}^{(K)}\mathbf{w}\right)\left[g(\mathbf{L})\mathbf{X}^{(K)}\right]^{\top} \boldsymbol{\delta}_{\mathbf{x}}\right\|_F \leq \alpha_{\sigma} \|\mathbf{X}^{(K)}\|_F \|\boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L})\|_2.$$

Combining (A.5) and  $\|\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})\|_{2} \leq C_{q}$ , we have

$$\|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta)\|_{\mathcal{F}} \leq B^K \alpha_{\sigma}^{K+1} C_{\sigma}^{K+1} C_{\mathbf{x}}.$$

Furthermore, for  $k=1,2,\ldots,K$ , by applying  $\|\mathbf{A}_1\mathbf{A}_2\|_F \leq \|\mathbf{A}_1\|_2\|\mathbf{A}_2\|_F$  twice, it yields

$$\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta)\|_F = \|[g(\mathbf{L})\mathbf{X}^{(k-1)}]^\top \left(\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)}\right)\|_F \le \|g(\mathbf{L})\|_2 \|\mathbf{X}^{(k-1)}\|_F \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)}\|_F.$$

Since the absolute value of the elements in  $\mathbf{R}^{(k)}$  is less than  $\alpha_{\sigma}$ , we have

$$\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta)\|_F \le \alpha_{\sigma} C_g \|\mathbf{X}^{(k-1)}\|_F \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F \le B^K \alpha_{\sigma}^{K+1} C_g^{K+1} C_{\mathbf{X}},$$

which holds by combining (A.5) and (A.6). This completes the proof of (A.7).

### APPENDIX B:PROOF OF LEMMA 2

To prove Lemma 2, we first provide the following lemma to show the variation of output in each layer for two GCNs with different learned parameters  $\theta = \{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \dots, \mathbf{W}^{(K)}, \mathbf{w}\}$  and  $\theta' = \{\mathbf{W}^{(1)'}, \mathbf{W}^{(2)'}, \dots, \mathbf{W}^{(K)'}, \mathbf{w}'\}$ . Let  $\mathbf{X}^{(k)}$  and  $\mathbf{X}^{(k)'}$  be their output of the hidden layer, as well as  $f(\mathbf{x}|\theta)$  and  $f(\mathbf{x}|\theta')$  the final output of node  $\mathbf{x}$ . The following lemma provides a bound of  $\mathbf{X}^{(k)} - \mathbf{X}^{(k)'}$  and  $f(\mathbf{x}|\theta) - f(\mathbf{x}|\theta')$  based on  $\Delta \theta = \{\Delta \mathbf{W}^{(1)}, \dots, \Delta \mathbf{W}^{(K)}, \Delta \mathbf{w}\}$ .

**Lemma 6.** Consider two GCNs with parameters  $\theta$  and  $\theta'$ , respectively. Then, we obtain the following results for their variations.

• Their variation of outputs in hidden layers  $\triangle \mathbf{X}^{(k)} := \mathbf{X}^{(k)} - \mathbf{X}^{(k)'}$  (k = 1, 2, ..., K) satisfies

$$\|\Delta \mathbf{X}^{(k)}\|_F \le B^{k-1} \alpha_{\sigma}^k C_g^k C_{\mathbf{X}} \Big( \sum_{i=1}^k \|\Delta \mathbf{W}^{(j)}\|_2 \Big).$$
 (A.8)

• Furthermore, for the final output of node x,

$$|f(\mathbf{x}|\theta) - f(\mathbf{x}|\theta')| \le B^K \alpha_{\sigma}^{K+1} C_{\mathbf{x}}^{K+1} C_{\mathbf{x}} ||\Delta\theta||_*. \tag{A.9}$$

*Proof*: To prove (A.8), we first have that for k = 1, 2, ..., K,

$$\|\triangle \mathbf{X}^{(k)}\|_{F} = \|\mathbf{X}^{(k)} - \mathbf{X}^{(k)'}\|_{F} = \|\sigma(g(\mathbf{L})\mathbf{X}^{(k-1)}\mathbf{W}^{(k)}) - \sigma(g(\mathbf{L})\mathbf{X}^{(k-1)'}\mathbf{W}^{(k)'})\|_{F}.$$

Since  $\|\sigma(\mathbf{Z})\|_F \leq \alpha_{\sigma} \|\mathbf{Z}\|_F$  holds for any matrix  $\mathbf{Z}$ , we have

$$\|\Delta \mathbf{X}^{(k)}\|_{F} \leq \alpha_{\sigma} \|g(\mathbf{L}) (\mathbf{X}^{(k-1)} \mathbf{W}^{(k)} - \mathbf{X}^{(k-1)'} \mathbf{W}^{(k)'})\|_{F} \leq \alpha_{\sigma} \|g(\mathbf{L})\|_{2} \cdot \|\mathbf{X}^{(k-1)} \mathbf{W}^{(k)} - \mathbf{X}^{(k-1)'} \mathbf{W}^{(k)'}\|_{F}.$$

Note that

$$\begin{aligned} \|\mathbf{X}^{(k-1)}\mathbf{W}^{(k)} - \mathbf{X}^{(k-1)'}\mathbf{W}^{(k)'}\|_{F} &\leq \|\mathbf{X}^{(k-1)}\|_{F} \|\mathbf{W}^{(k)} - \mathbf{W}^{(k)'}\|_{2} + \|\mathbf{X}^{(k-1)} - \mathbf{X}^{(k-1)'}\|_{F} \|\mathbf{W}^{(k)'}\|_{2} \\ &= \|\mathbf{X}^{(k-1)}\|_{F} \|\triangle\mathbf{W}^{(k)}\|_{2} + \|\triangle\mathbf{X}^{(k-1)}\|_{F} \|\mathbf{W}^{(k)'}\|_{2}. \end{aligned}$$

Then, combining (A.5) and  $\|\mathbf{W}^{(k)'}\|_2 \leq B$ , we obtain

$$\|\mathbf{X}^{(k-1)}\mathbf{W}^{(k)} - \mathbf{X}^{(k-1)'}\mathbf{W}^{(k)'}\|_{F} \le B^{k-1}\alpha_{\sigma}^{k-1}C_{q}^{k-1}C_{\mathbf{X}}\|\triangle\mathbf{W}^{(k)}\|_{2} + B\|\triangle\mathbf{X}^{(k-1)}\|_{F}.$$

Thus,

$$\|\triangle \mathbf{X}^{(k)}\|_{F} \leq \alpha_{\sigma} \|g(\mathbf{L})\|_{2} \cdot \|\mathbf{X}^{(k-1)}\mathbf{W}^{(k)} - \mathbf{X}^{(k-1)'}\mathbf{W}^{(k)'}\|_{F} \leq B^{k-1}\alpha_{\sigma}^{k}C_{q}^{k}C_{\mathbf{X}}\|\triangle \mathbf{W}^{(k)}\|_{2} + B\alpha_{\sigma}C_{q}\|\triangle \mathbf{X}^{(k-1)}\|_{F}.$$

Then, since  $\|\triangle \mathbf{X}^{(1)}\|_F \leq \alpha_{\sigma} C_q C_{\mathbf{X}} \|\triangle \mathbf{W}^{(1)}\|_2$ , we have

$$\|\triangle \mathbf{X}^{(k)}\|_F \le B^{k-1} \alpha_{\sigma}^k C_g^k C_{\mathbf{X}} \Big( \sum_{i=1}^k \|\triangle \mathbf{W}^{(j)}\|_2 \Big),$$

holds for any k = 1, 2, ..., K. This completely proves (A.8).

Furthermore, for the final output, using the Lipschitz property of  $\sigma(\cdot)$ , we have

$$|f(\mathbf{x}|\theta) - f(\mathbf{x}|\theta')| = |\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})\mathbf{X}^{(K)}\mathbf{w}) - \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})\mathbf{X}^{(K)'}\mathbf{w}')| \le \alpha_{\sigma}|\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})(\mathbf{X}^{(K)}\mathbf{w} - \mathbf{X}^{(K)'}\mathbf{w}')|.$$

Note that

$$|\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})(\mathbf{X}^{(K)}\mathbf{w} - \mathbf{X}^{(K)'}\mathbf{w}')| \leq \|\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})\|_{2} \cdot \|\mathbf{X}^{(K)}\mathbf{w} - \mathbf{X}^{(K)'}\mathbf{w}'\|_{F} \leq C_{\sigma}(\|\mathbf{X}^{(K)}\|_{F}\|\triangle\mathbf{w}\|_{2} + \|\triangle\mathbf{X}^{(K)}\|_{F}\|\mathbf{w}'\|_{2}).$$

Combining (A.5) and (A.8), we further have

$$|\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})(\mathbf{X}^{(K)}\mathbf{w} - \mathbf{X}^{(K)'}\mathbf{w}')| \leq B^{K}\alpha_{\sigma}^{K}C_{g}^{K+1}C_{\mathbf{X}}(\|\triangle\mathbf{w}\|_{2} + \sum_{j=1}^{K}\|\triangle\mathbf{W}^{(j)}\|_{2}) = B^{K}\alpha_{\sigma}^{K}C_{g}^{K+1}C_{\mathbf{X}}\|\triangle\theta\|_{*}.$$

Thus,

$$|f(\mathbf{x}|\theta) - f(\mathbf{x}|\theta')| \le \alpha_{\sigma} |\mathbf{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) (\mathbf{X}^{(K)} \mathbf{w} - \mathbf{X}^{(K)'} \mathbf{w}')| \le B^{K} \alpha_{\sigma}^{K+} C_{q}^{K+1} C_{\mathbf{X}} ||\Delta \theta||_{*},$$

which completes the proof of (A.9).

*Proof of Lemma* 2: Now, we are ready to prove Lemma 2 based on Lemma 6. For any  $\mathbf{z} = (\mathbf{x}, y)$  taken from  $\mathcal{D}$ , we denote by  $\hat{y} = f(\mathbf{x}|\theta_T)$  and  $\hat{y}' = f(\mathbf{x}|\theta_T')$ . Firstly, using the Lipschitz property of loss function  $\ell(\cdot, \cdot)$ , we have

$$\sup_{\mathcal{S}, \mathbf{z}} \left| \mathbb{E}_{\mathcal{A}} \left[ \ell(\hat{y}, y) \right] - \mathbb{E}_{\mathcal{A}} \left[ \ell(\hat{y}', y) \right] \right| = \sup_{\mathcal{S}, z} \left| \mathbb{E}_{\mathcal{A}} \left[ \ell(f(\mathbf{x} | \theta_T), y) - \ell(f(\mathbf{x} | \theta_T'), y) \right] \right| \le \alpha_{\ell} \sup_{\mathbf{x}} \mathbb{E}_{\mathcal{A}} \left[ \left| f(\mathbf{x} | \theta_T) - f(\mathbf{x} | \theta_T') \right| \right]$$

Then, according to (A.9),

$$\sup_{S_{\mathcal{A}}} \left| \mathbb{E}_{\mathcal{A}} \left[ \ell(\hat{y}, y) \right] - \mathbb{E}_{\mathcal{A}} \left[ \ell(\hat{y}', y) \right] \right| \leq \alpha_{\ell} B^{K} \alpha_{\sigma}^{K+1} C_{g}^{K+1} C_{\mathbf{X}} \cdot \mathbb{E}_{\mathcal{A}} \left[ \| \triangle \theta_{T} \|_{*} \right].$$

This completes the proof of Lemma 2.

# APPENDIX C: PROOF OF LEMMA 3 AND LEMMA 4

To prove Lemma 3 and Lemma 4, we should first prove the following lemma.

**Lemma 7.** Consider two GCNs with parameters  $\theta$  and  $\theta'$ , respectively. Then, their variation of gradients of f with respect to  $\{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(K)}, \mathbf{w}\}$  satisfies

$$\left\| \nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta') \right\|_{F} \le \left( \upsilon_{\sigma} B^{2K} \alpha_{\sigma}^{2K} C_{g}^{2K+2} C_{\mathbf{X}}^{2} + B^{K-1} \alpha_{\sigma}^{K+1} C_{g}^{K+1} C_{\mathbf{X}} \right) \|\Delta \theta\|_{*}, \tag{A.10}$$

and for k = 1, 2, ..., K,

$$\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_{F} \le \left(\nu_{\sigma} B^{2K} \alpha_{\sigma}^{2K} C_{g}^{2K+2} C_{\mathbf{X}}^{2} + B^{K-1} \alpha_{\sigma}^{K+1} C_{g}^{K+1} C_{\mathbf{X}}\right) \|\Delta\theta\|_{*} + \rho_{k} \|\Delta\theta\|_{*}, \quad (A.11)$$

where

$$\rho_k := \nu_{\sigma} (B\alpha_{\sigma} C_g)^{K+k-1} C_g^2 C_{\mathbf{X}}^2 \Big( \sum_{j=0}^{K-k} (B\alpha_{\sigma} C_g)^j \Big). \tag{A.12}$$

*Proof.* First, according to the proof of (A.8) and (A.9), the following holds true for k = 1, 2, ..., K + 1:

$$\|\mathbf{X}^{(k-1)}\mathbf{W}^{(k)} - \mathbf{X}^{(k-1)'}\mathbf{W}^{(k)'}\|_{F} \leq B^{k-1}\alpha_{\sigma}^{k-1}C_{g}^{k-1}C_{\mathbf{X}}\|\Delta\mathbf{W}^{(k)}\|_{2} + B\|\Delta\mathbf{X}^{(k-1)}\|_{F}$$

$$\leq B^{k-1}\alpha_{\sigma}^{k-1}C_{g}^{k-1}C_{\mathbf{X}}\left(\sum_{j=1}^{k}\|\Delta\mathbf{W}^{(j)}\|_{2}\right), \tag{A.13}$$

where  $\mathbf{W}^{(K+1)} = \mathbf{w}$ .

We now prove (A.10). First, applying  $A_1A_2 - A_1'A_2' = (A_1 - A_1')A_2 + A_1'(A_2 - A_2')$ , we have

$$\begin{split} \|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta')\|_{F} &= \left\| \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} \right) [g(\mathbf{L}) \mathbf{X}^{(K)}]^{\top} \boldsymbol{\delta}_{\mathbf{x}} - \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}' \right) [g(\mathbf{L}) \mathbf{X}^{(K)'}]^{\top} \boldsymbol{\delta}_{\mathbf{x}} \right\|_{F} \\ &\leq \left\| \left( \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} \right) - \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}' \right) \right) [g(\mathbf{L}) \mathbf{X}^{(K)}]^{\top} \boldsymbol{\delta}_{\mathbf{x}} \right\|_{F} \\ &+ \left\| \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}' \right) [g(\mathbf{L}) \triangle \mathbf{X}^{(K)}]^{\top} \boldsymbol{\delta}_{\mathbf{x}} \right\|_{F}. \end{split}$$

Using the  $\nu_{\sigma}$ -smooth property of  $\sigma(\cdot)$  and applying  $\|\mathbf{A}_1\mathbf{A}_2\|_F \leq \|\mathbf{A}_1\|_2 \|\mathbf{A}_2\|_F$ , we have

$$\begin{split} & \left\| \left( \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} \right) - \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}' \right) \right) [g(\mathbf{L}) \mathbf{X}^{(K)}]^{\top} \boldsymbol{\delta}_{\mathbf{x}} \right\|_{F} \\ \leq & \left| \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} \right) - \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}' \right) \right| \cdot \left\| [g(\mathbf{L}) \mathbf{X}^{(K)}]^{\top} \boldsymbol{\delta}_{\mathbf{x}} \right\|_{F} \\ \leq & v_{\sigma} |\boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}' | \cdot \| \mathbf{X}^{(K)} \|_{F} \| \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \|_{2} \\ \leq & v_{\sigma} C_{g} \| \mathbf{X}^{(K)} \mathbf{w} - \mathbf{X}^{(K)'} \mathbf{w}' \|_{F} \cdot \| \mathbf{X}^{(K)} \|_{F} \cdot C_{g}, \end{split}$$

and since  $|\nabla \sigma(\cdot)| \leq \alpha_{\sigma}$ ,  $\|\nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}') [g(\mathbf{L}) \triangle \mathbf{X}^{(K)}]^{\top} \boldsymbol{\delta}_{\mathbf{x}} \|_{F} \leq \alpha_{\sigma} C_{g} \|\triangle \mathbf{X}^{(K)}\|_{F}$ . Then, combining (A.5), (A.8) and (A.13), we have

$$\|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta')\|_F \le \left(\upsilon_{\sigma} B^{2K} \alpha_{\sigma}^{2K} C_g^{2K+2} C_{\mathbf{X}}^2 + B^{K-1} \alpha_{\sigma}^{K+1} C_g^{K+1} C_{\mathbf{X}}\right) \|\triangle \theta\|_*,$$

which completes the proof of (A.10).

Next, we turn to prove (A.11). First, for k = 1, 2, ..., K,

$$\begin{aligned} & \left\| \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta') \right\|_{F} \\ &= \left\| \left[ g(\mathbf{L}) \mathbf{X}^{(k-1)} \right]^{\top} \left( \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} \right) - \left[ g(\mathbf{L}) \mathbf{X}^{(k-1)'} \right]^{\top} \left( \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)'} \right) \right\|_{F} \\ &\leq \left\| g(\mathbf{L}) \triangle \mathbf{X}^{(k-1)} \right\|_{F} \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} \right\|_{F} + \left\| g(\mathbf{L}) \mathbf{X}^{(k-1)'} \right\|_{F} \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)'} \right\|_{F} \\ &\leq C_{g} \left\| \triangle \mathbf{X}^{(k-1)} \right\|_{F} \cdot \alpha_{\sigma} \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \right\|_{F} + C_{g} \left\| \mathbf{X}^{(k-1)'} \right\|_{F} \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)'} \right\|_{F}. \end{aligned}$$

Let

$$\gamma_k := \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)'} \right\|_F. \tag{A.14}$$

Then, combining (A.5), (A.6) and (A.8), we have

$$\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_{F} \le B^{K-1} \alpha_{\sigma}^{K+1} C_{g}^{K+1} C_{\mathbf{X}} \left( \sum_{j=1}^{k-1} \|\triangle \mathbf{W}^{(j)}\|_{2} \right) + B^{k-1} \alpha_{\sigma}^{k-1} C_{g}^{k} C_{\mathbf{X}} \cdot \gamma_{k}, \tag{A.15}$$

Next, we need to bound  $\gamma_k$ .

$$\gamma_{k} \leq \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \left( \mathbf{R}^{(k)} - \mathbf{R}^{(k)'} \right) \right\|_{F} + \left\| \left( \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \right) \odot \mathbf{R}^{(k)'} \right\|_{F} \\
\leq h_{k} + \alpha_{\sigma} \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \right\|_{F} \\
\leq h_{k} + \alpha_{\sigma} \left\| g(\mathbf{L})^{\top} \left( \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)} \right) \left[ \mathbf{W}^{(k+1)} \right]^{\top} - g(\mathbf{L})^{\top} \left( \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k+1)'} \right) \left[ \mathbf{W}^{(k+1)'} \right]^{\top} \right\|_{F} \\
\leq h_{k} + \alpha_{\sigma} \| g(\mathbf{L}) \|_{2} \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)} \right\|_{F} \| \Delta \mathbf{W}^{(k+1)} \|_{2} + \alpha_{\sigma} \| g(\mathbf{L}) \|_{2} \| \mathbf{W}^{(k+1)'} \|_{2} \gamma_{k+1} \\
\leq h_{k} + \alpha_{\sigma}^{2} C_{g} (B \alpha_{\sigma} C_{g})^{K-k} \| \Delta \mathbf{W}^{(k+1)} \|_{2} + B \alpha_{\sigma} C_{g} \gamma_{k+1}, \\$$

where  $h_k := \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \left( \mathbf{R}^{(k)} - \mathbf{R}^{(k)'} \right) \right\|_F$ . By (A.13),

$$\|\mathbf{R}^{(k)} - \mathbf{R}^{(k)'}\|_{F} = \|\nabla\sigma(g(\mathbf{L})\mathbf{X}^{(k-1)}\mathbf{W}^{(k)}) - \nabla\sigma(g(\mathbf{L})\mathbf{X}^{(k-1)'}\mathbf{W}^{(k)'})\|_{F}$$

$$\leq \nu_{\sigma}C_{g}\|\mathbf{X}^{(k-1)}\mathbf{W}^{(k)} - \mathbf{X}^{(k-1)'}\mathbf{W}^{(k)'}\|_{F}$$

$$\leq \nu_{\sigma}B^{k-1}\alpha_{\sigma}^{k-1}C_{g}^{k}C_{\mathbf{X}}\left(\sum_{i=1}^{k}\|\Delta\mathbf{W}^{(j)}\|_{2}\right). \tag{A.16}$$

Combining (A.6), we have

$$h_{k} = \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \left( \mathbf{R}^{(k)} - \mathbf{R}^{(k)'} \right) \right\|_{F} \le \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \right\|_{F} \cdot \| \mathbf{R}^{(k)} - \mathbf{R}^{(k)'} \|_{F}$$

$$\le \nu_{\sigma} B^{K} \alpha_{\sigma}^{K} C_{g}^{K+1} C_{\mathbf{X}} \left( \sum_{i=1}^{k} \| \triangle \mathbf{W}^{(j)} \|_{2} \right). \tag{A.17}$$

Let  $h_{\max} = \nu_{\sigma} B^K \alpha_{\sigma}^K C_q^{K+1} C_{\mathbf{X}} \|\Delta \theta\|_*$ . Then, it is easy to see that

$$h_k \le h_{\max}$$
 holds for all  $k = 1, 2, \dots, K$ . (A.18)

Therefore,

$$\gamma_k \le h_{\max} + \alpha_\sigma^2 C_g (B\alpha_\sigma C_g)^{K-k} \|\Delta \mathbf{W}^{(k+1)}\|_2 + B\alpha_\sigma C_g \cdot \gamma_{k+1}.$$

Furthermore, since

$$\begin{split} & \| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(K)}} \|_{F} \\ = & \| \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} \right) \left[ \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \right]^{\top} \mathbf{w}^{\top} - \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}' \right) \left[ \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \right]^{\top} \mathbf{w}'^{\top} \|_{F} \\ \leq & B C_{g} \| \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)} \mathbf{w} \right) - \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}' \right) \|_{F} + \| \nabla \sigma \left( \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \mathbf{X}^{(K)'} \mathbf{w}' \right) \left[ \boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L}) \right]^{\top} \triangle \mathbf{w}^{\top} \|_{F} \\ \leq & \alpha_{\sigma} C_{g} \| \Delta \mathbf{w} \|_{F} + \nu_{\sigma} B C_{g}^{2} \| \mathbf{X}^{(K)} \mathbf{w} - \mathbf{X}^{(K)'} \mathbf{w}' \|_{F} \\ \leq & \alpha_{\sigma} C_{g} \| \Delta \mathbf{w} \|_{2} + \nu_{\sigma} B^{K+1} \alpha_{\sigma}^{K} C_{g}^{K+2} C_{\mathbf{X}} \| \Delta \theta \|_{*}, \end{split}$$

we have

$$\begin{split} \gamma_{K} = & \| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} \odot \mathbf{R}^{(K)} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(K)}} \odot \mathbf{R}^{(K)'} \|_{F} \\ \leq & \| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} \odot (\mathbf{R}^{(K)} - \mathbf{R}^{(K)'}) \|_{F} + \| (\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(K)}}) \odot \mathbf{R}^{(K)'} \|_{F} \\ \leq & h_{K} + \alpha_{\sigma} \| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(K)}} \|_{F} \\ \leq & h_{\max} + \alpha_{\sigma}^{2} C_{g} \| \Delta \mathbf{w} \|_{2} + \nu_{\sigma} B^{K+1} \alpha_{\sigma}^{K+1} C_{g}^{K+2} C_{\mathbf{X}} \| \Delta \theta \|_{*}. \end{split}$$

Finally, based on the above recursive formula of  $\gamma_k$ , we have

$$\gamma_{k} \leq h_{\max} \left( \sum_{j=0}^{K-k} (B\alpha_{\sigma}C_{g})^{j} \right) + \alpha_{\sigma}^{2} C_{g} (B\alpha_{\sigma}C_{g})^{K-k} \left( \sum_{j=k+1}^{K+1} \|\triangle \mathbf{W}^{(j)}\|_{2} \right) \\
+ \nu_{\sigma} B^{K+1} \alpha_{\sigma}^{K+1} C_{g}^{K+2} C_{\mathbf{X}} (B\alpha_{\sigma}C_{g})^{K-k} \|\triangle \theta\|_{*} \\
\leq h_{\max} \left( \sum_{j=0}^{K-k} (B\alpha_{\sigma}C_{g})^{j} \right) + \alpha_{\sigma}^{2} C_{g} (B\alpha_{\sigma}C_{g})^{K-k} \left( \sum_{j=k+1}^{K+1} \|\triangle \mathbf{W}^{(j)}\|_{2} \right) \\
+ \nu_{\sigma} B^{2K+1-k} \alpha_{\sigma}^{2K+1-k} C_{g}^{2K+2-k} C_{\mathbf{X}} \|\triangle \theta\|_{*}, \tag{A.19}$$

where  $\triangle \mathbf{W}^{(K+1)} = \triangle \mathbf{w}$ . Finally, substituting (A.19) into (A.15),

$$\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_{F} \leq B^{K-1} \alpha_{\sigma}^{K+1} C_{g}^{K+1} C_{\mathbf{X}} \Big( \sum_{j=1}^{k-1} \|\Delta \mathbf{W}^{(j)}\|_{2} \Big) + B^{k-1} \alpha_{\sigma}^{k-1} C_{g}^{k} C_{\mathbf{X}} \cdot \gamma_{k}$$

$$\leq \Big( \nu_{\sigma} B^{2K} \alpha_{\sigma}^{2K} C_{g}^{2K+2} C_{\mathbf{X}}^{2} + B^{K-1} \alpha_{\sigma}^{K+1} C_{g}^{K+1} C_{\mathbf{X}} \Big) \|\Delta \theta\|_{*}$$

$$+ \nu_{\sigma} B^{K+k-1} \alpha_{\sigma}^{K+k-1} C_{g}^{K+k+1} C_{\mathbf{X}}^{2} \Big( \sum_{j=0}^{K-k} (B\alpha_{\sigma} C_{g})^{j} \Big) \|\Delta \theta\|_{*}$$

$$\leq (\kappa_{1} + \rho_{k}) \|\Delta \theta\|_{*},$$

which completes the proof of (A.11).

Up to now, the proof of Lemma 7 is complete. Then, we prepare to prove Lemma 3 and Lemma 4.

#### Proof of Lemma 3.

Now, we are ready to prove Eq. (15). Firstly, note that

$$\begin{split} & \|\nabla_{\mathbf{w}}\ell(f(\mathbf{x}_{t}|\theta_{t-1}), y_{t}) - \nabla_{\mathbf{w}}\ell(f(\mathbf{x}_{t}|\theta'_{t-1}), y_{t})\|_{F} = \left\|\frac{\partial\ell(\hat{y}, y_{t})}{\partial\hat{y}}\nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta_{t-1}) - \frac{\partial\ell(\hat{y}', y_{t})}{\partial\hat{y}}\nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta'_{t-1})\right\|_{F} \\ \leq & \left\|\left(\frac{\partial\ell(\hat{y}, y_{t})}{\partial\hat{y}} - \frac{\partial\ell(\hat{y}', y_{t})}{\partial\hat{y}}\right)\nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta_{t-1}) + \frac{\partial\ell(\hat{y}', y_{t})}{\partial\hat{y}}\left(\nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta_{t-1}) - \nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta'_{t-1})\right)\right\|_{F} \\ \leq & \left|\frac{\partial\ell(\hat{y}, y_{t})}{\partial\hat{y}} - \frac{\partial\ell(\hat{y}', y_{t})}{\partial\hat{y}}\right| \cdot \|\nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta_{t-1})\|_{F} + \left|\frac{\partial\ell(\hat{y}', y_{t})}{\partial\hat{y}}\right| \cdot \|\nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta_{t-1}) - \nabla_{\mathbf{w}}f(\mathbf{x}|\theta'_{t-1})\|_{F} \\ \leq & \upsilon_{\ell}\left|f(\mathbf{x}_{t}|\theta_{t-1}) - f(\mathbf{x}|\theta'_{t-1})\right| \cdot \|\nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta_{t-1})\|_{F} + \alpha_{\ell}\|\nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta_{t-1}) - \nabla_{\mathbf{w}}f(\mathbf{x}_{t}|\theta'_{t-1})\|_{F}, \end{split}$$

where  $\hat{y} = f(\mathbf{x}_t | \theta_{t-1})$  and  $\hat{y}' = f(\mathbf{x}_t | \theta_{t-1}')$ . Then, according to (A.7), (A.9) and (A.10), we have

$$\|\nabla_{\mathbf{w}}\ell(f(\mathbf{x}_{t}|\theta_{t-1}), y_{t}) - \nabla_{\mathbf{w}}\ell(f(\mathbf{x}_{t}|\theta'_{t-1}), y_{t})\|_{F}$$

$$\leq \left\{ v_{\ell}B^{2K}\alpha_{\sigma}^{2K+2}C_{g}^{2K+2}C_{\mathbf{X}}^{2} + \alpha_{\ell}\left(v_{\sigma}B^{2K}\alpha_{\sigma}^{2K}C_{g}^{2K+2}C_{\mathbf{X}}^{2} + B^{K-1}\alpha_{\sigma}^{K+1}C_{g}^{K+1}C_{\mathbf{X}}\right)\right\} \|\triangle\theta_{t-1}\|_{*}.$$

This proves Eq. (15).

Similarly, for  $k = 1, 2, \dots, K$ ,

$$\begin{aligned} &\|\nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_t|\theta_{t-1}), y_t) - \nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_t|\theta_{t-1}'), y_t)\|_F \\ \leq & \upsilon_{\ell} |f(\mathbf{x}_t|\theta_{t-1}) - f(\mathbf{x}_t|\theta_{t-1}')| \cdot \|\nabla_{\mathbf{W}^{(k)}}f(\mathbf{x}_t|\theta_{t-1})\|_F + \alpha_{\ell} \|\nabla_{\mathbf{W}^{(k)}}f(\mathbf{x}_t|\theta_{t-1}) - \nabla_{\mathbf{W}^{(k)}}f(\mathbf{x}_t|\theta_{t-1}')\|_F. \end{aligned}$$

Then, according to (A.7), (A.9) and (A.11),

$$\begin{split} & \|\nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_{t}|\theta_{t-1}),y_{t}) - \nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_{t}|\theta_{t-1}'),y_{t})\|_{F} \\ \leq & \Big\{ \upsilon_{\ell}B^{2K}\alpha_{\sigma}^{2K+2}C_{g}^{2K+2}C_{\mathbf{X}}^{2} + \alpha_{\ell}\Big\{ \Big(\nu_{\sigma}B^{2K}\alpha_{\sigma}^{2K}C_{g}^{2K+2}C_{\mathbf{X}}^{2} \\ & + B^{K-1}\alpha_{\sigma}^{K+1}C_{g}^{K+1}C_{\mathbf{X}}\Big) + \nu_{\sigma}B^{K+k-1}\alpha_{\sigma}^{K+k-1}C_{g}^{K+k+1}C_{\mathbf{X}}^{2}\Big(\sum_{j=0}^{K-k}(B\alpha_{\sigma}C_{g})^{j}\Big) \Big\} \Big\} \|\triangle\theta_{t-1}\|_{*}, \end{split}$$

which competes the proof of Eq. (16).

# Proof of Lemma 4.

Since  $|\frac{\partial \ell(\hat{y},y)}{\partial \hat{y}}| \leq \alpha_{\ell}$  for any  $\hat{y}$  and y, we first have that for  $k=1,2,\ldots,K+1$ ,

$$\|\nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_t|\theta_{t-1}),y_t) - \nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_t'|\theta_{t-1}'),y_t')\|_F = \left\|\frac{\partial\ell(\hat{y},y_t)}{\partial\hat{y}}\nabla_{\mathbf{W}^{(k)}}f(\mathbf{x}_t|\theta_{t-1}) - \frac{\partial\ell(\hat{y}',y_t')}{\partial\hat{y}}\nabla_{\mathbf{W}^{(k)}}f(\mathbf{x}_t'|\theta_{t-1}')\right\|_F$$

$$\leq \alpha_{\ell}\left(\|\nabla_{\mathbf{W}^{(k)}}f(\mathbf{x}_t|\theta_{t-1})\|_F + \|\nabla_{\mathbf{W}^{(k)}}f(\mathbf{x}_t'|\theta_{t-1}')\|_F\right),$$

where  $\hat{y} = f(\mathbf{x}_t | \theta_{t-1})$  and  $\hat{y}' = f(\mathbf{x}_t' | \theta_{t-1}')$  and  $\mathbf{W}^{(K+1)} = \mathbf{w}$ . Finally, according to (A.7),

$$\|\nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_t|\theta_{t-1}), y_t) - \nabla_{\mathbf{W}^{(k)}}\ell(f(\mathbf{x}_t'|\theta_{t-1}', y_t'))\|_F \le \alpha_\ell \left(\|\nabla_{\mathbf{W}^{(k)}}f(\mathbf{x}|\theta_{t-1})\|_F + \|\nabla_{\mathbf{W}^{(k)}}f(\mathbf{x}|\theta_{t-1}')\|_F\right)$$

$$\le 2\alpha_\ell B^K \alpha_\sigma^{K+1} C_a^{K+1} C_{\mathbf{X}},$$

holds for k = 1, 2, ..., K + 1.

## **APPENDIX D:PROOF OF THEOREM 3**

Based on Lemma 3 and Lemma 4, we detail the proof of Theorem 3 as follows.

Note that  $(\mathbf{x}_t, y_t) = (\mathbf{x}_t', y_t')$  with probability  $1 - \frac{1}{m}$  and  $(\mathbf{x}_t, y_t) \neq (\mathbf{x}_t', y_t')$  with probability  $\frac{1}{m}$ . By considering Eq. (3) (in Section 3.2) and incorporating the probability of the two scenarios presented in Lemma 3 and Lemma 4, using  $\mathcal{F}$  and  $\mathcal{F}'$  to denote  $f(\mathbf{x}_t|\theta_{t-1})$  and  $f(\mathbf{x}_t|\theta_{t-1}')$ , respectively, we have:

$$\mathbb{E}_{\mathcal{A}}\left[\|\triangle\mathbf{w}_{t}\|_{2}\right] = (1 - \frac{1}{m})\mathbb{E}_{\mathcal{A}}\left[\|\triangle\mathbf{w}_{t-1} - \eta(\nabla_{\mathbf{w}}\ell(\mathcal{F}, y_{t}) - \nabla_{\mathbf{w}}\ell(\mathcal{F}', y_{t}))\|_{2}\right]$$

$$+ \frac{1}{m}\mathbb{E}_{\mathcal{A}}\left[\|\triangle\mathbf{w}_{t-1} - \eta(\nabla_{\mathbf{w}}\ell(\mathcal{F}, y_{t}) - \nabla_{\mathbf{w}}\ell(\mathcal{F}', y_{t}'))\|_{2}\right]$$

$$\leq (1 - \frac{1}{m})\mathbb{E}_{\mathcal{A}}\left[\|\triangle\mathbf{w}_{t-1}\|_{2} + \eta\|\nabla_{\mathbf{w}}\ell(\mathcal{F}, y_{t}) - \nabla_{\mathbf{w}}\ell(\mathcal{F}', y_{t})\|_{2}\right]$$

$$+ \frac{1}{m}\mathbb{E}_{\mathcal{A}}\left[\|\triangle\mathbf{w}_{t-1}\|_{2} + \eta\|\nabla_{\mathbf{w}}\ell(\mathcal{F}, y_{t}) - \nabla_{\mathbf{w}}\ell(\mathcal{F}', y_{t}')\|_{2}\right]$$

$$\leq (1 - \frac{1}{m})\mathbb{E}_{\mathcal{A}}\left[\|\triangle\mathbf{w}_{t-1}\|_{2} + \eta\|\nabla_{\mathbf{w}}\ell(\mathcal{F}, y_{t}) - \nabla_{\mathbf{w}}\ell(\mathcal{F}', y_{t})\|_{F}\right]$$

$$+ \frac{1}{m}\mathbb{E}_{\mathcal{A}}\left[\|\triangle\mathbf{w}_{t-1}\|_{2} + \eta\|\nabla_{\mathbf{w}}\ell(\mathcal{F}, y_{t}) - \nabla_{\mathbf{w}}\ell(\mathcal{F}', y_{t}')\|_{F}\right].$$

Based on Lemma 3 and Lemma 4,

$$\mathbb{E}_{\mathcal{A}}[\|\triangle \mathbf{w}_{t}\|_{2}] \leq \mathbb{E}_{\mathcal{A}}[\|\triangle \mathbf{w}_{t-1}\|_{2}] + \eta \kappa_{1} \mathbb{E}_{\mathcal{A}}[\|\triangle \theta_{t-1}\|_{*}] + \frac{2\eta \alpha_{\ell} B^{K} \alpha_{\sigma}^{K+1} C_{g}^{K+1} C_{\mathbf{X}}}{m}$$

Similarly, for  $k = 1, 2, \dots, K$ ,

$$\mathbb{E}_{\mathcal{A}}\left[\|\triangle\mathbf{W}_{t}^{(k)}\|_{2}\right] \leq \mathbb{E}_{\mathcal{A}}\left[\|\triangle\mathbf{W}_{t-1}^{(k)}\|_{2}\right] + \eta(\kappa_{1} + \rho_{k})\mathbb{E}_{\mathcal{A}}\left[\|\triangle\theta_{t-1}\|_{*}\right] + \frac{2\eta\alpha_{\ell}B^{K}\alpha_{\sigma}^{K+1}C_{g}^{K+1}C_{\mathbf{X}}}{m}.$$

Then,

$$\mathbb{E}_{\mathcal{A}}[\|\triangle\theta_{t}\|_{*}] = \mathbb{E}_{\mathcal{A}}[\|\triangle\mathbf{w}_{t}\|_{2}] + \sum_{k=1}^{K} \mathbb{E}_{\mathcal{A}}[\|\triangle\mathbf{W}_{t}^{(k)}\|_{2}]$$

$$\leq \mathbb{E}_{\mathcal{A}}[\|\triangle\mathbf{w}_{t-1}\|_{2}] + \eta\kappa_{1}\mathbb{E}_{\mathcal{A}}[\|\triangle\theta_{t-1}\|_{*}] + \frac{2\eta\alpha_{\ell}B^{K}\alpha_{\sigma}^{K+1}C_{g}^{K+1}C_{\mathbf{X}}}{m}$$

$$+ \sum_{k}^{K} \mathbb{E}_{\mathcal{A}}[\|\triangle\mathbf{W}_{t-1}^{(k)}\|_{2}] + \eta(\kappa_{1} + \rho_{k})\mathbb{E}_{\mathcal{A}}[\|\triangle\theta_{t-1}\|_{*}] + \frac{2\eta\alpha_{\ell}B^{K}\alpha_{\sigma}^{K+1}C_{g}^{K+1}C_{\mathbf{X}}}{m}$$

$$= (1 + (K+1)\eta\kappa_{1} + \eta\kappa_{2})\mathbb{E}_{\mathcal{A}}[\|\triangle\theta_{t-1}\|_{*}] + \frac{2(K+1)\eta\alpha_{\ell}B^{K}\alpha_{\sigma}^{K+1}C_{g}^{K+1}C_{\mathbf{X}}}{m}.$$

where  $\kappa_2 = \sum_{k=1}^K \rho_k$ . By (A.12), we have  $\kappa_2 = \nu_\sigma \left(B\alpha_\sigma C_g\right)^K C_g^2 C_{\mathbf{X}}^2 \left(\sum_{j=0}^{K-1} (j+1)(B\alpha_\sigma C_g)^j\right)$ , as defined in (8)). Finally, since  $\|\Delta\theta_0\|_* = \|\theta_0 - \theta_0'\|_* = 0$ 

$$\mathbb{E}_{\mathcal{A}}\big[\|\triangle\theta_T\|_*\big] \le \frac{c}{m} \sum_{t=1}^T \Big(1 + (K+1)\eta\kappa_1 + \eta\kappa_2\Big)^{t-1}.$$

This completes the proof of Theorem 3.

# **APPENDIX E: Proof for Section 6.1**

Recall to the GCNII,

$$\begin{cases} \mathbf{X}^{(k)} = \sigma\Big(\big((1-a_k)g(\mathbf{L})\mathbf{X}^{(k-1)} + a_k\mathbf{X}^{(0)}\big)\big((1-b_k)\mathbf{I}_d + b_k\mathbf{W}^{(k)}\big)\Big), & k = 1, 2, \dots, K; \\ f(\mathbf{x}|\theta) = \sigma\Big(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\big((1-a_{K+1})g(\mathbf{L})\mathbf{X}^{(K)} + a_{K+1}\mathbf{X}^{(0)}\big)\mathbf{w}\Big) \end{cases}$$

Proof of Eq. (19) and Eq. (20):

We first bound the output of each layer, i.e., bound  $\|\mathbf{X}^{(k)}\|_F$ . Applying  $\|\sigma(\mathbf{Z})\|_F \leq \alpha_{\sigma} \|\mathbf{Z}\|_F$  holds for any matrix  $\mathbf{Z}$  and  $\|\mathbf{A}_1\mathbf{A}_2\|_F \leq \|\mathbf{A}_1\|_F \|\mathbf{A}_2\|_2$ , we have

$$\|\mathbf{X}^{(k)}\|_{F} = \sigma \Big( ((1 - a_{k})g(\mathbf{L})\mathbf{X}^{(k-1)} + a_{k}\mathbf{X}^{(0)}) ((1 - b_{k})\mathbf{I}_{d} + b_{k}\mathbf{W}^{(k)}) \Big)$$

$$\leq \alpha_{\sigma} \| ((1 - a_{k})g(\mathbf{L})\mathbf{X}^{(k-1)} + a_{k}\mathbf{X}^{(0)}) ((1 - b_{k})\mathbf{I}_{d} + b_{k}\mathbf{W}^{(k)}) \|_{F}$$

$$\leq \alpha_{\sigma} \| (1 - a_{k})g(\mathbf{L})\mathbf{X}^{(k-1)} + a_{k}\mathbf{X}^{(0)} \|_{F} \cdot \| (1 - b_{k})\mathbf{I}_{d} + b_{k}\mathbf{W}^{(k)} \|_{2}.$$

Furthermore, since  $C_{\mathbf{X}} = \|\mathbf{X}\|_F = \|\mathbf{X}^{(0)}\|_F$ ,  $C_g = \|g(\mathbf{L})\|_2$ , and  $\|\mathbf{A}_1\mathbf{A}_2\|_F \le \|\mathbf{A}_1\|_2 \|\mathbf{A}_2\|_F$ ,

$$\|(1 - a_k)g(\mathbf{L})\mathbf{X}^{(k-1)} + a_k\mathbf{X}^{(0)}\|_F \le \|(1 - a_k)g(\mathbf{L})\mathbf{X}^{(k-1)}\|_F + \|a_k\mathbf{X}^{(0)}\|_F$$
  
$$\le (1 - a_k)C_q\|\mathbf{X}^{(k-1)}\|_F + a_kC_q,$$

and since  $\|\mathbf{W}^{(k)}\|_2 \le B$ ,  $\|(1-b_k)\mathbf{I}_d + b_k\mathbf{W}^{(k)}\|_2 \le 1 - b_k + b_kB$ . Therefore,

$$\|\mathbf{X}^{(k)}\|_{F} \leq \alpha_{\sigma} ((1 - a_{k})C_{g}\|\mathbf{X}^{(k-1)}\|_{F} + a_{k}C_{g})(1 - b_{k} + b_{k}B)$$
  
=  $(1 - a_{k})(1 - b_{k} + b_{k}B)\alpha_{\sigma}C_{g}\|\mathbf{X}^{(k-1)}\|_{F} + (1 - b_{k} + b_{k}B)a_{k}\alpha_{\sigma}C_{g}$ .

Note that  $\|\mathbf{X}^{(0)}\|_F = C_{\mathbf{X}}$ , we thus have that for  $k = 1, 2, \dots, K$ ,

$$\|\mathbf{X}^{(k)}\|_{F} \leq \Big(\prod_{i=1}^{k} (1-a_{i})(1-b_{i}+b_{i}B)\alpha_{\sigma}C_{g}\Big)C_{\mathbf{X}} + \sum_{j=1}^{k} \Big(\prod_{i=j+1}^{k} (1-a_{i})(1-b_{i}+b_{i}B)C_{g}\Big)\Big((1-b_{j}+b_{j}B)a_{j}\alpha_{\sigma}C_{g}\Big).$$
(A.20)

For convenience, in the following text we denote

$$B_{\mathbf{X}}^{(k)} := \Big(\prod_{i=1}^{k} (1 - a_i)(1 - b_i + b_i B)\alpha_{\sigma} C_g\Big) C_{\mathbf{X}} + \sum_{j=1}^{k} \Big(\prod_{i=j+1}^{k} (1 - a_i)(1 - b_i + b_i B) C_g\Big) \Big((1 - b_j + b_j B)a_j \alpha_{\sigma} C_g\Big), \quad (A.21)$$

and thus  $\|\mathbf{X}^{(k)}\|_F \leq B_{\mathbf{X}}^{(k)}$ . When  $a_k = 0, b_k = 1$  for all k, GCNII degenerates into the traditional GCN, and  $B_{\mathbf{X}}^{(k)} = B^k \alpha_{\sigma}^k C_g^k C_{\mathbf{X}}$ , which is the same as shown in (A.5). The bound of  $\|\mathbf{X}^{(k)}\|_F$  implies

$$\|\mathbf{H}^{(k)}\|_F = \|(1 - a_{k+1})g(\mathbf{L})\mathbf{X}^{(k)} + a_{k+1}\mathbf{X}^{(0)}\|_F \le (1 - a_{k+1})C_gB_{\mathbf{X}}^{(k)} + a_{k+1}C_{\mathbf{X}}.$$

Then, we bound the perturbation of the output of each layer, i.e., bound  $\|\triangle \mathbf{X}^{(k)}\|_F$ . Note that

$$\Delta \mathbf{X}^{(k)} = \mathbf{X}^{(k)} - \mathbf{X}^{(k)'} = \sigma \Big( \mathbf{H}^{(k-1)} \big( (1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)} \big) \Big) - \sigma \Big( \mathbf{H}^{(k-1)'} \big( (1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)'} \big) \Big).$$

Thus, following a calculation similar to Lemma 5, we have

$$\begin{split} \|\triangle\mathbf{X}^{(k)}\|_{F} &= \left\| \sigma \Big( \mathbf{H}^{(k-1)} \big( (1-b_{k}) \mathbf{I}_{d} + b_{k} \mathbf{W}^{(k)} \big) \Big) - \sigma \Big( \mathbf{H}^{(k-1)'} \big( (1-b_{k}) \mathbf{I}_{d} + b_{k} \mathbf{W}^{(k)'} \big) \Big) \right\|_{F} \\ &\leq \alpha_{\sigma} \|\mathbf{H}^{(k-1)} \big( (1-b_{k}) \mathbf{I}_{d} + b_{k} \mathbf{W}^{(k)} \big) - \mathbf{H}^{(k-1)'} \big( (1-b_{k}) \mathbf{I}_{d} + b_{k} \mathbf{W}^{(k)'} \big) \right\|_{F} \\ &= \alpha_{\sigma} \Big( \|\mathbf{H}^{(k-1)} - \mathbf{H}^{(k-1)'} \|_{F} \cdot \|(1-b_{k}) \mathbf{I}_{d} + b_{k} \mathbf{W}^{(k)} \|_{2} + \|\mathbf{H}^{(k-1)'} \|_{F} \cdot \|b_{k} (\mathbf{W}^{(k)} - \mathbf{W}^{(k)'}) \|_{2} \Big). \end{split}$$

Since  $\mathbf{H}^{(k-1)} - \mathbf{H}^{(k-1)'} = (1 - a_k) g(\mathbf{L}) \triangle \mathbf{X}^{(k-1)}$ ,  $\|\mathbf{H}^{(k-1)} - \mathbf{H}^{(k-1)'}\|_F = (1 - a_k) C_g \|\triangle \mathbf{X}^{(k-1)}\|_F$ . Combining  $\|\mathbf{H}^{(k-1)'}\|_F \le (1 - a_k) C_g \mathbf{X}^{(k-1)} + a_k C_{\mathbf{X}}$ ,  $\triangle \mathbf{W}^{(k)} = \mathbf{W}^{(k)} - \mathbf{W}^{(k)'}$ , and  $\|(1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)}\|_2 \le (1 - b_k) + b_k B$ , we have

$$\|\triangle \mathbf{X}^{(k)}\|_{F} \le c_{1}^{(k)} \|\triangle \mathbf{X}^{(k-1)}\|_{F} + c_{2}^{(k)} \|\triangle \mathbf{W}^{(k)}\|_{2}, \tag{A.22}$$

where  $c_1^{(k)} = (1 - a_k)(1 - b_k + b_k B)\alpha_{\sigma}C_g$  and  $c_2^{(k)} = \alpha_{\sigma}b_k((1 - a_k)C_gB_{\mathbf{X}}^{(k-1)} + a_kC_{\mathbf{X}})$ . This completes the proof of Eq. (19).

Furthermore, since  $\|\Delta \mathbf{X}^{(1)}\|_F \leq c_2^{(1)} \|\Delta \mathbf{W}^{(1)}\|_2$ , we further have

$$\|\triangle \mathbf{X}^{(k)}\|_F \le e^{(k)} \cdot (\sum_{j=1}^k \|\triangle \mathbf{W}^{(j)}\|_2),$$
 (A.23)

where  $e^{(k)} = \max\{c_1^{(k)}e^{(k-1)},c_2^{(k)}\}$  with  $e^{(0)} = 0$ . When  $a_k = 0, b_k = 1$  for all k, GCNII degenerates into the traditional GCN, we have  $c_1^{(k)} = B\alpha_\sigma C_g, c_2^{(k)} = B^{k-1}\alpha_\sigma^k C_g^k C_{\mathbf{X}}$ , and thus  $e^{(k)} = B^{k-1}\alpha_\sigma^k C_g^k C_{\mathbf{X}}$ , which is the same as shown in (A.8). This conclusively proves Eq. (20).

# Proof of Eq. (21):

To bound  $|f(\mathbf{x}|\theta) - f(\mathbf{x}|\theta')|$ , we apply the  $\alpha_{\sigma}$ -Lipschitz property of  $\sigma(\cdot)$  and then have

$$|f(\mathbf{x}|\theta) - f(\mathbf{x}|\theta')| = |\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w}) - \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}')| \le \alpha_{\sigma} \cdot |\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}'|,$$

that is, we need to bound  $|\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}'|$ 

Since 
$$\|\mathbf{A}_1\mathbf{A}_2 - \mathbf{A}_1'\mathbf{A}_2'\|_F \le \|\mathbf{A}_1 - \mathbf{A}_1'\|_F \|\mathbf{A}_2\|_2 + \|\mathbf{A}_1'\|_F \|\mathbf{A}_2 - \mathbf{A}_2'\|_2$$
,  
 $\|\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}'\| \le \|\boldsymbol{\delta}_{\mathbf{x}}^{\top}(\mathbf{H}^{(K)} - \mathbf{H}^{(K)'})\|_F \cdot \|\mathbf{w}\|_2 + \|\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\|_F \cdot \|\mathbf{w} - \mathbf{w}'\|_2$ .

Since  $\|\mathbf{H}^{(K)} - \mathbf{H}^{(K)'}\|_F \le (1 - a_{K+1})C_q \|\Delta \mathbf{X}^{(K)}\|_F$  and  $\|\mathbf{w}\|_2 \le B$ ,

$$\|\boldsymbol{\delta}_{\mathbf{x}}^{\mathsf{T}}(\mathbf{H}^{(K)} - \mathbf{H}^{(K)'})\|_{F} \cdot \|\mathbf{w}\|_{2} \leq (1 - a_{K+1})BC_{q}\|\Delta \mathbf{X}^{(K)}\|_{F},$$

and

$$\|\boldsymbol{\delta}_{\mathbf{x}}^{\mathsf{T}}\mathbf{H}^{(K)'}\|_{F} \cdot \|(\mathbf{w} - \mathbf{w}')\|_{2} \le ((1 - a_{K+1})C_{g}B_{\mathbf{X}}^{(K)} + a_{K+1}C_{\mathbf{X}})\|\Delta\mathbf{w}\|_{2},$$

which holds true because  $\|\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\|_{F} \leq \|\mathbf{H}^{(K)'}\|_{F} \leq (1 - a_{K+1})C_{g}B_{\mathbf{X}}^{(K)} + a_{K+1}C_{\mathbf{X}}$ . That is,

$$|\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}'| \le (1 - a_{K+1})BC_q\|\Delta\mathbf{X}^{(K)}\|_F + ((1 - a_{K+1})C_qB_{\mathbf{X}}^{(K)} + a_{K+1}C_{\mathbf{X}})\|\Delta\mathbf{w}\|_2.$$

By (A.23), we further have

$$|\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}'| \leq (1 - a_{K+1})BC_g \sum_{j=1}^{K} \|\Delta\mathbf{W}^{(j)}\|_2 + ((1 - a_{K+1})C_gB_{\mathbf{X}}^{(K)} + a_{K+1}C_{\mathbf{X}})\|\Delta\mathbf{w}\|_2$$
$$\leq \varrho \cdot \left(\sum_{j=1}^{K} \|\Delta\mathbf{W}^{(j)}\|_2 + \|\mathbf{w}\|_2\right) = \varrho \cdot \|\Delta\theta\|_*, \tag{A.24}$$

where  $\varrho = \max \left\{ (1 - a_{K+1}) B C_g \cdot e^{(K)}, (1 - a_{K+1}) C_g B_{\mathbf{X}}^{(K)} + a_{K+1} C_{\mathbf{X}} \right\}$ . Therefore,

$$|f(\mathbf{x}|\theta) - f(\mathbf{x}|\theta')| \le \alpha_{\sigma} \cdot |\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)'} \mathbf{w}'| \le \alpha_{\sigma} \varrho \cdot ||\Delta \theta||_{*}. \tag{A.25}$$

Note that when  $a_k=0, b_k=1$  for all k, we have  $e^{(K)}=B^{K-1}\alpha_\sigma^K C_g^K C_{\mathbf{X}}$  and  $B_{\mathbf{X}}^{(K)}=B^K\alpha_\sigma^K C_g^K C_{\mathbf{X}}$ , then at this point,  $\varrho=B^K\alpha_\sigma^K C_g^{K+1}C_{\mathbf{X}}$ , and  $|f(\mathbf{x}|\theta)-f(\mathbf{x}|\theta')|\leq \alpha_\sigma\varrho\cdot\|\triangle\theta\|_*=B^K\alpha_\sigma^{K+1}C_g^{K+1}C_{\mathbf{X}}\cdot\|\triangle\theta\|_*$ , which is consistent with (A.9). Thus, we complete the proof of Eq. (21).

# Proof of Eq. (22):

To bound the perturbation of the gradient, we first follow the calculation technique used in Appendix A to obtain the gradient of  $f(\mathbf{x}|\theta)$  as follow:

i) For the final layer,

$$\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) = \nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w}) [\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)}]^{\top}.$$

ii) For the hidden layer  $k = 1, 2, \dots, K$ ,

$$\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) = \sum_{i,j} \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}_{ij}^{(k)}} \cdot \frac{\partial \mathbf{X}_{ij}^{(k)}}{\partial \mathbf{W}^{(k)}} = b_k [\mathbf{H}^{(k-1)}]^\top (\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)}),$$

where  $\mathbf{R}^{(k)} = \nabla \sigma \left( \mathbf{H}^{(k-1)} \left( (1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)} \right) \right)$ . Furthermore,

$$\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k-1)}} = \sum_{i,j} \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}_{ij}^{(k)}} \cdot \frac{\partial \mathbf{X}_{ij}^{(k)}}{\partial \mathbf{X}^{(k-1)}} = (1 - a_k)[g(\mathbf{L})]^{\top} (\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)})[(1 - b_k)\mathbf{I}_d + b_k \mathbf{W}^{(k)}]^{\top},$$

with

$$\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} = (1 - a_{K+1}) \nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w}) [\boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L})]^{\top} \mathbf{w}^{\top}.$$

We now bound  $\|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta')\|_F$ . Note that  $\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) = \nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w}) [\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)}]^{\top}$ , we apply  $\|\mathbf{A}_1 \mathbf{A}_2 - \mathbf{A}_1' \mathbf{A}_2'\|_F \le \|\mathbf{A}_1 - \mathbf{A}_1'\|_F \cdot \|\mathbf{A}_2\|_F + \|\mathbf{A}_1'\|_F \cdot \|\mathbf{A}_1 - \mathbf{A}_2'\|_F$  and have

$$\begin{aligned} &\|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta')\|_{F} = \|\nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w})[\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}]^{\top} - \nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}')[\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}]^{\top}\|_{F} \\ \leq &|\nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w}) - \nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}')| \cdot \|[\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}]^{\top}\|_{F} + |\nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}')| \cdot \|[\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}]^{\top} - [\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}]^{\top}\|_{F}. \end{aligned}$$

We further apply the property of  $\sigma(\cdot)$  and have

$$\|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta')\|_{F} \leq \nu_{\sigma} \cdot \left|\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)'} \mathbf{w}'\right| \cdot \left\| \left[\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)}\right]^{\top} \right\|_{F} + \alpha_{\sigma} \cdot \left\| \left[\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)}\right]^{\top} - \left[\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)'}\right]^{\top} \right\|_{F} \\ \leq \nu_{\sigma} \cdot \left|\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)'} \mathbf{w}'\right| \cdot \left( (1 - a_{K+1}) C_{g} B_{\mathbf{X}}^{(K)} + a_{K+1} C_{\mathbf{X}} \right) + \alpha_{\sigma} \cdot (1 - a_{K+1}) C_{g} \|\Delta \mathbf{X}^{(K)}\|_{F}.$$

Finally, combining (A.23) and (A.24), we have

$$\|\nabla_{\mathbf{w}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{w}} f(\mathbf{x}|\theta')\|_{F} \leq \nu_{\sigma} \varrho \cdot \|\triangle \theta\|_{*} \cdot \left( (1 - a_{K+1}) C_{g} B_{\mathbf{X}}^{(K)} + a_{K+1} C_{\mathbf{X}} \right) + \alpha_{\sigma} \cdot (1 - a_{K+1}) C_{g} e^{(K)} \cdot \left( \sum_{j=1}^{K} \|\triangle \mathbf{W}^{(j)}\|_{2} \right)$$

$$\leq \left( \nu_{\sigma} \varrho \cdot \left( (1 - a_{K+1}) C_{g} B_{\mathbf{X}}^{(K)} + a_{K+1} C_{\mathbf{X}} \right) + \alpha_{\sigma} \cdot (1 - a_{K+1}) C_{g} e^{(K)} \right) \cdot \|\triangle \theta\|_{*}.$$

Thus, we complete the proof of Eq. (22).

Proof of bounding  $\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_F$ .

Next, we bound  $\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_F$ . First,

$$\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_{F} = \|b_{k}[\mathbf{H}^{(k-1)}]^{\top} (\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)}) - b_{k}[\mathbf{H}^{(k-1)'}]^{\top} (\frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)'})\|_{F}$$

$$\leq b_{k} (\|\mathbf{H}^{(k-1)} - \mathbf{H}^{(k-1)'}\|_{F} \cdot \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)}\|_{F} + \|\mathbf{H}^{(k-1)'}\|_{F} \cdot \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)'}\|_{F}).$$

Since  $\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)}\|_F \le \alpha_{\sigma} \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F$  and  $\|\mathbf{H}^{(k-1)} - \mathbf{H}^{(k-1)'}\|_F \le (1 - a_k)C_g \|\Delta \mathbf{X}^{(k-1)}\|_F$ 

$$\|\mathbf{H}^{(k-1)} - \mathbf{H}^{(k-1)'}\|_F \cdot \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)}\|_F \le (1 - a_k)C_g \|\Delta \mathbf{X}^{(k-1)}\|_F \cdot \alpha_\sigma \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F.$$

Following (A.14), we denote

$$\gamma_k := \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)'} \right\|_F. \tag{A.26}$$

Since  $\|\mathbf{H}^{(k-1)'}\|_F \le (1-a_k)C_g B_{\mathbf{X}}^{(k-1)} + a_k C_{\mathbf{X}}$ , we further apply (A.23) and (A.28) to obtain

$$\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_F$$

$$\leq b_k \left\{ (1 - a_k) \alpha_\sigma C_g e^{(k-1)} \cdot \left( \sum_{j=1}^{k-1} \| \Delta \mathbf{W}^{(j)} \|_2 \right) \cdot \| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \|_F + \left( (1 - a_k) C_g B_{\mathbf{X}}^{(k-1)} + a_k C_{\mathbf{X}} \right) \cdot \gamma_k \right\}. \tag{A.27}$$

That is, to bound  $\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_F$ , we need to bound  $\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F$  and  $\gamma_k$ . We provide the following steps to the bound of  $\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F$  and  $\gamma_k$ . We provide the following  $\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta)\|_F$  and  $\gamma_k$ . Using these two bounds, we finally obtain the upper bound of  $\|\nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta) - \nabla_{\mathbf{W}^{(k)}} f(\mathbf{x}|\theta')\|_F$  by applying them to (A.27).

Step 1: we first bound  $\left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k-1)}} \right\|_F$ . According to the iterative formula of  $\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k-1)}}$ , we have

$$\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k-1)}}\|_F = \|(1-a_k)[g(\mathbf{L})]^\top (\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)})[(1-b_k)\mathbf{I}_d + b_k \mathbf{W}^{(k)}]^\top \|_F$$

$$\leq (1-a_k)\|g(\mathbf{L})\|_2 \cdot \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot \mathbf{R}^{(k)}\|_F \cdot \|(1-b_k)\mathbf{I}_d + b_k \mathbf{W}^{(k)}\|_2.$$

Since the absolute value of the elements in  $\mathbf{R}^{(k)}$  is less than  $\alpha_{\sigma}$ ,  $\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\odot \mathbf{R}^{(k)}\|_F \leq \alpha_{\sigma}\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F$ . Then, combining  $\|g(\mathbf{L})\|_2 = C_g$  and  $\|(1-b_k)\mathbf{I}_d + b_k\mathbf{W}^{(k)}\|_2 \leq (1-b_k) + b_kB$ , we obtain the following iterative formula

$$\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{x}^{(k-1)}}\|_F \le (1 - a_k)(1 - b_k + b_k B)\alpha_{\sigma} C_g \cdot \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{x}^{(k)}}\|_F.$$

Note that since  $|\nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w})| \leq \alpha_{\sigma}$  and  $\|\mathbf{w}\|_{2} \leq B$ ,

$$\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}}\|_F = \|(1 - a_{K+1})\nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w})[\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})]^{\top}\mathbf{w}^{\top}\|_F \le (1 - a_{K+1})B\alpha_{\sigma}C_g.$$

Thus,

$$\left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \right\|_{F} \le \left( \prod_{j=k+1}^{K} (1 - a_{j})(1 - b_{j} + b_{j}B) \right) (1 - a_{K+1}) B \alpha_{\sigma}^{K+1-k} C_{g}^{K+1-k}.$$

For simplicity, we denote  $B_{\partial \mathbf{X}}^{(k)} = \Big(\prod_{j=k+1}^K (1-a_j)(1-b_j+b_jB)\Big)(1-a_{K+1})B\alpha_\sigma^{K+1-k}C_g^{K+1-k}$ , and then

$$\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}}\|_F \le B_{\partial \mathbf{X}}^{(k)}, \ k = 1, 2, \dots, K.$$
(A.28)

When  $a_k=0, b_k=1$  for all k, we have  $B_{\partial \mathbf{X}}^{(k)}=B^{K+1-k}\alpha_\sigma^{K+1-k}C_g^{K+1-k}$ , which is the same as shown in (A.6).

# Step 2: We next bound $\gamma_k$ .

Following the proof of Lemma 7, we have by (A.26) that

$$\gamma_{k} \leq \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot (\mathbf{R}^{(k)} - \mathbf{R}^{(k)'}) \right\|_{F} + \alpha_{\sigma} \cdot \left\| \left( \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \right) \right\|_{F}. \tag{A.29}$$

Similarly, let  $h_k := \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot (\mathbf{R}^{(k)} - \mathbf{R}^{(k)'}) \right\|_F$ . Then, applying  $\|\mathbf{A}_1 \odot \mathbf{A}_2\|_F \le \|\mathbf{A}_1\|_F \|\mathbf{A}_2\|_F$ , we have

$$h_k = \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \odot (\mathbf{R}^{(k)} - \mathbf{R}^{(k)'}) \right\|_F \le \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} \right\|_F \cdot \left\| \mathbf{R}^{(k)} - \mathbf{R}^{(k)'} \right\|_F. \tag{A.30}$$

Note that  $\mathbf{R}^{(k)} = \nabla \sigma \Big( \mathbf{H}^{(k-1)} \big( (1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)} \big) \Big)$ , then the  $\nu_{\sigma}$ -smooth property of  $\sigma(\cdot)$  implies

$$\begin{aligned} \left\| \mathbf{R}^{(k)} - \mathbf{R}^{(k)'} \right\|_F &= \left\| \nabla \sigma \left( \mathbf{H}^{(k-1)} \left( (1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)} \right) \right) - \nabla \sigma \left( \mathbf{H}^{(k-1)'} \left( (1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)'} \right) \right) \right\|_F \\ &\leq \nu_{\sigma} \cdot \left\| \mathbf{H}^{(k-1)} \left( (1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)} \right) - \mathbf{H}^{(k-1)'} \left( (1 - b_k) \mathbf{I}_d + b_k \mathbf{W}^{(k)'} \right) \right\|_F. \end{aligned}$$

Applying  $\|\mathbf{A}_1\mathbf{A}_2 - \mathbf{A}_1'\mathbf{A}_2'\|_F \le \|\mathbf{A}_1 - \mathbf{A}_1'\|_F \cdot \|\mathbf{A}_2\|_2 + \|\mathbf{A}_1'\|_F \cdot \|\mathbf{A}_2 - \mathbf{A}_2'\|_2$ , we further have

$$\|\mathbf{R}^{(k)} - \mathbf{R}^{(k)'}\|_{F} \le \nu_{\sigma} \cdot \left( \|(\mathbf{H}^{(k-1)} - \mathbf{H}^{(k-1)'}\|_{F} \cdot \|(1 - b_{k})\mathbf{I}_{d} + b_{k}\mathbf{W}^{(k)}\|_{2} + \|\mathbf{H}^{(k-1)'}\|_{F} \cdot \|b_{k}(\mathbf{W}^{(k)} - \mathbf{W}^{(k)'})\|_{2} \right).$$

Note that  $\|\mathbf{H}^{(k-1)} - \mathbf{H}^{(k-1)'}\|_F \le (1 - a_k)C_g\|\Delta \mathbf{X}^{(k-1)}\|_F \le (1 - a_k)C_g \cdot e^{(k-1)} \left(\sum_{j=1}^{k-1} \|\Delta \mathbf{W}^{(j)}\|_2\right)$  (see (A.23)),  $\|\mathbf{H}^{(k-1)'}\|_F \le (1 - a_k)C_gB_{\mathbf{X}}^{(k-1)} + a_kC_{\mathbf{X}}$ , and  $\|(1 - b_k)\mathbf{I}_d + b_k\mathbf{W}^{(k)}\|_2 \le (1 - b_k) + b_kB$ . Thus,

$$\|\mathbf{R}^{(k)} - \mathbf{R}^{(k)'}\|_{F} \le \nu_{\sigma} \cdot \left( (1 - a_{k})(1 - b_{k} + b_{k}B)C_{g} \cdot e^{(k-1)} \left( \sum_{j=1}^{k-1} \|\triangle \mathbf{W}^{(j)}\|_{2} \right) + \left( (1 - a_{k})C_{g}B_{\mathbf{X}}^{(k-1)} + a_{k}C_{\mathbf{X}} \right) b_{k} \|\triangle \mathbf{W}^{(k)}\|_{2} \right)$$

$$\leq \nu_{\sigma} r_k \cdot \left(\sum_{i=1}^k \|\triangle \mathbf{W}^{(j)}\|_2\right),\tag{A.31}$$

where  $r_k := \max\left\{(1-a_k)(1-b_k+b_kB)C_g\cdot e^{(k-1)}, \left((1-a_k)C_gB_{\mathbf{X}}^{(k-1)}+a_kC_{\mathbf{X}}\right)b_k\right\}$ . When  $a_k=0,b_k=1$  for all  $k, e^{(k-1)}=B^{k-2}\alpha_\sigma^{k-1}C_g^{k-1}C_{\mathbf{X}}$  and  $B_{\mathbf{X}}^{(k-1)}=B^{k-1}\alpha_\sigma^{k-1}C_g^{k-1}C_{\mathbf{X}}$ , then at this point,  $r_k=B^{k-1}\alpha_\sigma^{k-1}C_g^kC_{\mathbf{X}}$ , and thus  $\left\|\mathbf{R}^{(k)}-\mathbf{R}^{(k)'}\right\|_F \leq \nu_\sigma B^{k-1}\alpha_\sigma^{k-1}C_g^kC_{\mathbf{X}}\left(\sum_{i=1}^k \|\triangle\mathbf{W}^{(j)}\|_2\right)$ , which is the same as shown in (A.16).

Combining (A.28), (A.30), and (A.31), we have

$$h_k \le B_{\partial \mathbf{X}}^{(k)} \cdot \nu_{\sigma} r_k \cdot \left( \sum_{j=1}^k \| \Delta \mathbf{W}^{(j)} \|_2 \right). \tag{A.32}$$

Next, we use the same technique as in (A.18) that uses an  $h_{max}$  to bound all  $h_k$ . Specifically, let

$$h_{\max} = \max_{k=1,2,\ldots,K} \left\{ B_{\partial \mathbf{X}}^{(k)} \cdot \nu_{\sigma} r_k \right\} \cdot \| \triangle \theta \|_{*}.$$

Then,

$$h_k \le h_{\text{max}} \text{ holds for all } k = 1, 2, \dots, K.$$
 (A.33)

One can prove that when  $a_k=0,\ b_k=1$  for all k, then  $h_{\max}=\nu_{\sigma}B^K\alpha_{\sigma}^KC_g^{K+1}C_{\mathbf{X}}\|\triangle\theta\|_*$ , which is the same as in the case of traditional GCN.

Applying (A.33) to (A.29), we have

$$\gamma_k \le h_{\max} + \alpha_{\sigma} \cdot \left\| \left( \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \right) \right\|_F, \tag{A.34}$$

and we can derive the iterative formula for the bound of  $\gamma_k$ . To do this, we utilize the iterative formula of  $\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k-1)}}$  and obtain

$$\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}}\|_{F} \leq \|(1 - a_{k+1})[g(\mathbf{L})]^{\top} (\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)})[(1 - b_{k+1})\mathbf{I}_{d} + b_{k+1}\mathbf{W}^{(k+1)}]^{\top} \\ - (1 - a_{k+1})[g(\mathbf{L})]^{\top} (\frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)'})[(1 - b_{k+1})\mathbf{I}_{d} + b_{k+1}\mathbf{W}^{(k+1)'}]^{\top}\|_{F} \\ \leq (1 - a_{k+1})C_{g} \cdot \left( \|(\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)})[(1 - b_{k+1})\mathbf{I}_{d} + b_{k+1}\mathbf{W}^{(k+1)}]^{\top} \\ - (\frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)'})[(1 - b_{k+1})\mathbf{I}_{d} + b_{k+1}\mathbf{W}^{(k+1)'}]^{\top}\|_{F} \right).$$

Applying  $\|\mathbf{A}_1\mathbf{A}_2 - \mathbf{A}_1'\mathbf{A}_2'\|_F \le \|\mathbf{A}_1 - \mathbf{A}_1'\|_F \|\mathbf{A}_2\|_2 + \|\mathbf{A}_1'\|_F \|\mathbf{A}_2 - \mathbf{A}_2'\|_2$  and  $\|(1 - b_{k+1})\mathbf{I}_d + b_{k+1}\mathbf{W}^{(k+1)}\|_2 \le 1 - b_{k+1} + b_{k+1}B$ , we have

$$\left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \right\|_{F} \le (1 - a_{k+1}) C_{g} \cdot \left( (1 - b_{k+1} + b_{k+1}B) \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)'} \right\|_{F} + \left\| \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k+1)}} \odot \mathbf{R}^{(k+1)'} \right\|_{F} \cdot \left\| b_{k+1} (\mathbf{W}^{(k+1)} - \mathbf{W}^{(k+1)'}) \right\|_{2} \right).$$

Since  $\left\|\frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k+1)}}\odot \mathbf{R}^{(k+1)'}\right\|_F \leq \alpha_\sigma \left\|\frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k+1)}}\right\|_F \leq \alpha_\sigma B_{\partial \mathbf{X}}^{(k+1)} \text{ and } \gamma_{k+1} = \left\|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k+1)}}\odot \mathbf{R}^{(k+1)} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k+1)}}\odot \mathbf{R}^{(k+1)'}\right\|_F$ , thus

$$\left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(k)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(k)}} \right\|_{F} \le (1 - a_{k+1}) C_g \cdot \left( (1 - b_{k+1} + b_{k+1}B) \gamma_{k+1} + b_{k+1} \alpha_{\sigma} B_{\partial \mathbf{X}}^{(k+1)} \cdot \left\| \triangle \mathbf{W}^{(k+1)} \right\|_{2} \right). \tag{A.35}$$

Combining (A.34) and (A.35), we obtain the iterative formula for the bound of  $\gamma_k$  as

$$\gamma_{k} \le h_{\max} + (1 - a_{k+1})\alpha_{\sigma}C_{g} \cdot \left( (1 - b_{k+1} + b_{k+1}B) \cdot \gamma_{k+1} + b_{k+1}\alpha_{\sigma} \cdot B_{\partial \mathbf{X}}^{(k+1)} \cdot \|\Delta \mathbf{W}^{(k+1)}\|_{2} \right). \tag{A.36}$$

Furthermore, since  $\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} = (1 - a_{K+1}) \nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w}) [\boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L})]^{\top} \mathbf{w}^{\top}$ ,

$$\begin{split} \|\frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(K)}}\|_F = & \|(1 - a_{K+1})\nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w})[\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})]^{\top}\mathbf{w}^{\top} - (1 - a_{K+1})\nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}')[\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})]^{\top}\mathbf{w}'^{\top}\|_F \\ = & (1 - a_{K+1})\|\nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)}\mathbf{w})[\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})]^{\top}\mathbf{w}^{\top} - \nabla\sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top}\mathbf{H}^{(K)'}\mathbf{w}')[\boldsymbol{\delta}_{\mathbf{x}}^{\top}g(\mathbf{L})]^{\top}\mathbf{w}'^{\top}\|_F. \end{split}$$

The inequation  $\|\mathbf{A}_1\mathbf{A}_2 - \mathbf{A}_1'\mathbf{A}_2'\|_F \le \|\mathbf{A}_1 - \mathbf{A}_1'\|_F \|\mathbf{A}_2\|_2 + \|\mathbf{A}_1'\|_F \|\mathbf{A}_2 - \mathbf{A}_2'\|_2$  further derives

$$\begin{split} \left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(K)}} \right\|_{F} &\leq \left( 1 - a_{K+1} \right) \left( \left| \nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w}) - \nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)'} \mathbf{w}') \right| \cdot \left\| [\boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L})]^{\top} \mathbf{w}^{\top} \right\|_{2} \\ &+ \left| \nabla \sigma(\boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)'} \mathbf{w}') \right| \cdot \left\| [\boldsymbol{\delta}_{\mathbf{x}}^{\top} g(\mathbf{L})]^{\top} (\mathbf{w} - \mathbf{w}')^{\top} \right\|_{2} \right) \\ &\leq \left( 1 - a_{K+1} \right) \left( \nu_{\sigma} \cdot \left| \boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)} \mathbf{w} - \boldsymbol{\delta}_{\mathbf{x}}^{\top} \mathbf{H}^{(K)'} \mathbf{w}' \right| \cdot BC_{g} + \alpha_{\sigma} C_{g} \| \Delta \mathbf{w} \|_{2} \right), \end{split}$$

where the last inequation holds true because of the  $\nu_{\sigma}$ -smooth property of  $\sigma(\cdot)$  and  $|\nabla \sigma| \leq \alpha_{\sigma}$ . Then, we apply (A.24) to obtain

$$\left\| \frac{\partial f(\mathbf{x}|\theta)}{\partial \mathbf{X}^{(K)}} - \frac{\partial f(\mathbf{x}|\theta')}{\partial \mathbf{X}^{(K)}} \right\|_{F} \le (1 - a_{K+1}) \left( \nu_{\sigma} B C_{g} \cdot \varrho \| \triangle \theta \|_{*} + \alpha_{\sigma} C_{g} \| \triangle \mathbf{w} \|_{2} \right), \tag{A.37}$$

where  $\varrho = \max\left\{(1 - a_{K+1})BC_g \cdot e^{(K)}, \ (1 - a_{K+1})C_gB_{\mathbf{X}}^{(K)} + a_{K+1}C_{\mathbf{X}}\right\}$ . Substituting (A.37) into (A.34),

$$\gamma_K \le h_{\max} + (1 - a_{K+1})\alpha_{\sigma} \cdot \left(\nu_{\sigma} B C_g \cdot \varrho \|\triangle \theta\|_* + \alpha_{\sigma} C_g \|\triangle \mathbf{w}\|_2\right). \tag{A.38}$$

Combining (A.36) and (A.38), we can further obtain the bound of  $\gamma_k$ .

# **APPENDIX F:Proof for Section 6.2**

Proof of Eq. (23):

For vectors  $\mathbf{z} = (z_1, z_2, \dots, z_p)$  and  $\mathbf{z}' = (z_1', z_2', \dots, z_p')$  (with  $\|\mathbf{z} - \mathbf{z}'\|_{\infty} < 1$ ), the softmax function is defined as:

$$\operatorname{softmax}(\mathbf{z}) = (Z_1, Z_2, \dots, Z_p), \quad \operatorname{softmax}(\mathbf{z}') = (Z_1', Z_2', \dots, Z_p'),$$

where

$$Z_k = \frac{e^{z_k}}{\sum_{j=1}^p e^{z_j}}, \quad Z'_k = \frac{e^{z'_k}}{\sum_{j=1}^p e^{z'_j}} \quad \forall k = 1, 2, \dots, p.$$

For each k, rewrite  $|Z_k - Z'_k|$  using the softmax definition:

$$|Z_k - Z'_k| = \left| \frac{e^{z_k}}{S} - \frac{e^{z'_k}}{S'} \right| = \left| \frac{e^{z_k}S' - e^{z'_k}S}{SS'} \right|,$$

where  $S = \sum_{j=1}^{p} e^{z_j}$  and  $S' = \sum_{j=1}^{p} e^{z_j'}$ . By the triangle inequality in the numerator:

$$|e^{z_k}S' - e^{z'_k}S| \le e^{z_k}|S' - S| + S|e^{z_k} - e^{z'_k}|.$$

Summing  $|Z_k - Z'_k|$  over k gives the 1-norm:

$$\|\operatorname{softmax}(\mathbf{z}) - \operatorname{softmax}(\mathbf{z}')\|_1 = \sum_{k=1}^p |Z_k - Z_k'| \le \sum_{k=1}^p \left(\frac{e^{z_k}}{SS'}|S' - S| + \frac{1}{S'}|e^{z_k} - e^{z_k'}|\right).$$

Since  $\sum_{k=1}^{p} \frac{e^{z_k}}{S} = 1$  , this simplifies to:

$$\|\operatorname{softmax}(\mathbf{z}) - \operatorname{softmax}(\mathbf{z}')\|_{1} \le \frac{|S' - S|}{S'} + \frac{1}{S'} \sum_{k=1}^{p} |e^{z_{k}} - e^{z'_{k}}|.$$
 (A.39)

Notice that for any k, by the mean value theorem,

$$|e^{z_k} - e^{z'_k}| \le e^{z'_k} |e^{z_k - z'_k} - 1| \le e^{z'_k} \cdot e \cdot |z_k - z'_k|,$$

where the last inequation holds true because  $\|\mathbf{z} - \mathbf{z}'\|_{\infty} < 1$ , and, by the triangle inequality,

$$|S' - S| = \left| \sum_{j=1}^{p} (e^{z'_j} - e^{z_j}) \right| \le \sum_{j=1}^{p} |e^{z'_j} - e^{z_j}|.$$

Substituting into (A.39) gives:

$$\|\operatorname{softmax}(\mathbf{z}) - \operatorname{softmax}(\mathbf{z}')\|_1 \le \frac{2}{S'} \sum_{j=1}^p |e^{z'_j} - e^{z_j}| \le 2e \cdot \max |z_k - z'_k| = 2e \cdot \|\mathbf{z} - \mathbf{z}'\|_{\infty},$$

where the second inequation holds true because  $\sum_{k=1}^{p} \frac{e^{z'_k}}{S'} = 1$ . Thus, we complete the proof of Eq. (23).

# Proof of Eq. (24):

Recall that we denote B a constant which bounds all original and perturbed parameters, i.e,

$$\|\mathbf{W}_{K}\|_{2}, \|\mathbf{W}_{K}'\|_{2}, \|\mathbf{W}_{Q}\|_{2}, \|\mathbf{W}_{Q}'\|_{2}, \|\mathbf{W}_{V}\|_{2}, \|\mathbf{W}_{V}'\|_{2}, \|\mathbf{W}_{O}\|_{2}, \|\mathbf{W}_{O}'\|_{2} \le B,$$

and  $\|\mathbf{a}\|_2$ ,  $\|\mathbf{a}'\|_2 \leq B$  (output vector norm). And for  $\theta = \{\mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V, \mathbf{W}_O, \mathbf{a}\}$ ,

$$\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) = \mathbf{a}^T \mathbb{I}_{\geq 0}(\mathbf{Z}_\theta) \cdot \mathbf{W}_O \cdot \mathbf{W}_V \cdot \mathbf{M}_\theta \cdot (\mathbf{W}_O \mathbf{x}_n)^\top, \tag{A.40}$$

where:  $\mathbf{Z}_{\theta} = \mathbf{W}_O \cdot \operatorname{Attn}(\mathbf{x}_n; \theta)$ ,  $\operatorname{Attn}(\mathbf{x}_n; \theta) = \sum_{s \in \mathcal{T}^n} A_{s,n}(\theta) (\mathbf{W}_V \mathbf{x}_s)$ ,  $A_{s,n}(\theta) = \operatorname{softmax}(S_{\cdot,n}(\theta))_s$  with  $S_{s,n}(\theta) = (\mathbf{W}_K \mathbf{x}_s)^\top (\mathbf{W}_Q \mathbf{x}_n)$ ,  $\mathbf{M}_{\theta} = \sum_{s \in \mathcal{T}^n} A_{s,n}(\theta) (\mathbf{x}_s - \bar{\mathbf{x}}_n) \mathbf{x}_s^\top$  (aggregate neighbor term). Then the gradient perturbation  $\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta')$  arises from differences in  $\theta$  and  $\theta'$ . According to (A.40), we

Then the gradient perturbation  $\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta')$  arises from differences in  $\theta$  and  $\theta'$ . According to (A.40), we apply the triangle inequality and Lipschitz continuity of matrix multiplication/activation functions and then decompose the perturbation into contributions from each parameter:

$$\|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta')\|_2 \le \sum_{\phi \in \{\theta\}} \|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta_{\phi \to \phi'})\|_2$$
(A.41)

where  $\theta_{\phi \to \phi'}$  replaces parameter  $\phi$  with  $\phi'$  while keeping others fixed.

• Contribution from  $\mathbf{a} - \mathbf{a}'$ : The term  $\mathbf{a}^{\top}$  in (A.41) introduces a perturbation bounded by:

$$\|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n|\theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n|\theta_{\mathbf{a} \to \mathbf{a}'})\|_2 \le \|\mathbf{a} - \mathbf{a}'\|_2 \cdot \|\mathbb{I}_{\ge 0}(\mathbf{Z}_\theta)\|_2 \cdot \|\mathbf{W}_O\|_2 \cdot \|\mathbf{W}_V\|_2 \cdot \|\mathbf{M}_\theta\|_2 \cdot \|\mathbf{W}_Q\mathbf{x}_n\|_2,$$

Using

$$\|\mathbb{I}_{>0}(\mathbf{Z}_{\theta})\|_{2} \le 1, \quad \|\mathbf{M}_{\theta}\|_{2} \le K_{\max}C_{\mathbf{x}}^{2}, \quad \|\mathbf{W}_{Q}\mathbf{x}_{n}\|_{2} \le BC_{\mathbf{x}},$$
 (A.42)

where  $K_{\text{max}}$  is the maximum neighborhood size, we have

$$\|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta_{\mathbf{a} \to \mathbf{a}'})\|_2 \le \|\mathbf{a} - \mathbf{a}'\|_2 \cdot 1 \cdot B \cdot B \cdot K_{\max} C_{\mathbf{x}}^2 \cdot BC_{\mathbf{x}} = \|\mathbf{a} - \mathbf{a}'\|_2 \cdot K_{\max} B^3 C_{\mathbf{x}}^3$$

• Contribution from  $\mathbf{W}_O - \mathbf{W}_O'$ : The term  $\mathbf{W}_O$  affects both  $\mathbf{Z}_{\theta}$  and the gradient product. By Lipschitz continuity of matrix multiplication:

$$\|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n|\theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n|\theta_{\mathbf{W}_O \to \mathbf{W}_O'})\|_2 \le \|\mathbf{W}_O - \mathbf{W}_O'\|_2 \cdot \|\mathbf{a}\|_2 \cdot \|\mathbf{I}_{\geq 0}(\mathbf{Z}_\theta)\|_2 \cdot \|\mathbf{W}_V\|_2 \cdot \|\mathbf{M}_\theta\|_2 \cdot \|\mathbf{W}_Q\mathbf{x}_n\|_2.$$
Substituting (A.42):

$$\|\nabla_{\mathbf{W}_K}F(\mathbf{x}_n|\theta) - \nabla_{\mathbf{W}_K}F(\mathbf{x}_n|\theta_{\mathbf{W}_O \to \mathbf{W}_O'})\|_2 \leq \|\mathbf{W}_O - \mathbf{W}_O'\| \cdot B \cdot 1 \cdot B \cdot K_{\max}C_{\mathbf{x}}^2 \cdot BC_{\mathbf{x}} = \|\mathbf{W}_O - \mathbf{W}_O'\|_2 \cdot K_{\max}B^3C_{\mathbf{x}}^3.$$

Contribution from  $\mathbf{W}_V - \mathbf{W}_V'$ :  $\mathbf{W}_V$  is independent of  $\mathrm{Attn}(\mathbf{x}_n)$  and  $\mathbf{M}_{\theta}$ . The perturbation bound is:

$$\|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n|\theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n|\theta_{\mathbf{W}_V \to \mathbf{W}_V'})\|_2 \le \|\mathbf{W}_V - \mathbf{W}_V'\|_2 \cdot \|\mathbf{a}\|_2 \cdot \|\mathbf{I}_{\ge 0}(\mathbf{Z}_\theta)\|_2 \cdot \|\mathbf{W}_O\|_2 \cdot \|\mathbf{M}_\theta\|_2 \cdot \|\mathbf{W}_Q \mathbf{x}_n\|_2.$$

By symmetry with  $\mathbf{W}_{O}$ :

$$\|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta_{\mathbf{W}_V \to \mathbf{W}_V'})\|_2 \le \|\mathbf{W}_V - \mathbf{W}_V'\|_2 \cdot K_{\max} B^3 C_{\mathbf{x}}^3$$

Contribution from  $\mathbf{W}_Q - \mathbf{W}_Q'$ :  $\mathbf{W}_Q$  affects attention scores  $S_{s,n}$  and thus  $A_{s,n}$  and  $\mathbf{M}_{\theta}$ . Using Lipschitzness of softmax and matrix multiplication:

$$\|\nabla_{\mathbf{W}_K}F(\mathbf{x}_n|\theta) - \nabla_{\mathbf{W}_K}F(\mathbf{x}_n|\theta_{\mathbf{W}_Q \to \mathbf{W}_Q'})\|_2 \leq 2e\|\mathbf{W}_Q - \mathbf{W}_Q'\|_2 \cdot \|\mathbf{a}\|_2 \cdot \|\mathbf{I}_{\geq 0}(\mathbf{Z}_\theta)\|_2 \cdot \|\mathbf{W}_O\|_2 \cdot \|\mathbf{W}_V\|_2 \cdot \|\mathbf{M}_\theta\|_2 \cdot C_{\mathbf{x}}.$$
 Substituting bounds:

$$\|\nabla_{\mathbf{W}_K}F(\mathbf{x}_n|\theta) - \nabla_{\mathbf{W}_K}F(\mathbf{x}_n|\theta_{\mathbf{W}_Q \to \mathbf{W}_Q'})\|_2 \leq 2e\|\mathbf{W}_Q - \mathbf{W}_Q'\|_2 \cdot K_{\max}B^3C_{\mathbf{x}}^3.$$

Contribution from  $\mathbf{W}_K - \mathbf{W}_K'$ :  $\mathbf{W}_K$  directly impacts  $S_{s,n}$ ,  $A_{s,n}$ , and  $\mathbf{M}_{\theta}$ . By analogous reasoning:

$$\|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta_{\mathbf{W}_K \to \mathbf{W}_K'}) \|_2 \le 2e \|\mathbf{W}_K - \mathbf{W}_K' \|_2 \cdot K_{\max} B^3 C_{\mathbf{x}}^3.$$

According to (A.41), total gradient perturbation bound summing all contributions, and we get:

$$\|\nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta) - \nabla_{\mathbf{W}_K} F(\mathbf{x}_n | \theta')\|_2 \le 2eK_{\max} B^3 C_{\mathbf{x}}^3 \|\triangle \theta\|_*,$$

where 
$$\|\Delta\theta\|_* = \|\mathbf{W}_K - \mathbf{W}_K'\|_2 + \|\mathbf{W}_V - \mathbf{W}_V'\|_2 + \|\mathbf{W}_O - \mathbf{W}_O'\|_2 + \|\mathbf{W}_Q - \mathbf{W}_Q'\|_2 + \|\mathbf{a} - \mathbf{a}'\|_2$$
.

# REFERENCES

- Y. Liang, F. Meng, Y. Zhang, Y. Chen, J. Xu, and J. Zhou, "Emotional conversation generation with heterogeneous graph neural network," Artificial Intelligence, vol. 308, p. 103714, 2022.
- Y. Ma and J. Tang, Deep learning on graphs. Cambridge University Press, 2021.
- M. Gori, G. Monfardini, and F. Scarselli, "A new model for learning in graph domains," in Proceedings of the IEEE International Joint Conference on Neural Networks. IEEE, 2005, pp. 729-734.
- F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," IEEE Transactions on Neural Networks, vol. 20, no. 1, pp. 61-80, 2008.
- A. Micheli, "Neural network for graphs: A contextual constructive approach," IEEE Transactions on Neural Networks, vol. 20, no. 3, pp. 498-511,
- K. Yao, J. Liang, J. Liang, M. Li, and F. Cao, "Multi-view graph convolutional networks with attention mechanism," Artificial Intelligence, vol. 307, p. 103708, 2022.
- W. L. Hamilton, Graph representation learning. Morgan & Claypool, 2020.
- L. Wu, P. Cui, J. Pei, and L. Zhao, Graph Neural Networks: Foundations, Frontiers, and Applications. Springer, 2022.
- F. M. Bianchi, D. Grattarola, L. Livi, and C. Alippi, "Graph neural networks with convolutional arma filters," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 44, no. 7, pp. 3496–3507, 2021.

  B. Jiang, B. Wang, S. Chen, J. Tang, and B. Luo, "Graph neural network meets sparse representation: Graph sparse neural networks via
- exclusive group lasso," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 45, no. 10, pp. 12692-12698, 2023.
- H. Zhang, Y. Zhu, and X. Li, "Decouple graph neural networks: Train multiple simple gnns simultaneously instead of one," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 46, no. 11, pp. 7451–7462, 2024.
- [12] D. Bacciu, F. Errica, A. Micheli, and M. Podda, "A gentle introduction to deep learning for graphs," Neural Networks, vol. 129, pp. 203–221,
- [13] J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun, "Graph neural networks: A review of methods and applications," AI Open, vol. 1, pp. 57–81, 2020.
- [14] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu, "A comprehensive survey on graph neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 1, pp. 4–24, 2021.
   [15] Z. Zhang, P. Cui, and W. Zhu, "Deep learning on graphs: A survey," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 1, pp.
- 249-270, 2022.
- [16] Q. Li, Z. Han, and X.-M. Wu, "Deeper insights into graph convolutional networks for semi-supervised learning," in Proceedings of the 32nd AAAI Conference on Artificial Intelligence, 2018, pp. 3538–3545.
- [17] L. Zhao and L. Akoglu, "PairNorm: Tackling oversmoothing in GNNs," in International Conference on Learning Representations, 2020.
- [18] K. Oono and T. Suzuki, "Graph neural networks exponentially lose expressive power for node classification," in International Conference on Learning Representations, 2020.
- [19] Y. Rong, W. Huang, T. Xu, and J. Huang, "DropEdge: Towards deep graph convolutional networks on node classification," in International Conference on Learning Representations, 2020.
- [20] H. Yuan, J. Tang, X. Hu, and S. Ji, "XGNN: Towards model-level explanations of graph neural networks," in Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2020, pp. 430–438.
- [21] H. Yuan, H. Yu, J. Wang, K. Li, and S. Ji, "On explainability of graph neural networks via subgraph explorations," in Proceedings of the 38th
- International Conference on Machine Learning, 2021, pp. 12241–12252.
  [22] H. Yuan, H. Yu, S. Gui, and S. Ji, "Explainability in graph neural networks: A taxonomic survey," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 45, no. 5, pp. 5782–5799, 2022.

  T. Schnake, O. Eberle, J. Lederer, S. Nakajima, K. T. Schütt, K.-R. Müller, and G. Montavon, "Higher-order explanations of graph neural
- networks via relevant walks," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 44, no. 11, pp. 7581–7596, 2021.

- [24] G. Bouritsas, F. Frasca, S. Zafeiriou, and M. M. Bronstein, "Improving graph neural network expressivity via subgraph isomorphism counting," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 45, no. 1, pp. 657-668, 2022.
- [25] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" in International Conference on Learning Representations,
- [26] Z. Chen, S. Villar, L. Chen, and J. Bruna, "On the equivalence between graph isomorphism testing and function approximation with GNNs," in Proceedings of the 33rd International Conference on Neural Information Processing Systems, 2019, pp. 15868–15876.
- [27] N. Dehmamy, A.-L. Barabási, and R. Yu, "Understanding the representation power of graph neural networks in learning graph topology," in Proceedings of the 33rd International Conference on Neural Information Processing Systems, 2019, pp. 15413-15423.
- [28] S. S. Du, K. Hou, R. R. Salakhutdinov, B. Poczos, R. Wang, and K. Xu, "Graph neural tangent kernel: Fusing graph neural networks with graph kernels," in Proceedings of the 33rd International Conference on Neural Information Processing Systems, 2019, pp. 5723-5733.
- [29] S. Zhang, M. Wang, S. Liu, P.-Y. Chen, and J. Xiong, "Fast learning of graph neural networks with guaranteed generalizability: one-hidden-layer case," in *Proceedings of the 37th International Conference on Machine Learning*, 2020, pp. 11268–11277.
- [30] F. Scarselli, A. C. Tsoi, and M. Hagenbuchner, "The Vapnik-Chervonenkis dimension of graph and recursive neural networks," Neural Networks, vol. 108, pp. 248-259, 2018.
- [31] V. Garg, S. Jegelka, and T. Jaakkola, "Generalization and representational limits of graph neural networks," in Proceedings of the 37 th International Conference on Machine Learning, 2020, pp. 3419–3430.
- [32] K. Oono and T. Suzuki, "Optimization and generalization analysis of transduction through gradient boosting and application to multi-scale graph neural networks," in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020, pp. 18917–18930.
- Š. Lv, "Generalization bounds for graph convolutional neural networks via Rademacher complexity," arXiv preprint arXiv:2102.10234, 2021.
- [34] P. Esser, L. Chennuru Vankadara, and D. Ghoshdastidar, "Learning theory can (sometimes) explain generalisation in graph neural networks," in Proceedings of the 35th International Conference on Neural Information Processing Systems, 2021, pp. 27 043–27 056.
- [35] H. Tang and Y. Liu, "Towards understanding the generalization of graph neural networks," in Proceedings of the 40th International Conference on Machine Learning, 2023, pp. 33 674-33 719.
- S. Verma and Z.-L. Zhang, "Stability and generalization of graph convolutional neural networks," in Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2019, pp. 1539-1548.
- [37] M. K. Ng and A. Yip, "Stability and generalization of graph convolutional networks in eigen-domains," Analysis and Applications, vol. 21, no. 03, pp. 819-840, 2023.
- [38] R. Liao, R. Urtasun, and R. Zemel, "A PAC-Bayesian approach to generalization bounds for graph neural networks," in International Conference on Learning Representations, 2021.
- [39] H. Ju, D. Li, A. Sharma, and H. R. Zhang, "Generalization in graph neural networks: Improved PAC-Bayesian bounds on graph diffusion," in Proceedings of the 26th International Conference on Artificial Intelligence and Statistics, 2023, pp. 6314–6341.
- A. Jacot, F. Gabriel, and C. Hongler, "Neural tangent kernel: Convergence and generalization in neural networks," in *Proceedings of the 32nd* International Conference on Neural Information Processing Systems, 2018, pp. 8580–8589.
- [41] S. S. Du, K. Hou, R. R. Salakhutdinov, B. Poczos, R. Wang, and K. Xu, "Graph neural tangent kernel: Fusing graph neural networks with graph kernels," in Proceedings of the 33rd International Conference on Neural Information Processing Systems, 2019, pp. 5723–5733.
- [42] M. Chen, Z. Wei, Z. Huang, B. Ding, and Y. Li, "Simple and deep graph convolutional networks," in Proceedings of the 37th International Conference on Machine Learning, 2020, pp. 1725-1735.
- [43] H. Li, M. Wang, T. Ma, S. Liu, Z. Zhang, and P.-Y. Chen, "What improves the generalization of graph transformers? a theoretical dive into the self-attention and positional encoding," in Proceedings of the 41th International Conference on Machine Learning, 2024, pp. 28784–28829.
- S. Liu, L. Wei, S. Lv, and M. Li, "Stability and generalization of  $\ell_p$ -regularized stochastic learning for GCN," in *Proceedings of the 32nd* International Joint Conference on Artificial Intelligence, 2023, pp. 5685-5693.
- [45] C. Huang, M. Li, F. Cao, H. Fujita, Z. Li, and X. Wu, "Are graph convolutional networks with random weights feasible?" IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 45, no. 3, pp. 2751–2768, 2023.
- [46] K. Xu, J. Li, M. Zhang, S. S. Du, K.-i. Kawarabayashi, and S. Jegelka, "What can neural networks reason about?" in International Conference on Learning Representations, 2020.
- [47] K. Xu, M. Zhang, J. Li, S. S. Du, K.-i. Kawarabayashi, and S. Jegelka, "How neural networks extrapolate: From feedforward to graph neural networks," in International Conference on Learning Representations, 2021.
- C. Shi, L. Pan, H. Hu, and I. Dokmanić, "Homophily modulates double descent generalization in graph convolution networks," Proceedings of the National Academy of Sciences, vol. 121, no. 8, p. e2309504121, 2024.
- [49] A. Vasileiou, S. Jegelka, R. Levie, and C. Morris, "Survey on generalization theory for graph neural networks," arXiv preprint arXiv:2503.15650,
- [50] M. Hardt, B. Recht, and Y. Singer, "Train faster, generalize better: Stability of stochastic gradient descent," in Proceedings of The 33rd International Conference on Machine Learning, 2016, pp. 1225–1234.
- [51] W. Cong, M. Ramezani, and M. Mahdavi, "On provable benefits of depth in training graph convolutional networks," in Proceedings of the 35rd International Conference on Neural Information Processing Systems,, 2021, pp. 9936–9949.
- [52] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," in International Conference on Learning Representations, 2014.
- [53] H. Li, M. Wang, S. Liu, P.-Y. Chen, and J. Xiong, "Generalization guarantee of training graph convolutional networks with graph topology sampling," in Proceedings of The 39th International Conference on Machine Learning, 2022, pp. 13014–13051.
- [54] N. Keriven, A. Bietti, and S. Vaiter, "Convergence and stability of graph convolutional networks on large random graphs," in Proceedings of the 34th International Conference on Neural Information Processing Systems, 2020, pp. 21512-21523.
- X. Zhou, K. Hu, and H. Wang, "A tighter generalization error bound for wide gcn based on loss landscape," Applied and Computational Harmonic Analysis, p. 101777, 2025.
- [56] S. Jegelka, "Theory of graph neural networks: Representation and learning," arXiv preprint arXiv:2204.07697, 2022.
- [57] A. Elisseeff, T. Evgeniou, M. Pontil, and L. P. Kaelbing, "Stability of randomized learning algorithms." Journal of Machine Learning Research, vol. 6, no. 1, 2005.
- [58] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad, "Collective classification in network data," AI Magazine, vol. 29, no. 3, pp. 93-93, 2008.
- [59] Z. Yang, W. Cohen, and R. Salakhudinov, "Revisiting semi-supervised learning with graph embeddings," in Proceedings of the International Conference on Machine Learning, 2016, pp. 40–48.

  T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning*
- Representations, 2017.
- [61] G. Li, C. Xiong, G. Qian, A. Thabet, and B. Ghanem, "Deepergcn: training deeper gcns with generalized aggregation functions," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 45, no. 11, pp. 13 024–13 034, 2023.
- [62] T. K. Rusch, M. M. Bronstein, and S. Mishra, "A survey on oversmoothing in graph neural networks," arXiv preprint arXiv:2303.10993, 2023.
- [63] T. Chen, K. Zhou, K. Duan, W. Zheng, P. Wang, X. Hu, and Z. Wang, "Bag of tricks for training deeper graph neural networks: A comprehensive benchmark study," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 45, no. 3, pp. 2769–2781, 2022.

- [64] X. Zheng, Y. Liu, S. Pan, M. Zhang, D. Jin, and P. S. Yu, "Graph neural networks for graphs with heterophily: A survey," arXiv preprint arXiv:2202.07082, 2022.
- [65] J. Zhu, Y. Yan, M. Heimann, L. Zhao, L. Akoglu, and D. Koutra, "Heterophily and graph neural networks: Past, present and future," IEEE
- Data Engineering Bulletin, vol. 47, no. 2, pp. 10–32, 2023.
  [66] H. Li, X. Wang, Z. Zhang, and W. Zhu, "OOD-GNN: Out-of-distribution generalized graph neural network," *IEEE Transactions on Knowledge* and Data Engineering, vol. 35, no. 7, pp. 7328-7340, 2022.
- —, "Out-of-distribution generalization on graphs: A survey," IEEE Transactions on Pattern Analysis and Machine Intelligence, 2025, DOI: 10.1109/TPAMI.2025.3593897.
- [68] A. Baranwal, K. Fountoulakis, and A. Jagannath, "Graph convolution for semi-supervised classification: Improved linear separability and out-of-distribution generalization," arXiv preprint arXiv:2102.06966, 2021.