

# DeepGRAND: Deep Graph Neural Diffusion

Khang Nguyen\*

*FSOFT AI Center*

Ho Chi Minh City, Vietnam

khang.nguyenhoang.vn@gmail.com

Hieu Nong\*

*FSOFT AI Center*

Hanoi, Vietnam

hieubkvn123@gmail.com

Khuong Nguyen

*FSOFT AI Center*

Ho Chi Minh City, Vietnam

khuongnd@gmail.com

Tan M. Nguyen\*\*

*University of California, Los Angeles*

Los Angeles, USA

tanmnguyen89@ucla.edu

Vinh Nguyen\*\*

*FSOFT AI Center*

Ho Chi Minh City, Vietnam

vinhnguyen.maths@gmail.com

**Abstract**—Graph neural networks (GNNs) have achieved remarkable success in numerous domains. Nevertheless, many popular GNNs are known to suffer from over-smoothing. This phenomenon causes node features to become indistinguishable and hurts the classification accuracy as layer depth increases, which limits their effectiveness at capturing complex and long range graph interactions. We propose the Deep Graph Neural Diffusion (DeepGRAND), a continuous-depth graph neural network that is based on the diffusion process on graphs that theoretically alleviates the over-smoothing issue. DeepGRAND perturbs the learnable graph diffusivity and re-scales the underlying diffusion equation by a data-dependent term. We empirically show that DeepGRAND mitigates the accuracy drop-off caused by over-smoothing and surpasses the best accuracy achieved by popular GNNs on various graph deep learning benchmarks. We further demonstrate the advantage of DeepGRAND over many existing graph neural networks in the low label rate regimes.

**Index Terms**—deep learning, graph neural networks, diffusion

## I. INTRODUCTION

Graph neural networks (GNNs) and machine learning on graphs [1] have been successfully applied in a wide range of applications including physical modeling [2], recommender systems [3], and social networks [4]. Recently, advanced GNNs have been developed to further improve the performance of the models and extend their application, which include graph convolutional networks (GCNs) [5], GraphSAGE [6], message passing neural networks (MPNNs) [2], and graph attention networks (GATs) [7]. A well-known problem of GNNs is that the performance of the model decreases significantly with increasing depth. This phenomenon is a common plight of most GNN architectures, and it is referred to as the over-smoothing issue of GNNs [8], [9].

Partial differential equations (PDEs) have been extensively studied and used in a variety of applications, including image processing [10] and computer graphics [11]. Interpreting GNNs as discretization schemes of the underlying diffusion PDE, [12] proposes Graph Neural Diffusion (GRAND), a novel class of continuous-depth GNNs. This framework allows the vast literature on PDEs to be used to improve

GNN performance. Despite its strengths, GRAND does not escape the plight of over-smoothing.

**Main contributions.** In this paper, we propose the Deep Graph Neural Diffusion (DeepGRAND), a novel continuous-depth graph neural network that improves on various aspects of the baseline GRAND [12]. At its core, DeepGRAND introduces a data-dependent scaling term and a perturbation to the diffusion dynamics to alter the spectral and convergent characteristics of the process. With this design, DeepGRAND attains the following advantages:

- 1) DeepGRAND inherits the diffusive characteristic of GRAND while significantly mitigating the over-smoothing issue.
- 2) DeepGRAND achieves remarkably better performance than popular GNNs and GRAND variants when fewer nodes are labeled as training data, meriting its use in low-labeling rates situations.
- 3) Feature representation under the dynamics of DeepGRAND is guaranteed to remain bounded, ensuring numerical stability.

## II. BACKGROUND

### A. Graph neural networks and the over-smoothing issue

The majority of the literature on GNNs can be derived from a common message passing scheme, forming a large class of MPNNs that differ from each other by how information is propagated and aggregated in each model [13]. An update rule for all MPNNs can be written in the form

$$\mathbf{H}_u = \xi \left( \mathbf{X}_u, \bigoplus_{v \in \mathcal{N}_u} \mu(\mathbf{X}_u, \mathbf{X}_v) \right),$$

where  $\mu$  is a learnable message function,  $\bigoplus$  is a permutation-invariant aggregate function, and  $\xi$  is the update function. Three popular MPNNs are GCN [5], GraphSAGE [6], and GAT [7]. In these specific GNN designs, the update rule can be further simplified to

$$\mathbf{H}_u = \sigma \left( \sum_{v \in \mathcal{N}_u \cup \{u\}} a_{uv} \mu(\mathbf{X}_v) \right), \quad (1)$$

\*: co-first authors, \*\*: co-last authors

where  $a$  is either given by the normalized augmented adjacency matrix (GCN), a random sampling scheme (GraphSAGE), or the attention mechanism (GAT),  $\mu$  is a linear transformation, and  $\sigma$  is an activation function.

It is noted that a GNN with more message passing layers suffers from over-smoothing [8], which makes learned features become indistinguishable and hurts the classification accuracy. This phenomenon goes against the common understanding that the deeper the model, the better its learning capacity is [14]. Due to over-smoothing, GNNs are constrained to having small depths and encounter difficulty in capturing long-range dependencies.

Over the past few years, a flurry of research into understanding and alleviating the over-smoothing issue has been conducted [9], [15]–[17].

### B. Graph Neural Diffusion

GRAND is a continuous-depth architecture for deep learning on graphs proposed by [12]. It draws inspiration from the heat diffusion process in mathematical physics and follows the same vein as other PDE-inspired neural networks [15], [18], [19]. Central to the formulation of GRAND is the diffusion equation on graphs

$$\frac{\partial \mathbf{X}(t)}{\partial t} = \text{div}[\mathbf{G}(\mathbf{X}(t), t)\nabla \mathbf{X}(t)], \quad (2)$$

where  $\mathbf{G} = \text{diag}(a(\mathbf{X}_i(t), \mathbf{X}_j(t), t))$  is a diagonal matrix giving the diffusivity between connected vertices, which describes the thermal conductance property of the graph.

Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be the matrix representing the initial  $d$  features of all  $n$  nodes in the graph. The architecture utilises the encoder-decoder design given by  $\mathbf{Y} = \psi(\mathbf{X}(T))$ , where for all  $T$ ,  $\mathbf{X}(T) \in \mathbb{R}^{n \times d}$  is given by

$$\mathbf{X}(T) = \mathbf{X}(0) + \int_0^T \frac{\partial \mathbf{X}}{\partial t}(t) dt, \quad \text{with } \mathbf{X}(0) = \phi(\mathbf{X}). \quad (3)$$

In the simplest case when  $\mathbf{G}$  is only dependent on the initial node features, the differential in (3) simplifies to

$$\frac{\partial \mathbf{X}}{\partial t}(t) = (\mathbf{A}(\mathbf{X}) - \mathbf{I})\mathbf{X}(t), \quad (4)$$

where  $\mathbf{A}(\mathbf{X}) = [(a(\mathbf{X}_i(t), \mathbf{X}_j(t), t))]$  is an  $n \times n$  matrix with the same structure as the adjacency matrix of the graph. From now on, we omit the term  $\mathbf{X}$  when writing  $\mathbf{A}(\mathbf{X})$ . The entries of  $\mathbf{A}$  are exactly those in  $\mathbf{G}$ , and thus determine the diffusivity. Furthermore,  $\mathbf{A}$  can be informally thought of as the attention weight between vertices. Building upon this heuristic, GRAND models the attention matrix  $\mathbf{A}$  in (4) by the multi-head self-attention mechanism, where

$$\mathbf{A} = \frac{1}{h} \sum_{l=1}^h \mathbf{A}^l(\mathbf{X}) \quad (5)$$

with  $h$  being the number of heads and the attention matrix  $\mathbf{A}^l(\mathbf{X}) = (a^l(\mathbf{X}_i, \mathbf{X}_j))$ , for  $l = 1, \dots, h$ . This specific implementation is called GRAND-1 since it becomes a linear matrix differential equation once  $\mathbf{A}$  is calculated from the initial node

features. The authors of [12] also proposed GRAND-nl, an implementation of GRAND where  $\mathbf{G}$  is dependent on the current node features at each step.

As has been pointed out by [12], [20], many GNN architectures, including GAT and GCN, can be formalized as a discretization scheme of (2) if no non-linearity is used between the layers. The intuition behind the connection can readily be seen from (1), where each subsequent node data is computed similarly or equal to a weighted average of the neighboring node features. This gives rise to the diffusive nature of these architectures.

### III. DOES GRAND SUFFER FROM OVER-SMOOTHING?

Various works [9], [15] have attributed the occurrence of the over-smoothing issue to the exponential convergence of node representations. We formally define this phenomenon for continuous depth graph neural networks.

**Definition 1.** Let  $\mathbf{X}(t) \in \mathbb{R}^{n \times d}$  denote the feature representation at time  $t \geq 0$ .  $\mathbf{X}$  is said to experience over-smoothing if there exists a vector  $\mathbf{v} \in \mathbb{R}^d$  and constants  $C_1, C_2 > 0$  such that for  $\mathbf{V} = (\mathbf{v}, \mathbf{v}, \dots, \mathbf{v})^\top$

$$\|\mathbf{X}(t) - \mathbf{V}\|_\infty \leq C_1 e^{-C_2 t}. \quad (6)$$

This definition is similar in spirit to the one given by [15].

**Proposition 1.** The GRAND-1 dynamics given by equations (3), (4), and (5) cause  $\mathbf{X}$  to experience over-smoothing.

*Sketch of proof.*  $\mathbf{A}$  is a right-stochastic matrix with positive entries, hence its Perron-Frobenius eigenvalue is  $\alpha_1 = 1$ . Moreover,  $\alpha_1$  is a simple eigenvalue and its one-dimensional eigenspace has the basis  $\mathbf{u} = (1, 1, \dots, 1)$ . Suppose  $\{\alpha_1, \alpha_2, \dots, \alpha_k\}$  is the complex spectrum of  $\mathbf{A}$ . That is, they are all complex eigenvalues of  $\mathbf{A}$ . The matrix  $\mathbf{A} - \mathbf{I}$  has eigenvalues  $\beta_i = \alpha_i - 1$  for all  $i = \overline{1, k}$ , so  $\beta_1 = 0$  and  $\text{Re } \beta_i < 0$  for all  $i = \overline{2, k}$ . A Jordan canonical form  $\mathbf{J}$  of  $\mathbf{A} - \mathbf{I}$  consists of a single entry 0 on the top left most corner, followed by several Jordan blocks  $\mathbf{J}_i$  corresponding to eigenvalues  $\beta_i$  on the diagonal of the matrix. Let  $\mathbf{P}$  be the change of basis matrix such that  $\mathbf{J} = \mathbf{P}^{-1}(\mathbf{A} - \mathbf{I})\mathbf{P}$  and the first column of  $\mathbf{P}$  is  $\mathbf{u}$  and let  $\mathbf{Z}(t) = \mathbf{P}^{-1}\mathbf{X}(t)$ . We can rewrite (4) as

$$\frac{\partial \mathbf{Z}}{\partial t}(t) = \mathbf{J}\mathbf{Z}(t).$$

The solution to this matrix differential equation is

$$\mathbf{Z}(t) = \exp(t\mathbf{J})\mathbf{Z}(0).$$

Since  $\text{Re } \beta_i < 0$  for all  $i = \overline{2, k}$ , we readily obtain

$$\lim_{t \rightarrow \infty} \mathbf{Z}(t) = \lim_{t \rightarrow \infty} \exp(t\mathbf{J})\mathbf{Z}(0) = (\mathbf{Z}(0)_{1,1}, \dots, \mathbf{Z}(0)_{1,d})^\top$$

where the convergence is exponential in  $\infty$ -norm with regard to  $t$ . Recall that the first column of  $\mathbf{P} = \mathbf{u}$ , we have

$$\lim_{t \rightarrow \infty} \mathbf{X}(t) = \mathbf{P} \lim_{t \rightarrow \infty} \mathbf{Z}(t) = \mathbf{V}$$

for

$$\mathbf{V} = \begin{pmatrix} \mathbf{Z}(0)_{1,1} & \mathbf{Z}(0)_{1,2} & \dots & \mathbf{Z}(0)_{1,d} \\ \mathbf{Z}(0)_{1,1} & \mathbf{Z}(0)_{1,2} & \dots & \mathbf{Z}(0)_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{Z}(0)_{1,1} & \mathbf{Z}(0)_{1,2} & \dots & \mathbf{Z}(0)_{1,d} \end{pmatrix}.$$

We deduce that  $\|\mathbf{X}(t) - \mathbf{V}\|_\infty$  exponentially converges to 0, and we can choose  $C_1, C_2$  suitably so that the inequality (6) is satisfied.  $\square$

Proposition 1 conclusively shows that GRAND-1 still suffers from over-smoothing. Of course, if we consider more general variants of GRAND, the above arguments no longer hold. Although we can not rigorously assert that over-smoothing affect all implementations of GRAND, it can be argued that the occurrence of this phenomenon in GRAND intuitively makes sense since diffusion have the tendency to "even out" distribution over time. Hence, a purely diffusion based model like GRAND is inherently ill-suited for deep networks.

#### IV. DEEPGRAND: DEEPER GRAPH NEURAL DIFFUSION

##### A. Model formulation

We propose a new model of continuous Graph Neural Network based on GRAND capable of learning at much higher depth (DeepGRAND). It leverages a perturbation to the diffusivity matrix and a scaling factor to make the model both more stable and more resilient to the over-smoothing issue.

Denote the  $i$ -th column of  $\mathbf{X}$  by  $\mathbf{X}_i$ . With the usual 2-norm, we define the column-wise norm matrix  $\langle \mathbf{X}(t) \rangle^\alpha \in \mathbb{R}^{n \times d}$  as

$$\langle \mathbf{X}(t) \rangle^\alpha = \begin{pmatrix} \|\mathbf{X}_1\|^\alpha & \|\mathbf{X}_2\|^\alpha & \dots & \|\mathbf{X}_d\|^\alpha \\ \|\mathbf{X}_1\|^\alpha & \|\mathbf{X}_2\|^\alpha & \dots & \|\mathbf{X}_d\|^\alpha \\ \vdots & \vdots & \ddots & \vdots \\ \|\mathbf{X}_1\|^\alpha & \|\mathbf{X}_2\|^\alpha & \dots & \|\mathbf{X}_d\|^\alpha \end{pmatrix}, \quad (7)$$

for some constant  $\alpha > 0$ . The dynamics of DeepGRAND is given by (3), with  $\frac{\partial \mathbf{X}}{\partial t}$  given by

$$\frac{\partial \mathbf{X}}{\partial t}(t) = (\mathbf{A} - (1 + \epsilon)\mathbf{I}) \mathbf{X}(t) \odot \langle \mathbf{X}(t) \rangle^\alpha, \quad (8)$$

where  $\odot$  is the Hadamard product. Hence, our proposed dynamics differs from GRAND's dynamics by the perturbation  $\epsilon$  and the data-dependent scaling term  $\langle \mathbf{X}(t) \rangle^\alpha$ .

##### B. How DeepGRAND overcomes the over-smoothing issue

To explain the motivation behind DeepGRAND, we first note that the convergence property of all nodes to a pre-determined feature vector is not necessarily an undesirable trait. It guarantees that our model will remain bounded and not "explode", which helps to ensure numerical stability. Our perturbation by  $\epsilon$  serves to strengthen this behavior. Its purpose is to reduce the real part of all eigenvalues to negative values, making all node representations converge to 0 in the long run.

We observe that the rate of convergence is the chief factor in determining the range of effective depth. If the model converges very slowly, it is clear that we can train it at high depth without ever having to worry about over-smoothing. As

TABLE I  
PERFORMANCE OF DEEPGRAND VS. GRAND AT DIFFERENT DEPTHS.  
THE BEST RESULTS ARE HIGHLIGHTED IN BOLD.

Dataset	Depth	GRAND-1	GRAND-nl	DeepGRAND
Cora	4	83.5 $\pm$ 0.4	83.3 $\pm$ 0.3	<b>84.7 <math>\pm</math> 0.2</b>
	16	83.5 $\pm$ 0.5	82.5 $\pm$ 0.4	<b>84.5 <math>\pm</math> 0.3</b>
	32	81.8 $\pm$ 0.5	81.7 $\pm$ 0.3	<b>84.3 <math>\pm</math> 0.4</b>
	64	76.5 $\pm$ 0.6	78.7 $\pm$ 1.1	<b>82.1 <math>\pm</math> 0.7</b>
	128	68.6 $\pm$ 0.4	74.6 $\pm$ 0.5	<b>76.4 <math>\pm</math> 0.9</b>
Citeseer	4	68.5 $\pm$ 1.7	69.0 $\pm$ 1.5	<b>75.4 <math>\pm</math> 0.6</b>
	16	60.4 $\pm$ 2.0	60.7 $\pm$ 1.8	<b>75.0 <math>\pm</math> 0.6</b>
	32	60.1 $\pm$ 3.7	58.8 $\pm$ 1.6	<b>74.4 <math>\pm</math> 1.1</b>
	64	55.1 $\pm$ 2.6	55.7 $\pm$ 2.3	<b>71.0 <math>\pm</math> 1.1</b>
	128	50.5 $\pm$ 2.6	52.9 $\pm$ 1.8	<b>67.9 <math>\pm</math> 0.8</b>
Pubmed	4	77.8 $\pm$ 0.1	77.5 $\pm$ 0.3	<b>79.6 <math>\pm</math> 1.0</b>
	16	77.9 $\pm$ 0.4	77.8 $\pm$ 0.2	<b>79.5 <math>\pm</math> 0.6</b>
	32	76.7 $\pm$ 0.3	78.4 $\pm$ 0.4	<b>79.3 <math>\pm</math> 0.4</b>
	64	74.7 $\pm$ 0.5	73.3 $\pm$ 1.2	<b>79.1 <math>\pm</math> 0.6</b>
	128	NA	NA	<b>78.7 <math>\pm</math> 0.5</b>
Computers	1	85.3 $\pm$ 0.6	85.7 $\pm$ 0.1	<b>85.9 <math>\pm</math> 0.3</b>
	2	85.6 $\pm$ 2.9	85.0 $\pm$ 0.7	<b>85.8 <math>\pm</math> 0.3</b>
	4	85.4 $\pm$ 0.5	84.2 $\pm$ 0.5	<b>85.5 <math>\pm</math> 0.5</b>
	8	55.9 $\pm$ 0.3	52.1 $\pm$ 0.4	<b>85.1 <math>\pm</math> 0.3</b>
	16	40.4 $\pm$ 1.0	42.1 $\pm$ 3.0	<b>78.7 <math>\pm</math> 2.2</b>
	32	4.2 $\pm$ 4.1	8.2 $\pm$ 5.1	<b>48.9 <math>\pm</math> 0.3</b>
Photo	1	93.6 $\pm$ 0.2	93.4 $\pm$ 0.3	<b>93.9 <math>\pm</math> 0.3</b>
	2	93.7 $\pm$ 0.2	93.3 $\pm$ 0.3	<b>93.8 <math>\pm</math> 0.3</b>
	4	93.5 $\pm$ 0.3	92.3 $\pm$ 0.4	<b>93.8 <math>\pm</math> 0.2</b>
	8	92.6 $\pm$ 0.5	92.5 $\pm$ 0.6	<b>93.5 <math>\pm</math> 0.4</b>
	16	83.6 $\pm$ 1.7	93.3 $\pm$ 0.4	<b>93.3 <math>\pm</math> 0.4</b>
	32	25.5 $\pm$ 5.6	45.2 $\pm$ 7.3	<b>91.1 <math>\pm</math> 1.0</b>
CoauthorCS	1	90.3 $\pm$ 1.1	91.1 $\pm$ 0.4	<b>91.4 <math>\pm</math> 0.2</b>
	2	90.3 $\pm$ 0.9	90.8 $\pm$ 0.1	<b>91.2 <math>\pm</math> 0.7</b>
	4	90.4 $\pm$ 0.7	89.9 $\pm$ 0.8	<b>90.9 <math>\pm</math> 0.6</b>
	8	85.1 $\pm$ 1.8	87.2 $\pm$ 1.0	<b>90.7 <math>\pm</math> 0.7</b>
	16	65.2 $\pm$ 14.5	77.7 $\pm$ 3.7	<b>88.4 <math>\pm</math> 1.1</b>
	32	48.0 $\pm$ 8.9	49.7 $\pm$ 12.1	<b>79.7 <math>\pm</math> 2.4</b>

such, if we can control the convergence rate, we would be able to alleviate the over-smoothing issue. The addition of the term  $\langle \mathbf{X}(t) \rangle^\alpha \in \mathbb{R}^{n \times d}$  serves this purpose. As the node representations come close to their limits, this term acts as a scaling factor to slow down the convergent process.

We present an exact bound for the convergent rate of our model. We heuristically note that a symmetric matrix will almost surely have a simple spectrum. This characteristic has been made precise in some classes of matrices with certain distributions [21]. We will utilize this assumption in the next proposition.

**Proposition 2.** *Assuming  $\mathbf{A}$  is right-stochastic, symmetric and has a simple spectrum. With the dynamics given in (8), we have the bound for each column  $\mathbf{X}_i$*

$$\|\mathbf{X}_i(T)\| \geq ((2 + \epsilon)\alpha T + \|\mathbf{X}_i(0)\|^{-\alpha})^{-\frac{1}{\alpha}}, \quad (9)$$

$$\|\mathbf{X}_i(T)\| \leq (\epsilon\alpha T + \|\mathbf{X}_i(0)\|^{-\alpha})^{-\frac{1}{\alpha}}. \quad (10)$$

*Proof.* Let  $\mathbf{Y}_i(t) = \|\mathbf{X}_i(t)\|^4$ , we have

$$\begin{aligned} \mathbf{Y}_i'(t) &= 4\|\mathbf{X}_i(t)\|^2 \langle \mathbf{X}_i'(t), \mathbf{X}_i(t) \rangle \\ &= 4\|\mathbf{X}_i(t)\|^2 \langle (\mathbf{A} - (1 + \epsilon)\mathbf{I})\mathbf{X}_i(t) \|\mathbf{X}_i\|^\alpha, \mathbf{X}_i(t) \rangle \\ &= 4\|\mathbf{X}_i(t)\|^{2+\alpha} \langle (\mathbf{A} - (1 + \epsilon)\mathbf{I})\mathbf{X}_i(t), \mathbf{X}_i(t) \rangle \\ &= 4\|\mathbf{X}_i(t)\|^{2+\alpha} \langle \mathbf{A}\mathbf{X}_i(t), \mathbf{X}_i(t) \rangle - 4(1 + \epsilon)\|\mathbf{X}_i(t)\|^{4+\alpha} \end{aligned}$$

TABLE II  
PERFORMANCE OF GRAND, OTHER POPULAR GNNs, AND DEEPGRAND AT DIFFERENT LABEL RATES. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD. (NOTE: THE RESULTS FOR GRAND++-L, GCN, GAT, AND GRAPHsAGE ARE IMPORTED FROM [20])

Dataset	# labeled	GRAND-l	GRAND-nl	GRAND++-l	GCN	GAT	GraphSage	DeepGRAND
Cora	20	82.9 ± 1.1	82.7 ± 2.5	83.0 ± 1.4	82.1 ± 2.0	79.9 ± 2.3	80.0 ± 2.5	<b>84.2 ± 0.7</b>
	10	80.7 ± 2.2	80.5 ± 1.1	80.9 ± 3.0	78.8 ± 5.4	76.3 ± 4.9	75.0 ± 5.0	<b>82.9 ± 0.8</b>
	5	77.1 ± 3.1	77.7 ± 2.9	77.8 ± 4.5	73.9 ± 8.0	71.0 ± 5.7	68.1 ± 7.0	<b>80.9 ± 1.1</b>
	2	74.2 ± 5.6	69.4 ± 4.3	66.9 ± 10.0	60.9 ± 14.0	58.3 ± 13.6	54.0 ± 12.2	<b>76.5 ± 2.0</b>
	1	57.9 ± 8.1	55.9 ± 10.0	54.9 ± 16.1	47.7 ± 15.3	47.9 ± 15.4	43.0 ± 14.0	<b>70.2 ± 3.3</b>
Citeseer	20	71.7 ± 2.9	73.2 ± 3.1	73.5 ± 3.3	74.2 ± 2.9	73.2 ± 2.9	72.0 ± 2.8	<b>74.7 ± 0.8</b>
	10	66.3 ± 4.2	67.8 ± 4.0	72.3 ± 2.4	72.2 ± 3.5	71.4 ± 4.9	68.9 ± 5.1	<b>73.5 ± 0.8</b>
	5	69.0 ± 3.7	66.9 ± 4.6	70.0 ± 3.6	67.2 ± 4.2	67.4 ± 5.1	64.8 ± 5.2	<b>72.0 ± 1.0</b>
	2	58.4 ± 9.0	56.4 ± 5.5	65.0 ± 8.3	58.1 ± 9.8	55.6 ± 9.2	54.4 ± 11.4	<b>69.7 ± 3.0</b>
	1	49.7 ± 8.7	47.3 ± 6.7	<b>59.0 ± 9.6</b>	48.9 ± 10.2	50.3 ± 14.3	48.8 ± 11.5	58.4 ± 2.5
Pubmed	20	78.4 ± 0.5	75.2 ± 1.8	79.2 ± 1.4	76.9 ± 3.3	75.6 ± 4.1	74.6 ± 3.1	<b>79.5 ± 0.6</b>
	10	74.1 ± 2.0	74.2 ± 2.0	75.1 ± 3.9	72.6 ± 3.2	72.4 ± 3.5	70.7 ± 3.1	<b>78.9 ± 1.5</b>
	5	71.1 ± 1.9	72.1 ± 2.2	72.0 ± 1.9	68.7 ± 7.9	68.5 ± 5.8	66.1 ± 6.2	<b>77.1 ± 1.1</b>
	2	71.4 ± 3.9	65.5 ± 9.5	69.3 ± 4.9	60.5 ± 16.2	60.2 ± 14.4	59.0 ± 12.7	<b>72.3 ± 1.8</b>
	1	62.4 ± 7.6	63.5 ± 5.5	65.9 ± 4.9	58.6 ± 12.8	58.8 ± 12.8	55.5 ± 12.7	<b>70.0 ± 1.8</b>
Computers	20	84.0 ± 1.0	83.3 ± 1.2	85.7 ± 0.5	82.9 ± 1.5	80.1 ± 1.8	80.0 ± 1.0	<b>87.1 ± 0.5</b>
	10	82.3 ± 2.2	81.3 ± 3.0	83.0 ± 0.8	82.5 ± 0.7	76.0 ± 0.4	74.7 ± 1.3	<b>85.7 ± 1.1</b>
	5	78.7 ± 0.8	79.7 ± 2.0	<b>82.6 ± 0.6</b>	82.5 ± 1.0	71.4 ± 7.3	64.8 ± 1.6	82.4 ± 0.3
	2	66.2 ± 11.3	65.1 ± 8.3	76.5 ± 1.5	<b>76.9 ± 1.5</b>	65.1 ± 8.9	42.6 ± 42.9	76.6 ± 1.4
	1	49.8 ± 15.0	47.3 ± 11.2	67.7 ± 0.4	49.5 ± 1.7	37.1 ± 7.8	27.7 ± 2.4	<b>69.3 ± 2.7</b>
Photo	20	93.2 ± 0.4	91.8 ± 1.5	<b>93.6 ± 0.4</b>	92.0 ± 0.1	89.4 ± 2.5	91.3 ± 0.7	93.5 ± 0.7
	10	90.8 ± 1.4	89.0 ± 2.5	90.7 ± 1.2	90.4 ± 0.4	87.4 ± 2.4	84.4 ± 1.8	<b>92.3 ± 0.4</b>
	5	88.1 ± 2.6	88.3 ± 1.6	88.3 ± 1.2	88.9 ± 1.6	83.0 ± 3.6	78.3 ± 1.9	<b>90.5 ± 1.0</b>
	2	82.6 ± 3.0	80.6 ± 4.4	83.7 ± 0.9	83.6 ± 0.7	76.9 ± 4.9	51.9 ± 4.2	<b>85.1 ± 0.3</b>
	1	75.1 ± 5.4	76.3 ± 4.8	83.1 ± 0.8	82.9 ± 2.2	73.6 ± 8.2	45.4 ± 7.1	<b>83.3 ± 1.9</b>
CoauthorCS	20	91.0 ± 0.6	90.6 ± 1.0	90.8 ± 0.3	91.1 ± 0.4	80.0 ± 2.9	91.3 ± 0.4	<b>91.7 ± 0.6</b>
	10	89.0 ± 2.1	<b>90.0 ± 0.7</b>	86.9 ± 0.5	88.6 ± 0.5	74.7 ± 3.4	89.7 ± 0.4	89.8 ± 0.7
	5	84.2 ± 3.6	87.0 ± 2.0	84.8 ± 0.8	86.7 ± 0.4	71.7 ± 4.5	<b>89.1 ± 0.7</b>	88.2 ± 0.7
	2	75.2 ± 3.8	76.7 ± 6.85	76.5 ± 1.9	<b>83.6 ± 1.5</b>	63.1 ± 6.1	76.5 ± 1.3	82.1 ± 3.4
	1	56.6 ± 8.4	66.44 ± 8.17	60.3 ± 1.5	65.2 ± 2.3	51.1 ± 5.2	61.4 ± 1.4	<b>71.1 ± 1.7</b>

Since  $\mathbf{A}$  is a right-stochastic and symmetric matrix, we have

$$\langle \mathbf{A}\mathbf{X}_i(t), \mathbf{X}_i(t) \rangle \leq \|\mathbf{A}\| \|\mathbf{X}_i(t)\|^2 = \|\mathbf{X}_i(t)\|^2,$$

and

$$-\langle \mathbf{A}\mathbf{X}_i(t), \mathbf{X}_i(t) \rangle \leq \|\mathbf{A}\| \|\mathbf{X}_i(t)\|^2 = \|\mathbf{X}_i(t)\|^2,$$

so that

$$-\|\mathbf{X}_i(t)\|^2 \leq \langle \mathbf{A}\mathbf{X}_i(t), \mathbf{X}_i(t) \rangle \leq \|\mathbf{X}_i(t)\|^2.$$

Hence, we deduce that

$$4\|\mathbf{X}_i(t)\|^{4+\alpha}(-2-\epsilon) \leq \mathbf{Y}_i'(t) \leq 4\|\mathbf{X}_i(t)\|^{4+\alpha}(-\epsilon). \quad (11)$$

Multiply both sides of (11) with  $-\frac{\alpha}{4}\|\mathbf{X}_i(t)\|^{-4-\alpha} = -\frac{\alpha}{4}\mathbf{Y}_i(t)^{-1-\alpha/4}$ , and by noting that  $-\frac{\alpha}{4}\mathbf{Y}_i'\mathbf{Y}_i^{-1-\alpha/4} = (\mathbf{Y}_i^{-\alpha/4})'$ , we have

$$\alpha(2+\epsilon) \geq (\mathbf{Y}_i^{-\alpha/4})' \geq \alpha\epsilon.$$

Integrate from 0 to  $T$  and rearranging the appropriate terms, we get

$$(2+\epsilon)\alpha T + \mathbf{Y}_i(0)^{-\alpha/4} \geq \mathbf{Y}_i(T)^{-\alpha/4} \geq \epsilon\alpha T + \mathbf{Y}_i(0)^{-\alpha/4}.$$

Finally, by noting that  $\mathbf{Y}_i^{-\alpha/4} = \|\mathbf{X}_i\|^{-\alpha}$ , we obtain the desired bounds.  $\square$

The inequalities (9), (10) ensure that each column  $\mathbf{X}_i$  is bounded in norm on both sides by polynomial-like terms.

Hence, it no longer converges at exponential speed. The over-smoothing problem is thus alleviated.

## V. EXPERIMENTAL RESULTS

In this section, we conduct experiments to compare the performance of our proposed method DeepGRAND and the baseline GRAND. We show that DeepGRAND is able to achieve better accuracy when trained with higher depth and a limited number of labeled nodes per class. We also compare our results with those of GRAND++ [20] and popular GNNs architectures such as GCN [3], GAT [7], and GraphSage [6]. GRAND++ is a GRAND variant where a source term is added into the dynamics to reduce the effect of over-smoothing. In table II, we reported the results of GRAND++-l, the best-performing variant of GRAND++.

### A. Experiment design

We evaluate our proposed method using the following experiments:

- **Performance with different depths:** For all benchmarks, we trained variants of GRAND and DeepGRAND using the same integration limits  $T$  across a wide range of values and compare the respective test accuracies.
- **Performance with different label rates:** For all benchmarks, we trained GRAND-l, GRAND-nl, other popular GNNs, and DeepGRAND with different numbers of labeled nodes per class.

**Datasets.** We conduct our experiments on a wide range of popular node classification datasets, including Cora [22], Citeseer [23], Pubmed [24], Computers [25], Photo [25], and CoauthorCS [26].

### B. DeepGRAND is more adaptable to deep architectures

In table I, we provide empirical evidence for our argument that the dynamics of DeepGRAND is more resilient to deeper architectures. Specifically, we compared the change in performance of both GRAND and DeepGRAND as the integration limit  $T$  increases. Following [12], for Cora, Citeseer, and Pubmed, we used the default planetoid split with 10 random seeds per split. For Computers, Photo, and CoauthorCS, we split the datasets randomly and used 10 random seeds per split. For each seed, we used the default label rate of 20 labeled nodes per class. On all benchmarks, our proposed dynamics performed substantially better across all depths. Furthermore, the performance degradation as depth increases is less significant than that of GRAND, indicating that over-smoothing has been alleviated.

### C. DeepGRAND is more resilient under limited labeled training data

In Table II, we provide empirical results to demonstrate that DeepGRAND achieves superior accuracies compared to other GNNs with a limited number of labeled nodes per class. For all of the benchmarks: Cora, Citeseer, Pubmed, Computers, Photo, and CoauthorCS, we used grid search to find the optimal  $T$  values and evaluate the performance under different numbers of labeled nodes per class. For each dataset, we experimented with 1, 2, 5, 10, and 20 labeled nodes per class and compared the test accuracies between GRAND, GRAND++, GCN, GAT, GraphSage, and DeepGRAND. Our results show that DeepGRAND outperforms both GRAND variants and common GNN architectures like GCN, GAT, and GraphSage.

## VI. CONCLUSION

We propose DeepGRAND, a novel class of continuous-depth graph neural networks that leverage a data-dependent scaling term and a perturbation to the graph diffusivity to decrease the saturation rate of the underlying diffusion process, thus alleviating the over-smoothing issue. We also prove that the proposed method stabilizes the learning of the model by controlling the norm of the node representation. We theoretically and empirically show the advantage of DeepGRAND over GRAND and other popular GNNs in terms of resiliency to over-smoothing and overall performance. It is interesting to leverage advanced methods in improving the Neural ODEs [27] to further develop DeepGRAND.

## REFERENCES

- [1] Bronstein, M., Bruna, J., LeCun, Y., Szlam, A. & Vandergheynst, P. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*. **34**, 18-42 (2017)
- [2] Gilmer, J., Schoenholz, S., Riley, P., Vinyals, O. & Dahl, G. Neural message passing for quantum chemistry. *International Conference On Machine Learning*. pp. 1263-1272 (2017)
- [3] Ying, R., He, R., Chen, K., Eksombatchai, P., Hamilton, W. & Leskovec, J. Graph convolutional neural networks for web-scale recommender systems. *Proceedings Of The 24th ACM SIGKDD International Conference On Knowledge Discovery And Data Mining*. pp. 974-983 (2018)
- [4] Qiu, J., Tang, J., Ma, H., Dong, Y., Wang, K. & Tang, J. DeepInf: Social Influence Prediction with Deep Learning. *Proceedings Of The 24th ACM SIGKDD International Conference On Knowledge Discovery And Data Mining (KDD'18)*. (2018)
- [5] Kipf, T. & Welling, M. Semi-Supervised Classification with Graph Convolutional Networks. *Proceedings Of The 5th International Conference On Learning Representations*. (2017)
- [6] Hamilton, W., Ying, R. & Leskovec, J. Inductive Representation Learning on Large Graphs. *Proceedings Of The 31st International Conference On Neural Information Processing Systems*. pp. 1025-1035 (2017)
- [7] Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P. & Bengio, Y. Graph Attention Networks. *International Conference On Learning Representations*. (2018)
- [8] Li, Q., Han, Z. & Wu, X. Deeper Insights into Graph Convolutional Networks for Semi-Supervised Learning. (arXiv,2018)
- [9] Oono, K. & Suzuki, T. Graph Neural Networks Exponentially Lose Expressive Power for Node Classification. *International Conference On Learning Representations*. (2020)
- [10] Perona, P. & Malik, J. Scale-space and edge detection using anisotropic diffusion. *IEEE Transactions On Pattern Analysis And Machine Intelligence*. **12**, 629-639 (1990)
- [11] Sun, J., Ovsjanikov, M. & Guibas, L. A Concise and Provably Informative Multi-Scale Signature Based on Heat Diffusion. *Comput. Graph. Forum*. **28** pp. 1383-1392 (2009,7)
- [12] Chamberlain, B., Rowbottom, J., Goronova, M., Webb, S., Rossi, E. & Bronstein, M. GRAND: Graph Neural Diffusion. *Proceedings Of The 38th International Conference On Machine Learning, (ICML) 2021, 18-24 July 2021, Virtual Event*. (2021)
- [13] Bronstein, M., Bruna, J., Cohen, T. & Veličković, P. Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges.
- [14] Li, G., Müller, M., Ghanem, B. & Koltun, V. Training Graph Neural Networks with 1000 layers. *International Conference On Machine Learning (ICML)*. (2021)
- [15] Rusch, T., Chamberlain, B., Rowbottom, J., Mishra, S. & Bronstein, M. Graph-Coupled Oscillator Networks. *Proceedings Of The 39th International Conference On Machine Learning*. **162** pp. 18888-18909 (2022,7,17)
- [16] Zhao, L. & Akoglu, L. PairNorm: Tackling Oversmoothing in GNNs. *International Conference On Learning Representations*. (2020)
- [17] Nguyen, K., Nguyen, T., Nong, H., Nguyen, V., Ho, N. & Osher, S. Revisiting Over-smoothing and Over-squashing Using Ollivier-Ricci Curvature. *International Conference On Machine Learning (ICML)*. (2023), in press.
- [18] Chamberlain, B., Rowbottom, J., Eynard, D., Giovanni, F., Dong, X. & Bronstein, M. Beltrami Flow and Neural Diffusion on Graphs. *Advances In Neural Information Processing Systems*. (2021)
- [19] Eliasof, M., Haber, E. & Treister, E. PDE-GCN: Novel Architectures for Graph Neural Networks Motivated by Partial Differential Equations. *Advances In Neural Information Processing Systems*. (2021)
- [20] Thorpe, M., Nguyen, T., Xia, H., Strohmmer, T., Bertozzi, A., Osher, S. & Wang, B. GRAND++: Graph Neural Diffusion with A Source Term. *International Conference On Learning Representations*. (2022)
- [21] Tao, T. & Vu, V. Random matrices have simple spectrum. *Combinatorica*. (2014,12)
- [22] McCallum, A., Nigam, K., Rennie, J. & Seymore, K. Automating the Construction of Internet Portals with Machine Learning. *Information Retrieval*. **3**, 127-163 (2000,7)
- [23] Giles, C., Bollacker, K. & Lawrence, S. CiteSeer. *Proceedings Of The Third ACM Conference On Digital Libraries - DL '98*. (1998)
- [24] Sen, P., Namata, G., Bilgic, M., Getoor, L., Galligher, B. & Eliassirad, T. Collective Classification in Network Data. *AI Magazine*. **29**, 93 (2008,9)
- [25] McAuley, J., Targett, C., Shi, Q. & Hengel, A. Image-based Recommendations on Styles and Substitutes. (arXiv,2015)
- [26] Monti, F., Boscai, D., Masci, J., Rodolà, E., Svoboda, J. & Bronstein, M. Geometric deep learning on graphs and manifolds using mixture model CNNs. (arXiv,2016)
- [27] Chen, R., Rubanova, Y., Bettencourt, J. & Duvenaud, D. Neural Ordinary Differential Equations. *Advances In Neural Information Processing Systems*. **31** (2018)