

A Deep Graph Wavelet Convolutional Neural Network for Semi-supervised Node Classification

Jingyi Wang*, Zhidong Deng†

Institute for Artificial Intelligence at Tsinghua University (THUAI)

State Key Laboratory of Intelligent Technology and Systems

Beijing National Research Center for Information Science and Technology (BNRist)

Center for Intelligent Connected Vehicles and Transportation

Department of Computer Science, Tsinghua University, Beijing 100084, China

*wang-jy20@mails.tsinghua.edu.cn

†michael@mail.tsinghua.edu.cn

Abstract—Graph convolutional neural network provides good solutions for node classification and other tasks with non-Euclidean data. There are several graph convolutional models that attempt to develop deep networks but do not cause serious over-smoothing at the same time. Considering that the wavelet transform generally has a stronger ability to extract useful information than the Fourier transform, we propose a new deep graph wavelet convolutional network (DeepGWC) for semi-supervised node classification tasks. Based on the optimized static filtering matrix parameters of vanilla graph wavelet neural networks and the combination of Fourier bases and wavelet ones, DeepGWC is constructed together with the reuse of residual connection and identity mappings in network architectures. Extensive experiments on three benchmark datasets including Cora, Citeseer, and Pubmed are conducted. The experimental results demonstrate that our DeepGWC outperforms existing graph deep models with the help of additional wavelet bases and achieves new state-of-the-art performances eventually.

Index Terms—graph convolutional neural network, wavelet transform, filtering matrix, network architecture

I. INTRODUCTION

Convolutional neural networks (CNNs) achieve outstanding performance in a wide range of tasks that are based on Euclidean data such as computer vision [1] and recommender systems [2] [3]. However, graph-structured non-Euclidean data is also very common in real life. Many research areas like literature citation networks and knowledge graphs have the graph data structure. In this case, graph convolutional networks (GCNs) often have an outstanding exhibition, which are successfully applied in social analysis [5] [6], citation network [7] [8] [9], transport forecasting [10] [11], and other promising fields. For example, in a literature citation network, articles are usually represented as nodes and citation relationships as edges among nodes. When addressing the challenge of semi-supervised node classifications on such non-Euclidean graph data, classical GCNs first extract and aggregate features of articles using the graph Fourier transform and the convolution theorem, and then classify unlabeled articles according to output features of graph convolutions.

GCN and graph wavelet neural network (GWNN) [12] have their own characteristics. In fact, GCN-based methods that

employ the Fourier transform as their theoretical basis generally have high computational costs caused by eigendecomposition and frequent multiplications between dense matrices. ChebyNet attempts to reduce the computational burden via polynomial approximation but the locality of convolution is still not guaranteed. GWNN tries to address the limitations of GCN by using graph wavelet transform as the theoretical basis instead of graph Fourier transform. Graph wavelet transform has good sparsity in both spatial domain and spectral domain. Besides, wavelets in GWNN are highly localized in the vertex domain. Such localized property makes a great contribution to the performance of GWNN [12]. However, GWNN generally uses independent diagonal filtering matrices for every graph convolution layer. As a result, the memory occupied by the network would increase heavily when stacking more layers in GWNN. To avoid exceeding memory, the number of layers in GWNN models on a large dataset is quite limited. Therefore, GWNN fails to extract deep and high-level features of nodes.

The lately published GCN models suffer from so-called over-smoothing when more layers are stacked. In other words, after repeatedly applying Laplacian smoothing, the features of nodes within each connected component of the graph would converge to the same values [7]. There are several models that expand shallow graph convolution methods in order to get better results. For example, APPNP [13] and GDC [14] relieve the over-smoothing problem by using the Personalized PageRank matrix. Moreover, GCNII [15] makes use of residual skills and thus surpasses the performance of the previous trials, although only the graph Fourier transform is still utilized as its bases.

In this paper, we propose a new deep graph wavelet convolutional network called DeepGWC that improves the deep GCNII model. First, the diagonal filtering matrix in graph wavelet convolution of vanilla GWNN models is modified. We no longer set separate learnable convolution kernels that are initialized randomly for each convolution layer. By contrast, we specify the filtering matrices of all layers with selected static elements, i.e., the filtering matrix is not updated during training. As a result, the memory required by the model is greatly reduced. Second, we make use of graph wavelet

transform and attach the optimized wavelet bases to previous graph Fourier bases of GCNs to enhance the feature extraction capability. Considering the high efficiency and high sparsity of graph wavelet transform [12], the import of wavelet bases does not bring too much computational complexity. Following GCNII, we leverage residual connection and identity mappings to deepen the above graph wavelet convolutional network with the optimized filtering matrix. Extensive experiments on semi-supervised node classification are finished. On datasets Cora, Citeseer, and Pubmed, our experimental results demonstrate that the new DeepGWC model improves GCNII. DeepGWC is able to extract useful features and effectively relieve the over-smoothing problem at the same time. This allows the proposed method to outperform the baselines with 2, 4, 8, 16, 32, and 64 layers and achieve new state-of-the-art results.

The contributions of our work are summarized as follows.

- 1) We optimize the diagonal filtering matrices in graph wavelet convolution of vanilla GWNNs by specifying them with selected static elements that are not updated during training instead of learnable parameters. In this way, the memory requirement caused by increasing parameters when stacking more graph wavelet convolutional layers is solved.
- 2) We attach wavelet bases to the previous graph Fourier bases of classical GCNs as new bases of our DeepGWC model to improve the capability of extracting features of GCNII. With the reuse of residual connections and identity mappings of GCNII, DeepGWC becomes a deep graph wavelet convolutional network with the optimized filtering matrix.
- 3) We perform experiments on three benchmark datasets, i.e. Cora, Citeseer, and Pubmed. Our DeepGWC model achieves state-of-the-art results in the semi-supervised node classification task.

II. RELATED WORK

A. Graph Convolutional Neural Network

CNN achieves great performance on computer vision [1], natural language processing [16], and other fields [17] [18] [19] [20]. Its success motivates researchers to define convolution operators on graphs. In this way, CNN can be generalized to graphs like social networks and citation networks. Related methods are classified into spatial methods and spectral ones.

Spatial methods perform extraction of spatial features for topological graphs based on similar ideas as weighted summation in CNN. Using different sampling methods for neighborhood nodes, such methods calculate the weighted sum of features of sampled nodes and then update node features. For example, MoNet [21] computes weighted neighborhood node features instead of using average values. GraphSAGE [22] presents an inductive framework that leverages node feature information to efficiently give rise to node embeddings for previously unseen data. GraphsGAN [23] generates fake samples to improve performance on graph-based semi-supervised learning.

Spectral methods carry out the convolution on graphs with the help of spectral theory. They perform convolution operations on topological graphs in accordance with the theory of spectral graph. In fact, the concept of graph Fourier transform is derived from graph signal processing. According to graph Fourier transform, graph convolution is defined. With the rapid development of deep learning methods, GCN in spectral methods also appear. Spectral CNN [24] first implements CNN on graphs using graph Fourier transform. GCN [25] motivates the choice of convolutional architecture via a localized first-order approximation of spectral graph convolutions.

B. Graph Wavelet Neural Networks

Wavelet transform develops the idea of localization of short-time Fourier transform (STFT) and overcomes the weakness that the window size does not change with frequency at the same time. The wavelet transform is a signal encoding algorithm that is sparser than Fourier transform and has a very strong capability of information expression.

D.K. Hammond proposes a method for constructing wavelet transforms of functions defined on the vertices of an arbitrary finite weighted graph [26]. GraphWave [27] learns a multidimensional structural embedding for each node based on the diffusion of a spectral graph wavelet centered at the node. It provides mathematical guarantees on the optimality of learned structural embeddings. Furthermore, GWNN [12] addresses the shortcomings of previous spectral graph CNN methods with graph Fourier transform. It takes graph wavelets as a set of bases instead of eigenvectors of graph Laplacian. GWNN redefines graph convolution based on graph wavelets, achieving high efficiency and high sparsity. Besides, GWNN also detaches the feature transformation from graph convolution and thus yields better results.

C. Deep Graph Convolution Models

It has been proved that features of nodes in a graph would converge to the same values when repeatedly applying graph convolution operations on graphs [7]. In order to relieve such an over-smoothing problem and further enhance the feature extraction capability, several modified models for citation networks have been proposed. For example, JKNet [28] makes use of the jump connection and the adaptive aggregation mechanism to adapt to local neighborhood properties and tasks. IncepGCN takes advantage of the inception network and could be optimized by Dropedge [29]. In particular, GCNII [15] is viewed as an effective deep graph convolutional network, which is analyzed that a K -layer GCNII model can express a polynomial spectral filtering of order K with arbitrary coefficients. Though various efforts have been made to improve the performance of deep graph convolution models, it is still shallow models that usually obtain the best result.

III. METHODS

Considering that graph wavelets have a very strong ability to express information, this paper proposes a new deep graph wavelet convolutional network (DeepGWC) to improve the

deep graph convolution models [15]. Our DeepGWC model not only achieves the new state-of-the-art performance, but also could be converted into vanilla GCN, GWNN, and other shallow models by simply adjusting hyperparameters for different application scenarios.

A. Preliminary

Let $G = (V, E)$ be a simple and connected undirected graph, where V represents the set of nodes with $|V| = n$ and E denotes the set of edges among nodes in V with $|E| = m$. Assume \mathbf{A} to express the adjacency matrix of G . \mathbf{D} indicates the diagonal degree matrix with $\mathbf{D}_{i,i} = \sum_j \mathbf{A}_{i,j}$. Let $\mathbf{L} = \mathbf{I}_n - \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$ stand for the normalized Laplacian matrix, where \mathbf{I}_n is the identity matrix with dimensions $n \times n$. Obviously, \mathbf{L} is a symmetric positive semidefinite matrix with eigendecomposition $\mathbf{L} = \mathbf{U} \Lambda \mathbf{U}^T$, where $\mathbf{U} \in \mathbb{R}^{n \times n}$ denotes the complete set of orthonormal eigenvectors and Λ is a diagonal matrix of real, non-negative eigenvalues.

We use $\tilde{G} = (V, \tilde{E})$ to indicate the cyclic graph of G , i.e., the graph with a cycle attached to each node in G . Correspondingly, the adjacency matrix of \tilde{G} is $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}_n$ and the diagonal degree matrix is $\tilde{\mathbf{D}} = \mathbf{D} + \mathbf{I}_n$. Let $\tilde{\mathbf{L}} = \mathbf{I}_n - \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$ expresses the normalized Laplacian matrix of the cyclic graph \tilde{G} .

B. GCN

1) *Graph Fourier Transform*: Graph Fourier transform takes the eigenvectors of the normalized Laplacian matrix, i.e., \mathbf{U} as a set of bases. $\hat{\mathbf{x}} = \mathbf{U}^T \mathbf{x}$ performs graph Fourier transform on a signal $\mathbf{x} \in \mathbb{R}^n$ of graph G .

2) *Graph Convolution Operation*: Based on graph Fourier transform, the graph convolution operation with a filter $g_\gamma(\Lambda) = \text{diag}(\gamma)$ is defined as

$$g_\gamma(\mathbf{L}) * \mathbf{x} = \mathbf{U} g_\gamma(\Lambda) \mathbf{U}^T \mathbf{x} \quad (1)$$

where $\gamma \in \mathbb{R}^n$ represents the vector of spectral filtering coefficients. ChebyNet [30] restricts the convolution kernel g to a polynomial expansion, i.e.,

$$g_\theta = \sum_{k=0}^K \theta_k \Lambda^k \quad (2)$$

where K is a hyperparameter to control the size of node neighborhoods. $\theta \in \mathbb{R}^{K+1}$ stands for a vector of polynomial coefficients. Furthermore, with $\mathbf{L} = \mathbf{U} \Lambda \mathbf{U}^T$, the graph convolution operation in Eq. (1) can be approximated by,

$$\mathbf{U} g_\theta(\Lambda) \mathbf{U}^T \mathbf{x} \approx \mathbf{U} \left(\sum_{k=0}^K \theta_k \Lambda^k \right) \mathbf{U}^T \mathbf{x} = \left(\sum_{k=0}^K \theta_k \mathbf{L}^k \right) \mathbf{x}. \quad (3)$$

With $K = 1, \theta_0 = 2\theta$, and $\theta_1 = -\theta$, and the renormalization trick [25], we denote $\tilde{\mathbf{P}} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$ for the l -th graph convolution layer. We have

$$\mathbf{H}^{l+1} = \sigma(\tilde{\mathbf{P}} \mathbf{H}^l \mathbf{W}^l) \quad (4)$$

where σ expresses the activation function and $\mathbf{W}^l \in \mathbb{R}^{p \times q}$ is the parameter matrix for feature transformation. The graph

convolution layer in Eq. (4) transforms the input feature matrix of $\mathbf{H}^l \in \mathbb{R}^{n \times p}$ into the output feature matrix of $\mathbf{H}^{l+1} \in \mathbb{R}^{n \times q}$.

C. Graph Wavelet Convolutional Network

1) *Graph Wavelet Transform*: Similar to graph Fourier transform, the graph wavelet transform need a proper set of bases to project graph signal from vertex domain into the spectral domain. Here we indicate the wavelet bases as ψ_s , where s is a scaling parameter. ψ_s can be defined by

$$\psi_s = \mathbf{U} \mathbf{G}_s \mathbf{U}^T \quad (5)$$

where \mathbf{G}_s is a scaling matrix that has similar effects to $g_\gamma(\Lambda)$ in graph Fourier transform. In this case, \mathbf{G}_s can be evaluated by $\mathbf{G}_s = \text{diag}(g(s\Lambda))$, where g is an exponential function. We can get ψ_s^{-1} by replacing s with $-s$. By using graph wavelets as bases, $\hat{\mathbf{x}} = \psi_s \mathbf{x}$ conducts graph Fourier transform on a signal $\mathbf{x} \in \mathbb{R}^n$ of graph G .

2) *Graph Wavelet Convolution Operation*: Substituting Fourier transform in graph convolution operation with wavelet transform, we can find the structure of the l -th graph wavelet convolution layer below:

$$\mathbf{H}^{l+1} = \sigma(\psi_s \mathbf{F} \psi_s^{-1} \mathbf{H}^l \mathbf{W}^l) \quad (6)$$

where \mathbf{F} indicates the diagonal matrix for graph convolution kernel [12]. In Eq. (6), σ , \mathbf{W}^l , \mathbf{H}^l , and \mathbf{H}^{l+1} have the same definitions as those in Eq. (4).

For vanilla GWNN, \mathbf{F} is a diagonal filtering matrix learned in the spectral domain and independent for every layer. In experiments done in [12], a two-layer graph wavelet neural network is just designed for three datasets. In this case, the memory consumption is not too large for such a shallow GWNN. However, when trying to increase the number of layers to 8 or 16, the memory usage would rise rapidly. Actually, the out-of-memory problem is particularly prone to appear in the experiments on the Pubmed dataset. In order to solve such difficulties, we simplify the use of diagonal filtering matrix \mathbf{F} . Instead of exploiting independent filtering matrices for every layer, \mathbf{F} with the same static parameters is employed. By optimizing the elements of diagonal elements of \mathbf{F} as a set of static parameters, the wavelet basis can also be obtained statically like the Fourier basis. In this way, we add no parameters to be trained.

D. Network Architecture against Over-Smoothing

1) *Residual Connection*: The initial residual connection constructs a link to the initial feature representation H^0 [15]. We make sure that the final representation still contains a portion of the initial feature no matter how many layers we stack [15].

2) *Identity Mapping*: Like the identity mapping in ResNet [4], we append an identity matrix to the weight matrix \mathbf{W}^l in the case of $p = q$. Thus a direct pathway is established between the input and the output of a convolutional layer through importing identity mapping [15].

E. Our DeepGWC Model

In this paper, we build a deep graph wavelet convolutional network (DeepGWC). On the basis of the above-mentioned, we define the l -th layer of DeepGWC as

$$\mathbf{H}^{l+1} = \sigma(\mathbf{H}^l \mathbf{W}^l) \quad (7)$$

where σ is the activation function, \mathbf{H}^l is the results of graph convolution on \mathbf{H} , and \mathbf{W}^l is the optimized feature transformation matrix. \mathbf{H}^l and \mathbf{W}^l in Eq. (7) are described as follows:

$$\mathbf{H}^l = (1 - \alpha)\tilde{\mathbf{P}}^l \mathbf{H}^l + \alpha \mathbf{H}^0, \quad (8)$$

$$\mathbf{W}^l = \beta_l \mathbf{W}^l + (1 - \beta_l) \mathbf{I} \quad (9)$$

where α represents the ratio of the initial residual term \mathbf{H}^0 , β_l is the ratio of the original feature transformation matrix \mathbf{W}^l of the l -th layer, and $0 \leq \alpha, \beta \leq 1$. $\tilde{\mathbf{P}}^l$ in Eq. (8) is given by

$$\tilde{\mathbf{P}}^l = \gamma(\psi_s \mathbf{F} \psi_s^{-1}) + (1 - \gamma)(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}) \quad (10)$$

where γ stands for the ratio of the graph wavelet term and $0 \leq \gamma \leq 1$. As described in subsection III-C2, \mathbf{F} is optimized with static elements, i.e., the modified \mathbf{F} is the product of a static constant f and the identity matrix. In fact, γ and f could be combined and the first item of Eq. (10) could be written as

$$\gamma(\psi_s \mathbf{F} \psi_s^{-1}) = \gamma f(\psi_s \psi_s^{-1}). \quad (11)$$

However, to emphasize that we combine the Fourier bases and the wavelet bases, we still treat γ and f as two separate parameters later, i.e., we perform experiments according to Eq. (10). In this way, we could keep consistency with the form of vanilla wavelet bases in Eq. (6) to emphasize our modification about the filtering matrix \mathbf{F} .

There are the three key hyperparameters, i.e., α , β , and γ in our DeepGWC model. Through adjusting such three parameters, the DeepGWC could be reduced to the following basic models like:

- With $\alpha = 0$, $\beta = 1$ and $\gamma = 0$, Eq. (7) can be simplified to Eq. (4) and our DeepGWC model is equivalent to the vanilla GCN model [25].
- With $\alpha = 0$, $\beta = 1$ and $\gamma = 1$, Eq. (7) can be simplified to Eq. (6) and our DeepGWC model is then reduced to the GWNN model [12].
- With $\alpha \neq 0$, $\beta = 0$ and $\gamma = 0$, the weight matrix \mathbf{W}^l is ignored and our DeepGWC model then functions like the APPNP model [13].
- With $\alpha \neq 0$, $\beta \neq 0$ and $\gamma = 0$, the weight matrix \mathbf{W}^l is ignored and our DeepGWC model then resembles the GCNII model [15].

IV. EXPERIMENTS

We compare the proposed DeepGWC model with competitive baselines, i.e., the existing state-of-the-art graph neural network models on the semi-supervised node classification task.

A. Datasets and Baselines

1) *Datasets*: We use three commonly-used citation network datasets, i.e., Cora, Citeseer, and Pubmed [31]. The statistics of datasets are summarized in Table I. For every dataset, 20 labeled nodes per class are adopted for training in the overall assessment.

TABLE I
THE STATISTICS OF THREE CITATION NETWORK BENCHMARKS

Dataset	#Nodes	#Edges	#Classes	Features	Label Rate
Cora	2,708	5,429	7	1,433	5.2%
CiteSeer	3,327	4,732	6	3,703	3.6%
Pubmed	19,717	44,338	3	500	0.3%

2) *Baselines*: Considering that the vanilla GCN [25] is one of the basic model graph convolutional methods, we exploit it as an important baseline. Meanwhile, spectral CNN [24], ChebyNet [30], GAT [32], APPNP [13], and GWNN [12] are also utilized for performance comparison. In the experiment on the number of stacked layers, we are more interested in comparison with JKNet [28], IncepGCN [29], and GCNII [15] because these three models are often employed to set up deeper graph convolutional networks. In the experiments on the label rate, we focus on M3S [33] and some different training strategies [7] that improve the generalization capability of GCNs on graphs with few labeled nodes.

B. Implementation Details

There are two hyperparameters in the calculation of graph wavelets in general. s denotes the scaling parameter and t stands for the threshold of ψ_s and ψ_s^{-1} . Elements that are less than t are ignored and reset to zero in ψ_s and ψ_s^{-1} . Following the parameter settings of [12], we select the values of s and t . For Cora, we set $s = 1.0$, $t = 1e - 4$, and the calculated density of wavelet bases is 2.81%. For Citeseer, we set $s = 0.7$, $t = 1e - 5$, and the density value is 1.52%. For Pubmed, we set $s = 0.5$, $t = 1e - 7$, and the density value is 5.03%. We also provide the density values of wavelet bases. The wavelet bases with low density would not give rise to the calculation burden too much.

Note that we employ a linear connection layer before the graph convolution layers of DeepGWC that projects the original feature dimension on the hidden embedding dimension \mathbf{H}^0 . Similarly, there is also a linear connection layer after the graph convolution layers that maps the hidden embedding dimension to the output layer \mathbf{Y} , which has the same number of neurons as classes delivered in the dataset. Then we perform the logarithm-softmax operation below:

$$\mathbf{Z}_{i,j} = \log\left(\frac{e^{\mathbf{Y}_{i,j}}}{\sum_k e^{\mathbf{Y}_{i,k}}}\right) \quad (12)$$

where \mathbf{Z} denotes the actual output of the whole DeepGWC model. In the experiments, the negative log likelihood (NLL) loss function is exploited. Moreover, the proportion of nodes classified correctly is used as the metric to evaluate the

accuracy of models. We adopt the Adam optimizer to train our DeepGWC. Following [15], we set $\beta_l = \log(1 + \eta/l) \approx \eta/l$ where η is a hyperparameter and l means the l -th layer of DeepGWC. Dropout is adopted to avoid the overfitting issue.

C. Overall Assessment

On Cora, Citeseer, and Pubmed, we follow the parameter settings of the vanilla GCN [25] to ensure impartial performance comparison. For every dataset, 20 labeled nodes per class are adopted for training. In order to optimize the parameters of the DeepGWC model, 500 samples are used for validation and 1,000 ones for testing so as to get final results. We report the experimental results of the proposed

TABLE II
THE PERFORMANCE COMPARISON OF DEEPGWC WITH EXISTING GRAPH NEURAL NETWORK METHODS

Methods	Cora	CiteSeer	Pubmed
Spectral CNN [24]	73.3%	58.9%	73.9%
ChebyNet [30]	81.2%	69.8%	74.4%
GCN [25]	81.5%	70.3%	79.0%
GAT [32]	83.1%	70.8%	78.5%
APPNP [13]	83.3%	71.8%	80.1%
JKNet [28]	81.1%	69.8%	78.1%
IncepGCN [29]	81.7%	70.2%	77.9%
GWNN [12]	82.8%	71.7%	79.1%
GCNII [15]	85.5%	73.4%	80.3%
DeepGWC	86.4%	75.0%	81.6%

model yielded on three datasets including Cora, CiteSeer, and Pubmed as shown in Table II. We quote the experimental results of Spectral CNN [24], ChebyNet [30], GCN [25], and GWNN [12] that are already reported in [12]. The results of GAT [32], APPNP [13], JKNet [28], and IncepGCN [29] are provided according to [15]. For Cora, CiteSeer, and Pubmed, the depths of JKNet [28] are 4, 16, and 32, respectively, while the depths of IncepGCN [29] are 64, 4, and 4, respectively. By contrast, our DeepGWC model has depths of 64, 64, and 32, respectively. The others are all shallow models.

For the DeepGWC, we utilize the Adam Optimizer with a learning rate of 0.001. The detailed hyperparameters of DeepGWC of Table II are given as follows. Let f stand for the element of the filtering matrix \mathbf{F} , which is searched over (0.4, 0.8, 1.2, 1.6). Suppose that L represents the number of graph wavelet convolution layers in DeepGWC, d represents the dimension of hidden layers. We get the results of Table II with the parameters of:

- For Cora, we set $L = 64$, $d = 64$, $\alpha = 0.3$, $\eta = 0.8$, $\gamma = 0.4$, $f = 0.4$.
- For Citeseer, we set $L = 64$, $d = 256$, $\alpha = 0.1$, $\eta = 0.8$, $\gamma = 0.4$, $f = 0.4$.
- For Pubmed, we set $L = 32$, $d = 512$, $\alpha = 0.1$, $\eta = 0.4$, $\gamma = 0.4$, $f = 0.6$.

As for the parameters above, larger datasets require higher hidden feature dimension d to express information as completely as possible. We follow the setting of α and η of GCNII [15] basically for the three datasets. The proportion of wavelet

bases in the combined bases is reflected in the value of $\gamma \cdot f$, which is 0.16 for Cora and Citeseer, and 0.24 for Pubmed. The proportion of wavelet bases for Pubmed is higher than that for Cora and Citeseer because there are more nodes in Pubmed than Cora and Citeseer obviously. The locality of wavelet bases is important in the case of a large number of nodes.

Analyzing Table II, the DeepGWC improves GCNII [15] and achieves the best results among all the shallow and deep graph models on the three benchmark datasets. The average accuracies on the three datasets of DeepGWC is higher than that of the other models by at least 1.3%. We believe that the adaptability of our DeepGWC plays an important role in obtaining the best results for semi-supervised node classification problems. DeepGWC outperforms GCNII on every dataset, which indicates the excellence of wavelet bases. In addition, for DeepGWC with 8 graph wavelet convolution layers, we utilize the t-SNE algorithm [35] to visualize the learned embeddings during training on the Cora dataset. As shown in Fig.1, different colors represent different classes of articles. The difference between features of different classes become more obvious as the training proceeds.

D. Deeper Network

We further investigate the influence of network depths on model performance. Table III shows the classification accuracy vs. the network depth. For every dataset, the best accuracy with a fixed number of layers is marked in bold. For specific datasets and methods, the cell with a grey background in each row gives the best result in that row, i.e., the best number of layers. We report the classification results of GCN [25], JKNet [28], IncepGCN [29], GWNN [12], GCNII [15], and our DeepGWC model with 2, 4, 8, 16, 32, and 64 layers, where Dropedge [29] is further equipped with GCN [25], JKNet [28], and IncepGCN [29], respectively.

It is easy to see from Table III that our DeepGWC model improves GCNII [15] and yields the best results with all numbers of layers. The accuracy of the DeepGWC model becomes better roughly as the number of layers increases. GCN [25] and dropped GCN [29] do have good results with 2 layers, but their performance gets worse rapidly when stacking much more layers. Meanwhile, JKNet [28] and IncepGCN [29] reach their best results with higher layers, but the accuracy is still not guaranteed when the number of layers exceeds 32 layers. GWNN [12] suffers more from the over-smoothing problem than the vanilla GCN [25]. Besides, there exists a serious memory issue on Pubmed when using GWNN [12] with more than 4 layers. Among the baselines, the performance of GCNII [15] is exceptionally good because of its success in relieving the over-smoothing problem. The proposed DeepGWC improves GCNII with wavelet bases further.

We accomplish the comparison of DeepGWC with all the converted models GCN [25], GWNN [12], and GCNII [15] models as described in Subsection III-E and accordingly complete the ablation study. Fig.2 shows the classification accuracy vs. the number of layers of our DeepGWC model and the three relevant baselines which motivate our study. It

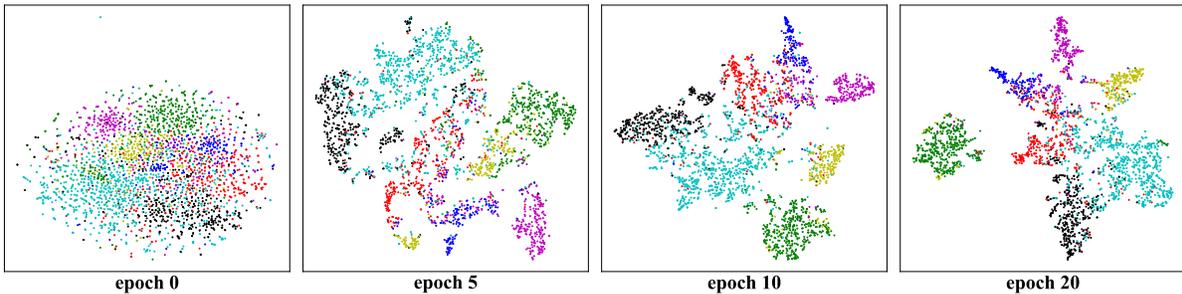


Fig. 1. The visualization of learned features of the Cora dataset during training DeepGWC.

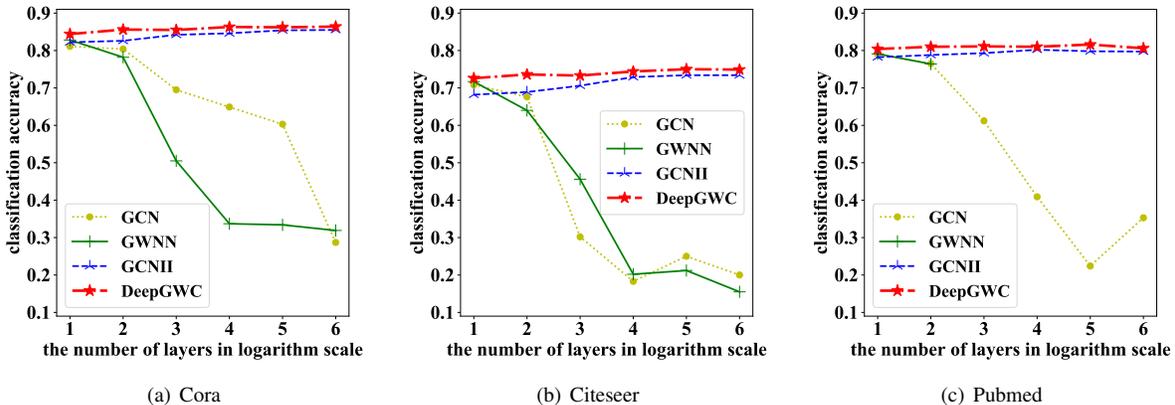


Fig. 2. Performance comparisons of DeepGWC with GCN, GWNN, and GCNII as ablation study.

is readily observed from Fig. 2 that the performance of both GCN [25] and GWNN [12] drops rapidly with the increase of the number of layers, while both GCNII [15] and DeepGWC steadily get better performance when stacking more layers. As a result of the addition of graph wavelet bases, our DeepGWC model achieves better performance than GCNII [15] in the semi-supervised node classification task.

E. Fewer Training Samples

We complete experiments on training with less labeled nodes. MultiStage [33] and M3S [33] are two training methods that enable GCNs to learn graph embeddings with few supervised signals. Following the experiments of M3S [33], we use Label Propagation (LP) [34] and GCN [25] as baselines. For GCN [25], we adopt four training strategies, i.e., Co-training, Self-training, Union, and Intersection additionally [7]. For Cora and Citeseer, we train our DeepGWC model and the baselines with the label rates of 0.5%, 1%, 2%, 3%, and 4%. For Pubmed, we train the models with the label rates of 0.03%, 0.05%, and 0.1% because Pubmed has more nodes. 500 samples are employed for validation and 1,000 ones for testing so as to get final results. We exploit the same number of layers and hyperparameters as that in Table II for GWNN [12], GCNII [15], and DeepGWC. We only adjust the label rates.

Table IV, Table V, and Table VI list the classification results with fewer training samples. Our DeepGWC model obtains

the best performance on most of the experiments with lower label rates. The accuracy of DeepGWC is 5.88% higher than that of M3S [33] on average, which focuses on improving the generalization capability of GCNs on graphs with few labeled nodes.

We compare DeepGWC with GCN [25], GWNN [12], and GCNII [15] in the case of less training data more carefully. The comparison results are shown in Fig. 3. It should be noted that DeepGWC performs better when there is less training data. In other words, the gap between DeepGWC and GCNII [15] is more obvious in the case of lower label rates. In experiments on Cora and Citeseer with the label rates of 0.5% and 1%, DeepGWC exceeds GCNII [15] by 4.1% on average while the gap is 0.73% for the label rates of 2%, 3%, and 4%. Besides, comparing GCN [25] and GWNN [12] that are both shallow models, we find that it is also in the case of lower label rates that GWNN [12] yields better results than that with normal label rates. The comparison between DeepGWC and GCNII [15] and the comparison between GWNN [12] and GCN [25] with fewer training nodes further prove the excellent capability of feature extraction of graph wavelets. It suggests that graph wavelets succeed in extracting useful information even with less labeled nodes.

V. CONCLUSION

In this paper, we present a deep graph wavelet convolutional network called DeepGWC. We modify the filtering matrix

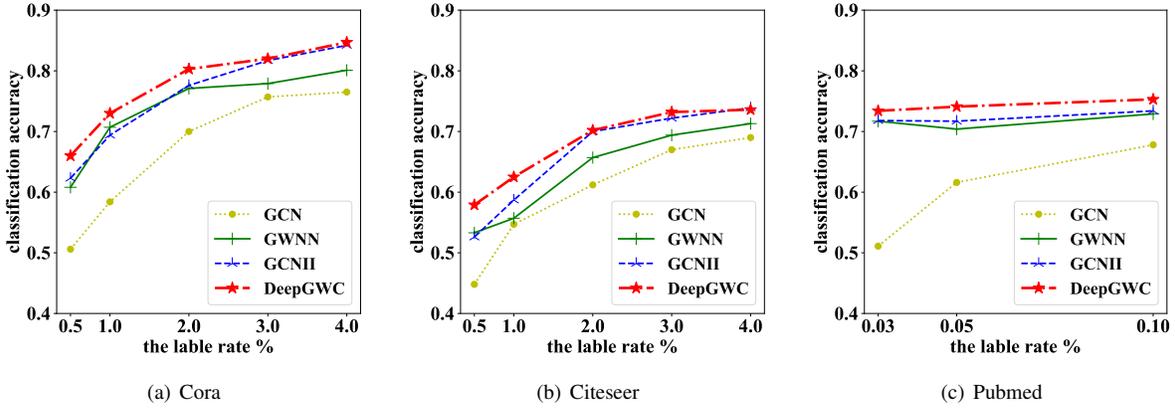


Fig. 3. The classification accuracy with less training samples.

TABLE III
THE RESULTS OF NODE CLASSIFICATIONS WITH DIFFERENT DEPTHS

Dataset	Method	#Layers					
		2	4	8	16	32	64
Cora [31]	GCN [25]	81.1%	80.4%	69.5%	64.9%	60.3%	28.7%
	GCN(drop) [25] [29]	82.8%	82.0%	75.8%	75.7%	62.5%	49.5%
	JKNet [28]	-	80.2%	80.7%	80.2%	81.1%	71.5%
	JKNet(drop) [28] [29]	-	83.3%	82.6%	83.0%	82.5%	83.2%
	IncepGCN [29]	-	77.6%	76.5%	81.7%	81.7%	80.0%
	IncepGCN(drop) [29]	-	82.9%	82.5%	83.1%	83.1%	83.5%
	GWNN [12]	82.8%	78.2%	50.5%	33.7%	33.4%	31.9%
	GCNII [15]	82.2%	82.6%	84.2%	84.6%	85.4%	85.5%
DeepGWC	84.4%	85.6%	85.5%	86.3%	86.2%	86.4%	
Citeseer [31]	GCN [25]	70.8%	67.6%	30.2%	18.3%	25.0%	20.0%
	GCN(drop) [25] [29]	72.3%	70.6%	61.4%	57.2%	41.6%	34.4%
	JKNet [28]	-	68.7%	67.7%	69.8%	68.2%	63.4%
	JKNet(drop) [28] [29]	-	72.6%	71.8%	72.6%	70.8%	72.2%
	IncepGCN [29]	-	69.3%	68.4%	70.2%	68.0%	67.5%
	IncepGCN(drop) [29]	-	72.7%	71.4%	72.5%	72.6%	71.0%
	GWNN [12]	71.7%	64.0%	45.6%	20.2%	21.2%	15.5%
	GCNII [15]	68.2%	68.9%	70.6%	72.9%	73.4%	73.4%
DeepGWC	72.6%	73.6%	73.3%	74.4%	75.0%	74.9%	
Pubmed [31]	GCN [25]	79.0%	76.5%	61.2%	40.9%	22.4%	35.3%
	GCN(drop) [25] [29]	79.6%	79.4%	78.1%	78.5%	77.0%	61.5%
	JKNet [28]	-	78.0%	78.1%	72.6%	72.4%	74.5%
	JKNet(drop) [28] [29]	-	78.7%	78.7%	79.1%	79.2%	78.9%
	IncepGCN [29]	-	77.7%	77.9%	74.9%	OOM	OOM
	IncepGCN(drop) [29]	-	79.5%	78.6%	79.0%	OOM	OOM
	GWNN [12]	79.1%	76.4%	OOM	OOM	OOM	OOM
	GCNII [15]	78.2%	78.8%	79.3%	80.2%	79.8%	79.7%
DeepGWC	80.4%	81.0%	81.1%	81.0%	81.6%	80.7%	

of graph wavelet convolution to make it adaptable to deep graph convolutional models. Applying the combination of Fourier bases and wavelet ones, the proposed DeepGWC model with the reuse of residual connection and identity achieves better performance than existing deep graph models. Extensive experiments on semi-supervised node classifications are conducted. On Cora, Citeseer, and Pubmed, the experimental results demonstrate that our DeepGWC model outperforms the baselines with 2, 4, 8, 16, 32, and 64 layers and yields new state-of-the-art results.

TABLE IV
THE CLASSIFICATION RESULTS WITH LESS TRAINING SAMPLES ON CORA

Method	Label Rate				
	0.5%	1%	2%	3%	4%
LP [34]	57.6%	61.0%	63.5%	64.3%	65.7%
GCN [25]	50.6%	58.4%	70.0%	75.7%	76.5%
Co-training [7]	53.9%	57.0%	69.7%	74.8%	75.6%
Self-training [7]	56.8%	60.4%	71.7%	76.8%	77.7%
Union [7]	55.3%	60.0%	71.7%	77.0%	77.5%
Intersection [7]	50.6%	60.4%	70.0%	74.6%	76.0%
MultiStage [33]	61.1%	63.7%	74.4%	76.1%	72.2%
M3S [33]	61.5%	67.2%	75.6%	77.8%	78.0%
GWNN [12]	60.8%	70.7%	77.1%	77.9%	80.1%
GCNII [15]	62.3%	69.4%	77.6%	81.7%	84.2%
DeepGWC	66.0%	73.0%	80.3%	82.0%	84.7%

TABLE V
THE CLASSIFICATION RESULTS WITH LESS TRAINING SAMPLES ON CITSEER

Method	Label Rate				
	0.5%	1%	2%	3%	4%
LP [34]	37.7%	41.6%	41.9%	44.4%	44.8%
GCN [25]	44.8%	54.7%	61.2%	67.0%	69.0%
Co-training [7]	42.0%	50.0%	58.3%	64.7%	65.3%
Self-training [7]	51.4%	57.1%	64.1%	67.8%	68.8%
Union [7]	48.5%	52.6%	61.8%	66.4%	66.7%
Intersection [7]	51.3%	61.1%	63.0%	69.5%	70.0%
MultiStage [33]	53.0%	57.8%	63.8%	68.0%	69.0%
M3S [33]	56.1%	62.1%	66.4%	70.3%	70.5%
GWNN [12]	53.3%	55.7%	65.7%	69.4%	71.3%
GCNII [15]	52.6%	58.7%	70.0%	72.2%	73.9%
DeepGWC	57.9%	62.5%	70.2%	73.2%	73.6%

ACKNOWLEDGMENT

This work was supported in part by the National Key R&D Program of China under Grant No.2017YFB1302200 and 2018YFB1600804, by TOYOTA TTAD 2020-08, by a grant from the Institute Guo Qiang (2019GQG0002), Tsinghua University, by Alibaba Group through Alibaba Innovative Research Program, by Tencent Group, and by Zhejiang Program in Innovation, Entrepreneurship and Leadership Team

TABLE VI
THE CLASSIFICATION RESULTS WITH LESS TRAINING SAMPLES ON
PUBMED

Method	Label Rate		
	0.03%	0.05%	0.1%
LP [34]	58.3%	61.3%	63.8%
GCN [25]	51.1%	58.0%	67.5%
Co-training [7]	55.5%	61.6%	67.8%
Self-training [7]	56.3%	63.6%	70.0%
Union [7]	57.2%	64.3%	70.0%
Intersection [7]	55.0%	58.2%	67.0%
MultiStage [33]	57.4%	64.3%	70.2%
M3S [33]	59.2%	64.4%	70.6%
GWNN [12]	71.7%	70.4%	72.9%
GCNII [15]	71.8%	71.7%	73.4%
DeepGWC	73.4%	74.1%	75.3%

(2018R01017).

REFERENCES

- [1] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks", in 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 2261–2269.
- [2] R. M. Qiao Liu, Yifu Zeng and H. Zhang, "STAMP: short-term attention/memory priority model for session-based recommendation", in Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018, ACM, 2018, pp. 1831–1839.
- [3] Q. Liu, F. Yu, S. Wu, and L. Wang, "A convolutional click prediction model." in CIKM. ACM, 2015, pp. 1743–1746.
- [4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2016, pp. 770–778.
- [5] C. Li and D. Goldwasser, "Encoding social information with graph convolutional networks for political perspective detection in news media." in ACL. Association for Computational Linguistics, 2019, pp. 2594–2604.
- [6] J. Qiu, J. Tang, H. Ma, Y. Dong, K. Wang, and J. Tang, "Deepin: Social influence prediction with deep learning." in KDD. ACM, 2018, pp. 2110–2119.
- [7] 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*, ser. AAAI '18. AAAI Press, 2018, pp. 3538–3545.
- [8] F. Hu, Y. Zhu, S. Wu, L. Wang, and T. Tan, "Hierarchical graph convolutional networks for semi-supervised node classification," in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI)*, 2019.
- [9] B. Hui, P. Zhu, and Q. Hu, "Collaborative graph convolutional networks: Unsupervised learning meets semi-supervised learning." in AAAI. AAAI Press, 2020, pp. 4215–4222.
- [10] Z. Diao, X. Wang, D. Zhang, Y. Liu, K. Xie, and S. He, "Dynamic spatial-temporal graph convolutional neural networks for traffic forecasting." in AAAI. AAAI Press, 2019, pp. 890–897.
- [11] Y. Han, S. Wang, Y. Ren, C. Wang, P. Gao, and G. Chen, "Predicting station-level short-term passenger flow in a citywide metro network using spatiotemporal graph convolutional neural networks." *ISPRS Int. J. Geo Inf.*, vol. 8, no. 6, p. 243, 2019.
- [12] B. Xu, H. Shen, Q. Cao, Y. Qiu, and X. Cheng, "Graph wavelet neural network." in ICLR, 2019.
- [13] J. Klicpera, A. Bojchevski, and S. Günnemann, "Predict then propagate: Graph neural networks meet personalized pagerank." in ICLR, 2019.
- [14] J. Klicpera, S. Weissenberger, and S. Günnemann, "Diffusion improves graph learning." in NeurIPS, 2019, pp. 13 333–13 345.
- [15] 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, ICML, Proceedings of Machine Learning Research, vol. 119. 2020, pp. 1725–1735.
- [16] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015.
- [17] J. Wang, Q. Liu, Z. Liu, and S. Wu, "Towards accurate and interpretable sequential prediction: A cnn & attention-based feature extractor." in CIKM. ACM, 2019, pp. 1703–1712.
- [18] O. Abdel-Hamid, A. rahman Mohamed, H. Jiang, and G. Penn, "Applying convolutional neural networks concepts to hybrid nn-hmm model for speech recognition." in ICASSP. IEEE, 2012, pp. 4277–4280.
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in Advances in neural information processing systems, 2012, pp. 1097–1105.
- [20] Y. Shen, X. He, J. Gao, L. Deng, and G. Mesnil, "A latent semantic model with convolutional-pooling structure for information retrieval," in Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management, CIKM 2014, Shanghai, China, November 3-7, 2014, pp. 101–110.
- [21] F. Monti, D. Boscaini, J. Masci, E. Rodolà, J. Svoboda, and M. M. Bronstein, "Geometric deep learning on graphs and manifolds using mixture model cnns." in CVPR. IEEE Computer Society, 2017, pp. 5425–5434.
- [22] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs." in NIPS, 2017, pp. 1024–1034.
- [23] M. Ding, J. Tang, and J. Zhang, "Semi-supervised learning on graphs with generative adversarial nets." in CIKM. ACM, 2018, pp. 913–922.
- [24] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs." in ICLR, 2014.
- [25] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks." in ICLR, 2017.
- [26] D. K. Hammond, P. Vandergheynst, and R. Gribonval, "Wavelets on graphs via spectral graph theory," Applied and Computational Harmonic Analysis, no. 2, pp. 129 – 150, 2011.
- [27] C. Donnat, M. Zitnik, D. Hallac, and J. Leskovec, "Learning structural node embeddings via diffusion wavelets." in KDD. ACM, 2018, pp. 1320–1329.
- [28] K. Xu, C. Li, Y. Tian, T. Sonobe, K. ichi Kawarabayashi, and S. Jegelka, "Representation learning on graphs with jumping knowledge networks." in ICML, ser. Proceedings of Machine Learning Research. PMLR, 2018, pp. 5449–5458.
- [29] Y. Rong, W. Huang, T. Xu, and J. Huang, "Dropedge: Towards deep graph convolutional networks on node classification." in ICLR, 2020.
- [30] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, 2016, pp. 3837–3845.
- [31] P. Sen, G. M. Namata, M. Bilgic, L. Getoor, B. Gallagher, and T. Eliassi-Rad, "Collective classification in network data," AI Magazine, vol. 29, no. 3, pp. 93–106, 2008.
- [32] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in ICLR, 2018.
- [33] K. Sun, Z. Lin, and Z. Zhu, "Multi-stage self-supervised learning for graph convolutional networks on graphs with few labeled nodes." in AAAI. AAAI Press, 2020, pp. 5892–5899.
- [34] X.-M. Wu, Z. Li, A. M.-C. So, J. Wright, and S.-F. Chang, "Learning with partially absorbing random walks." in NIPS, 2012, pp. 3086–3094.
- [35] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," Journal of Machine Learning Research, vol. 9, pp. 2579–2605, 2008.