# Graph-adaptive Rectified Linear Unit for Graph Neural Networks

Yifei Zhang
yfzhang@cse.cuhk.edu.hk
The Chinese University of Hong Kong
Hong Kong SAR, China

Hao Zhu
allenhaozhu@anu.edu.au
Australian National University and
Data61/CSIRO
Canberra, Australia

Ziqiao Meng
zqmeng@cse.cuhk.edu.hk
The Chinese University of Hong Kong
Hong Kong SAR, China

Piotr Koniusz
piotr.koniusz@data61.csiro.au
Data61/CSIRO and Australian
National University
Canberra, Australia

Irwin King
king@cse.cuhk.edu.hk
The Chinese University of Hong Kong
Hong Kong SAR, China

## ABSTRACT

Graph Neural Networks (GNNs) have achieved remarkable success by extending traditional convolution to learning on non-Euclidean data. The key to the GNNs is adopting the neural message-passing paradigm with two stages: aggregation and update. The current design of GNNs considers the topology information in the aggregation stage. However, in the updating stage, all nodes share the same updating function. The identical updating function treats each node embedding as i.i.d. random variables and thus ignores the implicit relationships between neighborhoods, which limits the capacity of the GNNs. The updating function is usually implemented with a linear transformation followed by a non-linear activation function. To make the updating function topology-aware, we inject the topological information into the non-linear activation function and propose Graph-adaptive Rectified Linear Unit (GReLU), which is a new parametric activation function incorporating the neighborhood information in a novel and efficient way. The parameters of GReLU are obtained from a hyperfunction based on both node features and the corresponding adjacent matrix. To reduce the risk of overfitting and the computational cost, we decompose the hyperfunction as two independent components for nodes and features respectively. We conduct comprehensive experiments to show that our plug-and-play GReLU method is efficient and effective given different GNN backbones and various downstream tasks.

## CCS CONCEPTS

• **Information systems → Data mining**.

## KEYWORDS

Graph Neural Networks, Rectified Linear Unit, Graph Representation Learning

## 1 INTRODUCTION

Deep learning has revolutionized many machine learning tasks in recent years, ranging from computer vision to speech and natural language understanding [21]. The data in these tasks is typically represented in the Euclidean space. However, there is an increasing number of applications where data is generated from non-Euclidean domains and is represented as graphs with complex relationships between objects [5, 9, 22, 24, 35, 41, 42]. By defining the convolution operators on the graph, graph neural networks (GNNs) extend convolution neural networks (CNNs) from the image domain into the graph domain. Thus, we can still adapt may tools from CNNs to GNNs, *i.e.*, ReLU [28], Dropout [32], or residual link [13].

GNN is now capable of mining the characteristics of the graph by adopting the so-called message passing (MP) framework, which iteratively aggregates and updates the node information following the edge path. Embracing the MP as a modeling tool, recent years witnessed an unprecedented achievement in downstream tasks, such as node/graph classification [30, 31] and link prediction [3, 4]. However, proven by many recent works, an inappropriately designed propagation step in MP violates the homophily assumption held by many real-world graphs, consequently causing critical performance degradation [15, 46].

The aforementioned nuisances motivate us to refine the MP and enhance the generalization ability of GNN. One promising line of exploration is to dynamically consider the contribution of each node during MP. Attention mechanisms are a good pointer as they have been explored in [15, 34]. One of the benefits of attention mechanisms is that they let network focus on the most relevant information to propagate by assigning the weight to different nodes accordingly *e.g.*, works [15, 34] show that propagating information by adaptive transformation yields impressive performance across several inductive benchmarks.

However, the current design of adaptive transformation is inefficient. Our reasoning arises from two critical standpoints. Firstly,
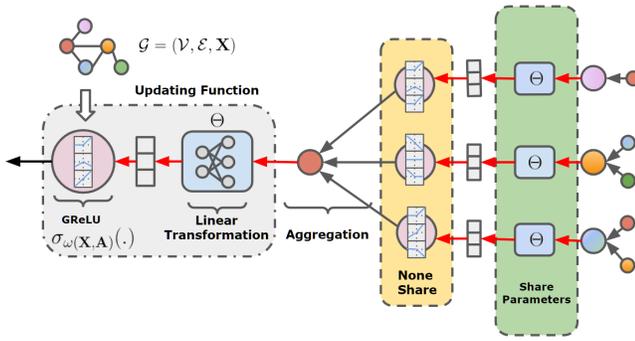
**Figure 1: Plugging GReLU in message passing with three channels. The grey box with dash lines shows details of the updating function which consists of a linear transformation (the blue box) and a non-linear activation function (the pink circle). The black and red arrows indicate the aggregation and updating processes. Note that the linear transformation is shared across different nodes while the activation functions are both node- and channel-specific which are adaptive to both features and the adjacency matrix.**

no topology information has been considered in inference of the adaptive weight for each input node, thus, failing to capture the homophily information of the graph [15]. Secondly, we suspect the limitation in improving model generalization ability and accuracy is due to sharing the adaptive weight across channels for each input node. It may be undesirable to transform all channels with no discrimination as some channels are less important than others.

Motivated by the issues above, we develop GReLU, a topology-aware adaptive activation function, amenable to any GNN backbone. GReLU refines the information propagation and improves the performance on node classification tasks. As illustrated in Figure 1, our mechanism is comprised of two components improving the training procedure of GNNs. On the one hand, the GReLU inherits the topology information by taking features and graphs as input. On the other hand, unlike standard rectifiers being the same function w.r.t. all inputs, GReLU assigns a set of parameters to dynamically control the piecewise non-linear function of each node and channel. Such a topology-adaptive property provides a better expressive power than the topology-unaware counterpart. Meanwhile, GReLU is computationally efficient *e.g.*, a two-layer GCN with GReLU and a two-layer GAT [34] have similar computational complexity, while the former is able to outperform GAT in node classification.

The main contributions of this work are listed below:

i. We incorporate the topology information in the updating function of GNNs by introducing a topology-aware activation function in the updating stage of message passing, which, to our knowledge, is novel.

ii. We propose a well-designed graph adaptive rectified linear unit for graph data, which can be used as a plug-and-play component in various types of spatial and spectral GNNs.

iii. We evaluate the proposed methods on various downstream tasks. We empirically show that the performance of various types of GNNs can be generally improved with the use of our GReLU.

## 2 RELATED WORKS

### 2.1 Graph Neural Networks

GNNs can be categorized into spectral or spatial domain type [37]. The spectral domain is exploring GNNs design from the perspective of graph signal processing *i.e.*, the Fourier transform and the eigen-decomposition of the Laplacian matrix [7] are used to define the convolution step. One can extract any frequency to design a graph filter with the theoretical guarantees. However, these models suffer from three main drawbacks: 1) a large amount of computational burden is induced by the eigen-decomposition of the given graph; 2) spatial non-localization; 3) the filter designed for a given graph cannot be applied to other graphs. Spatial GNNs generalize the architecture of GNNs by introducing the message passing mechanism, which is aggregating neighbor information based on some spatial characteristics of the graph such as adjacency matrix, node features, edge features, *etc.* [11]. Due to its high efficiency and scalability, spatial GNNs have become much more popular in recent works. The most representative model is graph attention network (GAT) [34], which weights the node features of different neighbors by using the attention mechanism in the aggregation step. However, for each node feature, GAT assigns the node a single weight shared across all channels, which ignores the fact that different channels may exhibit different importance.

### 2.2 Activation Functions

Activation functions have been long studied over the entire history of neural networks. One of the few important milestones is introducing the Rectified Linear Unit (ReLU) [28] into the neural network. Due to its simplicity and non-linearity, ReLU became the standard setting for many successful deep architectures. A series of research has been proposed to improve ReLU. Two variants of ReLU are proposed which adopt non-zero slopes $\alpha$ for the input less than zero: Absolute value rectification choose $\alpha = -1$ [14]; LeakyReLU fixes $\alpha$ to a small value [38]. PReLU [12] and [2] takes a further step by making the nonzero slope a trainable parameter. Maxout generalizes ReLU further, by dividing the input into groups and outputting the maximum [10]. Since ReLU is non-smooth, some smooth activation functions have been developed to address this issue, such as soft plus [8], ELU [6], SELU [17], Misc [27], Exp [25] and Sigmoid [20]. These rectifiers are all general functions that can be used in different types of neural networks. None of them has a special design for GNNs by considering the topology information.

### 2.3 Linear Graph Neural Network

In [33, 36, 43–45], the authors proposed simplified GNNs (SGC and $S^2GC$) by removing the activation function and collapsing the weight matrices between consecutive layers. These works hypothesize that the non-linearity between GCN layers is not critical but that the majority of the benefit arises from the topological smoothing. However, removing non-linearity between layers sacrifices the depth efficiency of GNNs. Moreover, the ReLU removed from SGC is not topology-aware. Hence, the rectifiers play a less important role in GNN. As stated in SGC, the topology information is crucial in the aggregation stage to make GNN work, which also motivates us to incorporate the topology information in the updating stage.

## 3 PRELIMINARIES

### 3.1 Notations

In this paper, a graph with node features is denoted as $\mathcal{G} = (X, A)$, where $\mathcal{V}$ is the vertex set, $\mathcal{E}$ is the edge set, and $X \in \mathbb{R}^{N \times C}$ is the feature matrix where $N = |\mathcal{V}|$ is the number of nodes and $C$ is the number of channels. $A \in \{0, 1\}^{N \times N}$ denotes the adjacency matrix of $\mathcal{G}$, i.e., the $(i, j)$-th entry in $A$, $a_{ij}$, is 1 if there is an edge between $v_i$ and $v_j$. We denote $\mathcal{N}(i) = \{j | a_{ij} = 1\}$ as the set of indexes of neighborhoods of node $i$. The degree of $v_i$, denoted as $d_i$, is the number of edges incident with $v_i$. For a $d$-dimensional vector $x \in \mathbb{R}^d$, we use $x_i$ to denote the $i$-th entry of $x$. We use $x_i$ to denote the row vector of $X$ and $x_{ij}$ for the $(i, j)$-th entry of $X$.

### 3.2 Message Passing (MP)

The success of spatial GNNs results from applying message passing (also called neighbor aggregation) between nodes in the graph. In the forward pass of GNNs, the hidden state of node $i$ is iteratively updated via aggregating the hidden state of node $j$ from $i$'s neighbors through edge $e_{ij}$. Suppose there are $l$ iterations, the message passing paradigm can be formalized as:

$$m_i^{(l)} = \text{AGGREGATE}^{(l)}(\{h_i^{(l)}, h_j^{(l)}, e_{ij} | j \in \mathcal{N}(i)\}), \quad (1)$$

$$h_i^{(l+1)} = \sigma(\Theta^{(l)}[h_i^{(l)}, m_i^{(l)}]), \quad (2)$$

where $h_i$ and $h_j$ are hidden states of nodes $i$ and $j$, $m_i$ denotes the message that the node $i$ receives from its neighbors $j \in \mathcal{N}(i)$. The updating function in (2) is modeled as a linear transformation followed by a non-linear activation function $\sigma(\cdot)$ where $\Theta \in \mathbb{R}^{C \times (2C)}$ matrix contains learned parameters for the linear transformation and $[\cdot]$ concatenates vectors.

### 3.3 Parametric ReLU

The vanilla ReLU is a piecewise linear function, $y = \max(x, 0)$, where $\max(\cdot)$ is an element-wise function applied between each channel of the input $x$ and 0. A more generalized way to extend ReLU is to let each channel enjoy different element-wise function:

$$y_c = \sigma_c(x_c) = \max_{1 < k \leq K}\{\alpha^k x_c + \beta^k\}, \quad (3)$$

where $x_c$ and $y_c$ are the $c$-th channels of $x$ and $y$ respectively, $k$ is the $k$-th segmentation in (3), and $\{(\alpha^k, \beta^k)\}_{k=1}^K$ is the set of parameters for the parametric ReLU. Note that when $\alpha^1 = 0, \beta^1 = 0, \alpha^2 = 1, \beta^2 = 0$, then (3) reduces to ReLU, whereas maxout [10] is a special case if $K = 2$. Furthermore, instead of learning parameters $(\alpha^k, \beta^k)$ directly, approach [2] proposes the hyperfunction $\omega(x)$ to estimate sample-specific $(\alpha^k, \beta^k)$ for different channels.

## 4 METHOD

We refer to the proposed adaptive rectified linear unit for GNNs as **GReLU**. The rest of the sections are organized as follows. Firstly, we introduce the parametric form of GReLU by defining two functions: the hyperfunction that generates parameters and activation function for computing the output. Then, we detail the architecture of hyperfunction of GReLU, and we link GReLU with GAT from the perspective of neighborhood weighting. We finally analyze the time complexity of GReLU and compare it with prior works.
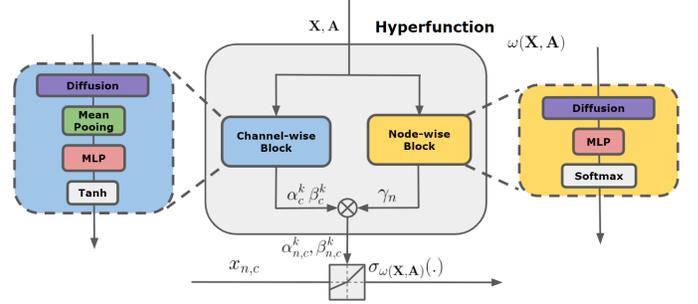


Figure 2: The architecture of GReLU. The arrows in the box indicate the flow for computing the parameter for hyperfunction. The arrow outside the box indicates the calculation of GReLU. The channel-wise block produces the $K \times C$ parameters for $C$ channels and $K$ segments. The node-wise block generates $N$ parameters for $N$ nodes. The total $K \times C \times N$ parameters are obtained via computing the outer product of channel-wise and node-wise parameters.

### 4.1 Graph-adaptive Rectified Linear Unit (GReLU)

Given the input feature of node $n$ where $x_n \in \mathbb{R}^C$, the GReLU is defined as the element-wise function:

$$\sigma_{\omega(X,A)}(x_{nc}) = \max_{1 \leq k \leq K}\{\alpha_{nc}^k x_{nc} + \beta_{nc}^k\}, \quad (4)$$

which consists of two parts:

- Hyperfunction $\omega(X, A)$ whose inputs are the adjacency matrix $A$ and the node feature matrix $X$, and whose output parameters $\{(\alpha_{nc}^k, \beta_{nc}^k)\}_{k=1}^K$ are used by the parametric activation function.
- The parametric activation function $\sigma(\cdot)$ generates activation with parameters $\alpha_{nc}^k$.

The hyperfunction encodes the adjacency matrix $A$ and node features matrix $X$ to determine the parameters of GReLU, adding the topology information to the non-linear activation function. By plugging GReLU into MP, the updating function becomes adaptive, gaining more representational power than its non-adaptive counterpart.

### 4.2 Design of Hyperfunction of GReLU

Note that to fully get GReLU, the hyperfunction needs to produce $K$ sets of parameters $\{(\alpha_{nc}^k, \beta_{nc}^k)\}_{k=1}^K$ for $C$ channels and $N$ nodes. The total number of parameters needed in GReLU is $2K \times C \times N$. We can simply model the hyperfunction $\omega(X, A)$ as a one-layer graph convolution neural network (GCN) with output dimensions $2K \times C \times N$. However, this implementation results in too many parameters in the hyperfunction, which poses the risk of overfitting, as we observe performance degradation in practice. To solve this issue, we decompose the parameters of nodes from channels by modeling them as two different blocks. One block learns the channel-wise parameters $\alpha_c^k, \beta_c^k$ and another block learns node-wise parameters $\gamma_n$. The final outputs $(\alpha_{nc}^k, \beta_{nc}^k)$ are computed as:

$$\alpha_{nc}^k = \gamma_n \alpha_c^k \quad \text{and} \quad \beta_{nc}^k = \gamma_n \beta_c^k. \quad (5)$$

Yifei Zhang, Hao Zhu, Ziqiao Meng, Piotr Koniusz, and Irwin King

**Table 1: Different Variants of ReLU.**

| Type | Formula | Learnable | Input | Plugin |
|------|---------|-----------|-------|--------|
| ReLU | $\max(x, 0)$ | No | $X$ | Yes |
| LReLU | $\max(x, \alpha x)$ | No | $X$ | Yes |
| ELU | $\max(x, \alpha(e^x - 1))$ | No | $X$ | Yes |
| PReLU | $\max(x, \alpha_p x)$ | Yes | $X$ | Yes |
| Maxout | $\max(w_1 x + b_1, w_2 x + b_2)$ | Yes | $X$ | No |
| GReLU | $\max_{1 \le k \le K}\{\alpha_{nc}^k x_{nc} + \beta_{nc}^k\}$ | Yes | $X, A$ | Yes |

This factorization limits parameters that need to be learned by the hyperfunction from $2K \times C \times N$ to $2K \times (C + N)$. As a result, the generality of the model and the computational efficiency are improved.

**Channel-wise Parameters.** To encode the global context of the graph $\mathcal{G}$ while reducing the number of learnable parameters as much as possible, we adopt the graph diffusion filter $S^{PPR}$ to squeeze and extract the topology information from $A$ and $X$. We propagate input matrix $X$ via the fully personalized PageRank scheme:

$$E = S^{PPR} X = \alpha \left(I_n - (1 - \alpha) D^{-1/2} A D^{-1/2}\right)^{-1} X, \qquad (6)$$

with teleport (or restart) probability $\alpha \in (0, 1]$, where $D$ is the degree matrix and $D^{-1/2} A D^{-1/2}$ denotes the normalized adjacency matrix $A$. We average the diffusion results and transform them into the channel-wise parameters $(\alpha_c^k, \beta_c^k)$ with the use of a linear transformation followed by a tanh function. We normalize $(\alpha_c^k, \beta_c^k)$ to be within $[-1, 1]$ range using $\tanh(\cdot)$ with the purpose of controlling the parameter scale. We define $h \in \mathbb{R}^C$ as the average of diffusion results and $P \in \mathbb{R}^{2 \times K \times C}$ as the channel-wise parameters. The channel-wise block is computed as:

$$P = \tanh(\text{MLP}(\bar{e})) \quad \text{where} \quad \bar{e} = \frac{1}{N} \sum_{n=1}^{N} e_n. \qquad (7)$$

**Node-wise parameters.** The node-wise block computes the parameter for each node. We also use the diffusion matrix to capture the graph information. To get the parameter $\gamma_n$ for each node $n$, instead of averaging the diffusion results, we directly squeeze each $e$ into one dimension using MLP. To reflect the importance of each node, softmax is applied to normalize $\gamma_n$ as:

$$\gamma_n = \frac{\exp(\text{MLP}(e_n))}{\sum_{n'=1}^{N} \exp(\text{MLP}(e_{n'}))}. \qquad (8)$$

**Plugging Step into the Updating Function.** By combining equations (8) and (7), we obtain the full parameters for GReLU via (5). We plug GReLU into the $l$-th layer updating function:

$$h_n^{(l+1)} = \sigma_n^{(l)}\left(\Theta^{(l)} [h_n^{(l)}, m_n^{(l)}]\right), \qquad (9)$$

where $\sigma_n^l(\cdot)$ is GReLU in the $l$-th layer. As illustrated in Figure 1, in the updating function, there is one set of parameters $\Theta^{(l)}$ per layer $l$, while GReLU is channel- and nose-adaptive.

**Application to spectral GNNs.** We also notice the proposed method is also applicable for Spectral GNNs not only spatial GNNs. Specifically, the graph convolution layer in the spectral domain can be
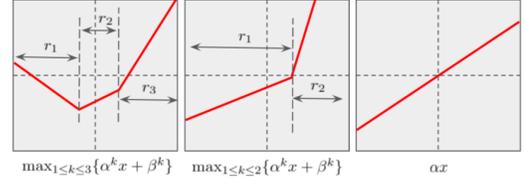


**Figure 3: Plots to the left and middle side are the piecewise functions in GReLU with $K = 3$ and $K = 2$. They cut the input value $x$ in different ranges $r_1$, $r_2$, $r_3$ and re-scale the $x$ w.r.t. $r$. The plot to the right is the weighted function used in GAT: the input is re-scaled regardless of the value range of $x$.**

written as a sum of filtered signals followed by an activation function as:

$$H^{(l+1)} = \sigma\left(\left(\sum_{i=0}^{k} \theta_i \lambda_i u_i u_i^\top\right) H^l\right) = \sigma(U g_\theta(\Lambda) U^\top H^{(l)}) \qquad (10)$$

Here, $\sigma(\cdot)$ is the activation function, $H^{(l)}$ is the hidden representation of nodes in the $l$-th layer, $U$ contains eigenvectors of the graph Laplacian $L = D - A$ (or $L = I - D^{-1/2} A D^{-1/2}$), where $D \in \mathbb{R}^{N \times N}$ is the diagonal degree matrix with entries, $\lambda_i$ are eigenvalues of $L$ and $g_\theta(\cdot)$ is the frequency filter function controlled by parameters $\theta$, where $k$ lowest frequency components are aggregated. Our method can be plugged by simply replacing the $\sigma(\cdot)$ with GReLU.

### 4.3 Understanding GReLU with Weighted Neighborhood and Attention Mechanism

Below, we analyze GReLU through the lens of Weighted Neighborhood and Attention Mechanism. To this end, we first introduce the so-call masked attention adopted in the graph attention network (GAT) [34], which uses the following form:

$$h_n' = \alpha_n h_n \quad \text{and} \quad \alpha_n = \text{MaskAttention}(H, A), \qquad (11)$$

where $\alpha_n$ is the attention coefficient assigned to node $n$. Note that (11) has the linear form of GReLU where the bias parameter and $\max(\cdot, \cdot)$ is omitted by setting $\beta = 0$ and $K = 1$. When $K = 1$, there is only one segment in the piecewise function of GReLU, the non-linear max operation reduces to a linear transformation and thus it can be removed. An illustration is shown in Figure 3. The success of GAT is due to selecting the important node information by weighting them differently during aggregation. While GAT is able to consider the difference of individual nodes, it fails to distinguish the difference between channels. GReLU solves this issue by making the parameter adaptive to both nodes and channels:

$$h_{nc}' = \max_{1 \le k \le K}\{\alpha_{nc}^k h_{nc} + \beta_{nc}^k\} \quad \text{where} \quad \{(\alpha^k, \beta^k)\}_{k=1}^K = \omega(X, A),$$
$$(12)$$

where $\alpha^k, \beta^k \in \mathbb{R}^{N \times C}$. Thus, instead of learning a single parameter shared across all channels, GReLU learns multiple weights to evaluate the importance of different channels. Moreover, adding $\max(\cdot)$ makes GReLU a piecewise linear function that cuts the input values in different ranges and rescales them differently. This further enables GReLU to consider the effect of different value ranges.

In particular, when $K = 1$, GReLU reduces to linear transformations that assign a single weight regardless of the range of input.

GReLU performs similar as attention mechanism described in global attention [23, 34]. In such a special case, these transformations are also related to Linear Graph Networks SGC [36] and S$^2$GC [43].

## 4.4 Computational Complexity.

Below, we analyze the Computational Complexity of GReLU, which is efficient as we decouple the node- and channel-wise parameters. Moreover, the computation cost of GReLU is better compared to the multi-head attention mechanism (MA) adopted in GAT.

*4.4.1 Time Complexity.* The complexity of GReLU consists of two parts: (i) the complexity of the hyperfunction $\omega_{(\mathbf{X},\mathbf{A})}$ which is used to infer the parameters of GReLU. (ii) the complexity of activation function $\sigma_{\omega_{(\mathbf{X},\mathbf{A})}}(\cdot)$ which outputs the activations.

i. Hyperfunction includes two components: (1) computation of the channel-wise parameters which consists of computing the average of diffusion output and MLP layer (with $\tanh(\cdot)$), and (2) the computation of node-wise parameters which includes the cost of diffusion and one-layer MLP with softmax. In the channel-wise block, GReLU spends $O(|\mathcal{E}|)$ in the diffusion layer and $O(2KC)$ in the MLP. In the node-wise block, GReLU spends $O(|\mathcal{E}|)$ in diffusion and $O(|\mathcal{V}|)$ in MLP with softmax.

ii. Activation function requires that GReLU is applied per channel of each node, which requires $O(C|\mathcal{V}|)$ to compute outputs.

Thus, GReLU has the complexity $O(|\mathcal{E}| + C|\mathcal{V}|)$. As the computation cost of activation function, $O(C|\mathcal{V}|)$, is the same for all rectifiers, the extra cost is mainly introduced by the computation of hyperfunction. Due to the light design of hyperfunction: 1) the averaging of diffusion output results in a single vector $\boldsymbol{e} \in \mathbb{R}^C$; 2) the MLP in node block contains only one output dimension and 3) the choice of $K$ is 2 (discussed in Section 5.4), the computation cost of hyperfunction is negligible compared to cost of GNN.

*4.4.2 Complexity of GReLU vs. Multi-head Attention Mechanism.* GReLU is faster than the multihead attention mechanism used in GAT. Note that MA needs to compute coefficient for $|\mathcal{E}|$ edges of $M$ heads, resulting in $O(M|\mathcal{E}|)$, whereas GReLU only computes coefficient for $C$ channel and $|\mathcal{V}|$ nodes, resulting in $O(|\mathcal{V}| + C)$ complexity. As $|\mathcal{V}| >> |\mathcal{E}|$ and $C$ is a relatively small number compared to both $|\mathcal{V}|$ and $|\mathcal{E}|$, adopting GReLU with a simple backbone (*e.g.*, GCN) is faster than GAT.

## 5 EXPERIMENTS

Below, we evaluate the proposed method GReLU on two tasks, node classification in Section 5.1 and graph classification in Section 5.2. We show the effectiveness of GReLU by comparing it with different ReLU variants with several GNN backbones. We show the contribution of different modules of GReLU in an ablation study in Section 5.3. The effect of choosing the hyper-parameter $K$ is discussed in Section 5.4. To verify the GReLU acts adaptively, an analysis of parameters is conducted in Section 5.5. We aim to answer the following Research Questions (RQ):

- **RQ1:** Does applying GReLU in the updating step help improve the capacity of GNNs for different downstream tasks?
- **RQ2:** What kind of factors does GReLU benefit from in its design and how do they contribute to performance?

- **RQ3:** How does the choice of parameter $K$ affect the performance of GReLU?
- **RQ4:** Does GReLU dynamically adapt to different nodes?

## 5.1 Node Classification (RQ1)

**Datasets.** GReLU is evaluated on five real-world datasets:

- Cora, PubMed, CiteSeer [16] are well-known citation network datasets, where nodes represent papers and edges represent their citations, and the nodes are labeled based on the paper topics.
- Flick [26] is an image and video hosting website, where users interact with each other via photo sharing. It is a social network where nodes represent users and edges represent their relationships, and all nodes are divided into 9 classes according to the interest groups of users.
- BlogCatalog[26] is a social network for bloggers and their social relationships from the BlogCatalog website. Node attributes are constructed by the keywords of user profiles, and the labels represent the topic categories provided by the authors, and all nodes are divided into 6 classes.

**Baselines.** We compare GReLU with different ReLU variants (as shown in Table 1) and choose several representative GNNs as backbone:

- **Chebyshev** [7] is a spectral GNN model utilizing Chebyshev filters.
- **GCN** [16] is a semi-supervised graph convolutional network model which transforms node representations into the spectral domain using the graph Fourier transform.
- **ARMA** [1] is a spectral GNN utilizing ARMA filters.
- **GraphSage** [11] is a spatial GNN model that learns node representations by aggregating information from neighbors.
- **GAT** [34] is a spatial GNN model using the attention mechanism in the aggregate function.
- **APPNP** [18] is a spectral GNN that propagates the message via a fully personalized PageRank scheme.

**Experimental Setting.** To evaluate our model and compare it with existing methods, we use the semi-supervised setting in the node classification task [16]. We use the well-studied data splitting schema in the previous work [16, 34] to reduce the risk of overfitting. To get more general results, in Cora, CiteSeer and Pubmed, we do not employ the public train and test split [16]. Instead, for each experiment, we **randomly** choose 20 nodes from each class for training and randomly select 1000 nodes for testing. All baselines are initialized with the parameters suggested by their papers, and we also further carefully tune the parameters to get the optimal performance. For our model, we choose $K = 2$ for the whole experiment. A detailed analysis of the effect of $K$ is in Section 5.4. We use a 2-layer GNN with 16 hidden units for the citation networks (Cora, PubMed, CiteSeer) and 64 units for social networks (Flick, BlogCatalog). In both networks, dropout rate 0.5 and L2-norm regularization are exploited to prevent overfitting. For each combination of ReLUs and GNNs, we run the experiments 10 times with random partitions and report the average results and the best runs. We use classification accuracy to evaluate the performance of different models.

**Table 2: The average node classification accuracy (over 10 runs). The training set and test set are randomly divided. The results in parentheses are the best results achieved among the 10 experiments. We omit the results of Cheb in BlogCatalog and Flickr dataset due to the out-of-memory issue during the training step.**

| Dataset | Model | ReLU | LReLU | ELU | PReLU | Maxout | GReLU |
|---------|-------|------|-------|-----|-------|--------|-------|
| Cora | Cheb | 76.5 ± 2.2(79.4) | 75.5 ± 2.2(80.0) | 74.9 ± 2.9(80.4) | 75.6 ± 3.0(80.6) | 75.4 ± 2.0(77.8) | **78.3 ± 1.7**(79.4) |
|  | GCN | 79.2 ± 1.4(81.6) | 79.4 ± 1.6(80.6) | 78.8 ± 1.3(80.2) | 78.9 ± 1.0(80.2) | 79.8 ± 1.5(81.5) | **81.8 ± 1.8**(83.0) |
|  | SAGE | 78.6 ± 1.7(81.3) | 77.6 ± 1.8(79.9) | 77.5 ± 1.4(79.9) | 77.6 ± 1.7(78.7) | 77.7 ± 1.7(79.9) | **79.3 ± 1.5**(80.7) |
|  | GAT | 81.2 ± 1.3(83.3) | 80.8 ± 2.3(83.1) | 80.5 ± 1.7(83.0) | 80.4 ± 2.1(83.1) | 80.2 ± 1.7(82.5) | **81.5 ± 2.1**(84.1) |
|  | ARMA | 79.0 ± 1.4(80.8) | 79.2 ± 2.1(82.9) | 79.7 ± 0.8(80.6) | 79.3 ± 2.0(81.9) | 79.8 ± 1.3(81.3) | **80.1 ± 1.5**(82.6) |
|  | APPNP | 82.0 ± 1.3(83.1) | 81.0 ± 2.3(83.1) | 81.5 ± 1.7(83.4) | 81.4 ± 2.1(83.4) | 81.0 ± 1.7(82.5) | **82.5 ± 2.1**(84.7) |
| PuMed | Cheb | 68.1 ± 4.1(74.9) | 70.5 ± 3.0(76.7) | 68.4 ± 2.8(72.7) | 71.5 ± 3.1(76.3) | 71.8 ± 3.9(75.8) | **73.4 ± 2.9**(75.6) |
|  | GCN | 77.6 ± 2.2(81.6) | 76.8 ± 1.6(79.4) | 76.8 ± 2.2(80.5) | 77.3 ± 3.7(82.5) | 77.3 ± 2.9(80.6) | **78.9 ± 1.7**(81.3) |
|  | SAGE | 75.7 ± 3.1(79.8) | 75.3 ± 3.3(79.8) | 74.8 ± 2.7(78.6) | 76.0 ± 2.5(80.5) | 74.5 ± 2.7(77.4) | **76.2 ± 1.6**(78.4) |
|  | GAT | 77.2 ± 3.1(81.3) | 78.7 ± 1.7(81.2) | 78.9 ± 2.3(82.1) | 76.2 ± 3.0(80.7) | 77.9 ± 1.7(80.3) | **79.1 ± 1.8**(80.6) |
|  | ARMA | 76.9 ± 2.6(80.7) | 76.5 ± 1.9(80.3) | 77.3 ± 2.5(80.4) | 76.5 ± 2.4(80.7) | 76.6 ± 2.9(81.3) | **77.4 ± 3.0**(80.2) |
|  | APPNP | 78.2 ± 3.1(82.3) | 78.2 ± 1.7(81.7) | 79.0 ± 2.3(82.1) | 79.2 ± 3.0(80.7) | 78.9 ± 1.7(81.3) | **79.8 ± 1.8**(81.2) |
| CiteSeer | Cheb | 67.8 ± 1.8(71.0) | 67.1 ± 2.9(71.1) | 67.8 ± 2.4(71.1) | 67.0 ± 2.3(70.5) | 67.4 ± 1.5(70.7) | **68.1 ± 1.3**(70.4) |
|  | GCN | 67.7 ± 2.3(72.1) | 68.4 ± 1.8(71.2) | 68.3 ± 1.4(70.1) | 67.3 ± 2.1(70.7) | 68.5 ± 2.2(72.5) | **68.5 ± 1.9**(71.7) |
|  | SAGE | 67.1 ± 2.8(70.1) | 67.3 ± 2.1(70.1) | 67.8 ± 1.7(70.2) | 66.2 ± 2.6(69.6) | 67.5 ± 1.8(71.5) | **68.0 ± 1.3**(69.7) |
|  | GAT | 68.6 ± 1.4(70.8) | 69.2 ± 1.9(71.7) | 68.4 ± 1.6(71.2) | 68.2 ± 1.6(69.7) | 68.6 ± 1.6(71.3) | **69.3 ± 1.7**(71.9) |
|  | ARMA | 67.7 ± 1.3(68.9) | 68.6 ± 2.4(71.5) | 67.9 ± 2.1(71.3) | 66.8 ± 1.5(69.4) | 68.5 ± 1.8(70.9) | **69.0 ± 2.2**(71.7) |
|  | APPNP | 68.7 ± 1.3(70.5) | 69.3 ± 1.6(71.2) | 69.4 ± 1.4(71.2) | 69.5 ± 1.6(70.7) | 69.2 ± 1.6(72.0) | **70.0 ± 1.7**(72.3) |
| BlogCatalog | GCN | 72.1 ± 1.9(75.5) | 72.6 ± 2.1(75.2) | 72.6 ± 1.8(75.3) | 71.4 ± 2.1(74.5) | 72.4 ± 1.4(75.3) | **73.7 ± 1.2**(74.2) |
|  | SAGE | 71.9 ± 1.3(76.0) | 72.2 ± 1.9(74.9) | 72.1 ± 1.8(74.3) | 72.0 ± 2.1(74.4) | 71.6 ± 1.4(74.3) | **73.3 ± 1.6**(75.3) |
|  | GAT | 41.7 ± 7.3(54.4) | 38.2 ± 3.3(41.8) | 46.6 ± 3.6(51.7) | 67.2 ± 2.6(71.8) | 54.2 ± 3.9(59.3) | **67.8 ± 3.9(72.3)** |
|  | ARMA | 72.5 ± 3.3(79.1) | 72.5 ± 5.9(78.7) | 77.2 ± 2.2(79.2) | 79.6 ± 3.0(84.5) | 84.4 ± 1.8(86.9) | **85.7 ± 2.7**(88.4) |
|  | APPNP | 71.1 ± 1.8(75.3) | 72.5 ± 2.0(75.2) | 72.4 ± 1.8(75.3) | 71.7 ± 2.1(74.8) | 72.8 ± 1.4(75.3) | **73.8 ± 1.5**(74.8) |
| Flickr | GCN | 50.7 ± 2.3(54.8) | 51.0 ± 2.0(53.8) | 52.8 ± 1.8(56.0) | 53.0 ± 1.6(54.9) | 54.0 ± 1.8(56.8) | **54.4 ± 1.6**(56.8) |
|  | SAGE | 49.7 ± 2.2(53.8) | 50.8 ± 2.1(54.0) | 52.6 ± 1.8(56.7) | 53.2 ± 1.3(56.1) | 54.0 ± 1.7(56.5) | **55.3 ± 1.4**(57.2) |
|  | GAT | 20.2 ± 3.1(25.2) | 20.3 ± 2.5(23.8) | 23.8 ± 2.9(28.1) | 32.8 ± 4.9(43.2) | 30.0 ± 2.6(34.6) | **33.7 ± 3.1**(36.3) |
|  | ARMA | 48.4 ± 4.7(56.1) | 52.0 ± 3.9(59.0) | 52.2 ± 4.1(58.1) | 45.9 ± 3.5(53.1) | 52.7 ± 4.0(59.4) | **56.5 ± 2.4**(59.1) |
|  | APPNP | 51.6 ± 2.0(54.0) | 51.6 ± 2.1(53.8) | 53.0 ± 1.8(56.0) | 53.2 ± 1.4(54.9) | 53.9 ± 1.9(56.6) | **54.8 ± 1.9**(56.7) |

**Observations.** The node classification results are reported in Table 2. We have the following observations:

- Compared with all baselines, the proposed GReLU generally achieves the best performance on all datasets in different GNNs. Especially, Compared with ReLU, GReLU achieves maximum improvements of 13.2% on BlogCatalog and 8.1% on Flickr using ARMA backbone. The results demonstrate the effectiveness of GReLU.
- We particularly notice that GReLU with a simple backbone (*i.e.*, GCN) outperforms the GAT with other ReLUs. This indicates that the mask edge attention adopted in GAT is not the optimal solution to weighting the information from the neighbors. As GReLU is both topology- and feature-aware, it can be regarded as a more effective and efficient way to weigh the information over both nodes and channels.

## 5.2 Graph Classification( RQ1)

**Dataset.** In the graph classification task, our proposed GReLU is evaluated on two real-world datasets MUTAG [40] and PROTEINS [40] which are common for graph classification tasks.

**Baselines.** Similarly to the node classification task, we compare GReLU with ReLU variants (Table 1) with various GNN-based pooing methods:

- GCN [16] is a multi-layer GCN model followed by mean pooling to produce the graph-level representation.
- GraphSage [11] is a message-passing model. The max aggregator is adopted to obtain the final graph representation.
- GIN [39] is a graph neural network with MLP, which is designed for graph classification.

**Experiment Setting.** The graph classification is evaluated on two datasets via 10-fold cross-validation. We reserved one fold as a test set and randomly sampled a validation set from the rest of the nine folds. We use the first eight folds to train the models. Hyper-parameter selection is performed for the number of hidden units {16, 32, 64, 128} and the number of layers {2, 3, 4, 5} w.r.t. the validation set. Similarly to the node classification task, we choose $p = 0.01$ for LeakyReLU and $K = 2$ for our model. We report the accuracy with the optimal parameters setting of each model. The classification accuracy is shown in Table 5.

**Observations.** We observe that GReLU achieves the best results on both datasets, indicating that the graph-level task can benefit

**Table 3: Ablation study: different GReLU variants (ABCD) evaluated on the Flickr dataset, the GCN backbone is used.**

| Model | Type | K | Hyperfunction | Activation | Classification Accuracy |
|---|---|---|---|---|---|
| GReLU-A | Non-linear | $K = 2$ | $\alpha_{nc}^k, \beta_{nc}^k \leftarrow \omega(X, A)$ | $\max_{1 \leq k \leq K} \{\alpha_{nc}^k x_{nc} + \beta_{nc}^k\}$ | $54.4 \pm 1.6$ |
| GReLU-B | Non-linear | $K = 2$ | $\alpha_{nc}^k, \beta_{nc}^k \leftarrow \omega(X)$ | $\max_{1 \leq k \leq K} \{\alpha_{nc}^k x_{nc} + \beta_{nc}^k\}$ | $53.4 \pm 1.8$ |
| GReLU-C | Linear | $K = 1$ | $\alpha_{nc}^k, \beta_{nc}^k \leftarrow \omega(X, A)$ | $\alpha_{nc}^k x_{nc} + \beta_{nc}^k$ | $53.7 \pm 1.4$ |
| GReLU-D | Linear | $K = 1$ | $\alpha_{nc}^k, \beta_{nc}^k \leftarrow \omega(X)$ | $\alpha_{nc}^k x_{nc} + \beta_{nc}^k$ | $52.2 \pm 1.6$ |
| ReLU | Non-linear | / | / | $\max\{x_n, 0\}$ | $50.7 \pm 2.3$ |
| SGC | Linear | / | / | $x_n$ | $48.6 \pm 1.3$ |

**Table 4: Ablation study: different GReLU variants (AEFG) are evaluated on the Flickr dataset, the GCN backbone is used.**

| Model | Parameter | Hyperfunction | Activation | Classification Accuracy |
|---|---|---|---|---|
| GReLU-A | Node/Channel-wise | $\alpha_{nc}^k, \beta_{nc}^k \leftarrow \omega(X, A)$ | $\max_{1 \leq k \leq K} \{\alpha_{nc}^k x_n + \beta_{nc}^k\}$ | $54.4 \pm 1.6$ |
| GReLU-E | Node/Channel-wise | $\alpha_{nc}^k \leftarrow \omega(X, A)$ | $\max_{1 \leq k \leq K} \{\alpha_{nc}^k x_{nc}\}$ | $54.0 \pm 1.4$ |
| GReLU-F | Channel-wise | $\alpha_c^k, \beta_c^k \leftarrow \omega(X, A)$ | $\max_{1 \leq k \leq K} \{\alpha_c^k x_{nc} + \beta_c^k\}$ | $53.0 \pm 1.9$ |
| GReLU-G | Node-wise | $\gamma_n \leftarrow \omega(X, A)$ | $\max_{1 \leq k \leq K} \{\gamma_n^k x_{n,c}\}$ | $53.8 \pm 2.1$ |

**Table 5: Graph classification results.**

| Dataset | Model | ReLU | LReLU | ELU | PReLU | Maxout | GReLU |
|---|---|---|---|---|---|---|---|
| PROTEINS | GCN | $76.0 \pm 3.2$ | $75.1 \pm 2.2$ | $76.3 \pm 3.0$ | $76.3 \pm 3.2$ | $76.5 \pm 3.2$ | $\mathbf{76.7 \pm 2.8}$ |
| | SAGE | $75.9 \pm 3.2$ | $75.7 \pm 1.2$ | $76.5 \pm 3.4$ | $75.9 \pm 3.4$ | $76.6 \pm 3.2$ | $\mathbf{76.9 \pm 1.6}$ |
| | GIN | $76.2 \pm 2.8$ | $76.1 \pm 2.9$ | $77.4 \pm 2.6$ | $76.4 \pm 3.8$ | $76.8 \pm 2.7$ | $\mathbf{77.8 \pm 3.1}$ |
| MUTAG | GCN | $85.6 \pm 5.8$ | $86.2 \pm 6.8$ | $85.7 \pm 5.5$ | $86.6 \pm 5.3$ | $85.9 \pm 5.8$ | $\mathbf{87.2 \pm 7.0}$ |
| | SAGE | $85.1 \pm 7.6$ | $84.5 \pm 7.7$ | $85.1 \pm 7.2$ | $85.7 \pm 7.4$ | $86.0 \pm 7.6$ | $\mathbf{86.7 \pm 6.4}$ |
| | GIN | $89.0 \pm 6.0$ | $89.2 \pm 6.2$ | $88.9 \pm 5.8$ | $89.2 \pm 6.1$ | $88.7 \pm 6.0$ | $\mathbf{89.5 \pm 7.4}$ |

from the global information extracted by GReLU. Compared to ReLU which is the default setting for most of GNNs, GReLU obtains the maximum improvement of 1.7% with GCN on PROTEINS. For other backbones, GReLU also improves the performance. Note that the non-linearity choice is orthogonal to the choice of pooling *i.e.*, combining GReLU with high-order pooling [19] may yield further improvements on the MUTAG and PROTEINS datasets.

## 5.3 Ablation Study (RQ2)

Below, we perform the ablation study (Flick dataset) on GReLU by removing different modules to validate the effectiveness of incorporating the graph topology into a non-linear activation function and the design of hyperfunction of GReLU. To further show the difference between GReLU and previous works, we also compare the variants with SGC. We define four variants as:

- GReLU-A: GReLU with $K = 2$ which is a non-linear function that rescales the input value in two different ranges.
- GReLU-B: No topology information due to omitting the adjacency matrix $A$ in the hyperfunction.
- GReLU-C: No non-linearity from GReLU-A due to setting $K = 1$. In this case, GReLU-C is a linear transformation.
- GReLU-D: Both non-linearity and the topology information are removed.

Table 3 shows the difference of GReLU variants and the classification accuracy evaluated on the Flickr dataset. The results of GReLU-A and GReLU-D are 1% and 1.5% better than GReLU-B and GReLU-D, indicating the topology information introduced by the adjacency matrix $A$ are useful in both linear and non-linear settings. GReLU-A outperforms the GReLU-C, verifying the necessity of adding the non-linearity. Note that GReLU-B and ReLU use adaptive and non-adaptive rectifiers, respectively. GReLU-B boosts the classification accuracy from 50.7% to 53.4% which implies the adaptive property plays a vital role in GReLU.

Based on these observations, we conclude that the major gain stems from three parts: introducing the topological structure, adding non-linearity, and making the rectifier be adaptive. Also, as SGC is the GNN without non-linear activation, it is somewhat related to the GReLU with $K = 1$. The major difference between them lies in the GReLU adopting the topology information to reweight the input $x$ using learnable parameters, whereas SGC uses an identical mapping. We compare SGC with GReLU-C and GReLU-D. Both GReLU variants with or without the adjacency matrix $A$ are better than the SGC, which again strengthens our conclusion.

We also study the need for the bias parameter $\beta$ and the use of node- and channel-wise functions. We further define three variants:
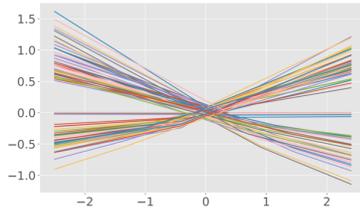
- GReLU-E: The parameters $\beta$ are dropped.

**Figure 4: The visualization of the GReLU function with $K = 2$. Each line in different colors denotes the piecewise function applied on a single channel of each node.**
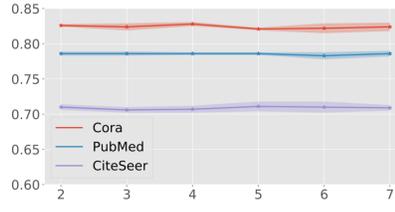
**Figure 5: The effect of different $K$ evaluated on Cora, PubMed, and CiteSeer. $x$ axis denotes the number of segments $K$ used in GReLU. $y$ axis denotes the classification accuracy.**
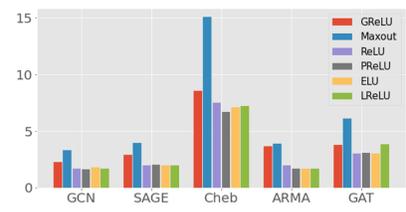
**Figure 6: Runtimes of different models with various ReLU variants on CiteSeer. y-axis denotes the runtime of training 200 epochs in seconds.**

- GReLU-F: The parameters of $\alpha$ and $\beta$ are channel-wise due to removing the node-wise block. All nodes share the same set of parameters.
- GReLU-G: The parameters of $\alpha$ and $\beta$ are node-wise due to removing the channel-wise block. All channels share the same set of parameters.

The results are shown in Table 4. We have observed different degrees of performance degradation by ablating different modules, which confirms the effectiveness of our design.

### 5.4 The Effect of Parameter $K$ (RQ3)

Note that $K$, which determines the number of segments in the non-linear function, needs to be predefined for GReLU. In this section, we discuss the effect of $K$ in the GReLU function $\max_{1 \le x \le K}\{\alpha_{nc}^k x_{nc} + \beta_{nc}^k\}$. We adopt GCN as the backbone model. To show the effect of $K$, we evaluate GReLU on the node classification task w.r.t. $K$. We choose $K \in \{2, 3, 4, 5, 6, 7\}$. Note that $K = 2$ is the default setting for the experiments. Figure 5 shows the fluctuation w.r.t. different $K$ is small. There is no obvious upward trend and a downward trend for all three datasets. The choice of $K$ is not a significant factor that influences performance. For the sake of saving computations, $K = 2$ is recommended.

### 5.5 Inspecting GReLU (RQ4)

We check if GReLU is adaptive by examining the input and output over different node attributes. We randomly choose 1,000 sets of parameters from GReLU on the node classification task with the CiteSeer dataset. We visualize the piecewise functions of the GReLU in Figure 4 which shows that they differ, which indicates GReLU is adaptive to different inputs. Some of the piecewise functions of GReLU are monotonically increasing, indicating both $\alpha_1$ and $\alpha_2$ are positive, which is consistent with other ReLU variants like PReLU and LeakyReLU. We observe that there are functions that are horizontal lines ($y = 0$) and lines with small slopes in GReLU. This indicates the parameters ($\alpha_1, \alpha_2, \beta_1, \beta_2$) for a particular channel of a node are zero. As a consequence, there is no activation for this channel $c$ in node $n$. This shows that GReLU can work as a sparse selector that filters nodes and channels. Moreover, there are monotonically decreasing piecewise functions indicating $\alpha_1$ and $\alpha_2$ are negative. This is interesting since it will flip the sign and take a controversial effect on its input. Such an effect is similar to

the negative part of the concatenated ReLU (CReLU)[29] defined as $\text{CReLU}(x) = [\text{ReLU}(x), \text{ReLU}(-x)]$. CReLU simply makes an identical copy of the linear responses, negates them, concatenates both parts of the activation, and then applies ReLU on both parts. The extra performance is gained by taking ReLU of the negative input $-x$ which is similar to setting $\alpha < 0$.

### 5.6 Runtime Comparisons

Below, runtimes of training different ReLU variants with 200 epochs on CiteSeer are reported. As shown in Figure 6, the runtime of the GReLU is shorter than that of the Maxout and slightly larger compared to other variants. The extra runtime is mainly caused by the hyperfunction that infers the parameters of GReLU, whereas other baselines (ReLU, leakyReLU, ELU, PReLU) do not use any hyperfunction. As discussed in Section 4, the hyperfunction of GReLU is simple. GReLU incurs a negligible computation cost.

## 6 CONCLUSIONS

We have proposed a topology adaptive activation function, GReLU, which incorporates the topological information of the graph and incurs a negligible computational cost. The model complexity is controlled by decoupling the parameters from nodes and channels. By making the activation function adaptive to nodes in GNN, the updating functions are more diverse than the typical counterparts in the message passing frameworks. Therefore, GReLU improves the capacity of the GNNs. Although GReLU is designed for the use in the spatial domain with message passing frameworks, GReLU can be regarded as a plug-and-play component in both spatial and spectral GNNs.

# REFERENCES

[1] Filippo Maria Bianchi, Daniele Grattarola, Lorenzo Livi, and Cesare Alippi. 2019. Graph neural networks with convolutional arma filters. *arXiv preprint arXiv:1901.01343* (2019).

[2] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. 2020. Dynamic ReLU. *arXiv preprint arXiv:2003.10027* (2020).

[3] Yankai Chen, Menglin Yang, Yingxue Zhang, Mengchen Zhao, Ziqiao Meng, Jianye Hao, and Irwin King. 2022. Modeling Scale-free Graphs with Hyperbolic Geometry for Knowledge-aware Recommendation. In *The Fifteenth ACM International Conference on Web Search and Data Mining*.

[4] Yankai Chen, Yaming Yang, Yujing Wang, Jing Bai, Xiangchen Song, and Irwin King. 2022. Attentive Knowledge-aware Graph Convolutional Networks with Collaborative Guidance for Personalized Recommendation. In *The 38th IEEE International Conference on Data Engineering*.

[5] Yankai Chen, Jie Zhang, Yixiang Fang, Xin Cao, and Irwin King. 2020. Efficient community search over large directed graphs: An augmented index-based approach. In *IJCAI*. 3544–3550.

[6] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. 2015. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289* (2015).

[7] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems* 29 (2016), 3844–3852.

[8] Charles Dugas, Yoshua Bengio, François Bélisle, Claude Nadeau, and René Garcia. 2000. Incorporating second-order functional knowledge for better option pricing. *Advances in neural information processing systems* 13 (2000), 472–478.

[9] Xinyu Fu, Jiani Zhang, Ziqiao Meng, and Irwin King. 2020. MAGNN: metapath aggregated graph neural network for heterogeneous graph embedding. In *Proceedings of The Web Conference 2020*. 2331–2341.

[10] Ian Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. 2013. Maxout networks. In *International conference on machine learning*. PMLR, 1319–1327.

[11] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *Advances in neural information processing systems*. 1024–1034.

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*. 1026–1034.

[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[14] Kevin Jarrett, Koray Kavukcuoglu, Marc'Aurelio Ranzato, and Yann LeCun. 2009. What is the best multi-stage architecture for object recognition?. In *2009 IEEE 12th international conference on computer vision*. IEEE, 2146–2153.

[15] Dongkwan Kim and Alice H. Oh. 2021. How to Find Your Friendly Neighborhood: Graph Attention Design with Self-Supervision. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. https://openreview.net/forum?id=Wi5KUNlqWty

[16] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).

[17] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. 2017. Self-normalizing neural networks. *Advances in neural information processing systems* 30 (2017), 971–980.

[18] Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. 2018. Predict then propagate: Graph neural networks meet personalized pagerank. *arXiv preprint arXiv:1810.05997* (2018).

[19] Piotr Koniusz and Hongguang Zhang. 2022. Power Normalizations in Fine-Grained Image, Few-Shot Image and Graph Classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 2 (2022), 591–609. https://doi.org/10.1109/TPAMI.2021.3107164

[20] Piotr Koniusz, Hongguang Zhang, and Fatih Porikli. 2018. A Deeper Look at Power Normalizations. In *CVPR*. 5774–5783.

[21] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature* 521, 7553 (2015), 436–444.

[22] Dongho Lee, Byungkook Oh, Seungmin Seo, and Kyong-Ho Lee. 2020. News Recommendation with Topic-Enriched Knowledge Graphs. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 695–704.

[23] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. 2015. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493* (2015).

[24] Nicholas Lim, Bryan Hooi, See-Kiong Ng, Xueou Wang, Yong Liang Goh, Renrong Weng, and Jagannadan Varadarajan. 2020. STP-UDGAT: Spatial-Temporal-Preference User Dimensional Graph Attention Network for Next POI Recommendation. In *Proceedings of the 29th ACM International Conference on Information &*

[25] J. Mairal, P. Koniusz, Z. Harchaoui, and C. Schmid. 2014. Convolutional Kernel Networks. *NIPS* (2014).

[26] Zaiqiao Meng, Shangsong Liang, Hongyan Bao, and Xiangliang Zhang. 2019. Co-embedding attributed networks. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*. 393–401.

[27] Diganta Misra. 2019. Mish: A self regularized non-monotonic neural activation function. *arXiv preprint arXiv:1908.08681* (2019).

[28] Vinod Nair and Geoffrey E Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *ICML*.

[29] Wenling Shang, Kihyuk Sohn, Diogo Almeida, and Honglak Lee. 2016. Understanding and improving convolutional neural networks via concatenated rectified linear units. In *international conference on machine learning*. PMLR, 2217–2225.

[30] Zixing Song, Ziqiao Meng, Yifei Zhang, and Irwin King. 2021. Semi-supervised Multi-label Learning for Graph-structured Data. In *CIKM '21: The 30th ACM International Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021*, Gianluca Demartini, Guido Zuccon, J. Shane Culpepper, Zi Huang, and Hanghang Tong (Eds.). ACM, 1723–1733. https://doi.org/10.1145/3459637.3482391

[31] Zixing Song, Xiangli Yang, Zenglin Xu, and Irwin King. 2021. Graph-based Semi-supervised Learning: A Comprehensive Review. *CoRR* abs/2102.13303 (2021). arXiv:2102.13303 https://arxiv.org/abs/2102.13303

[32] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research* 15, 1 (2014), 1929–1958.

[33] Ke Sun, Piotr Koniusz, and Zhen Wang. 2020. Fisher-Bures Adversary Graph Convolutional Networks. In *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference (Proceedings of Machine Learning Research, Vol. 115)*, Ryan P. Adams and Vibhav Gogate (Eds.). PMLR, 465–475. http://proceedings.mlr.press/v115/sun20a.html

[34] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903* (2017).

[35] Yaqing Wang, Fenglong Ma, and Jing Gao. 2020. Efficient Knowledge Graph Validation via Cross-Graph Representation Learning. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 1595–1604.

[36] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. In *International conference on machine learning*. PMLR, 6861–6871.

[37] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu. 2021. A Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems* 32, 1 (2021), 4–24. https://doi.org/10.1109/TNNLS.2020.2978386

[38] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. 2015. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853* (2015).

[39] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2018. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826* (2018).

[40] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How Powerful are Graph Neural Networks?. In *International Conference on Learning Representations*. https://openreview.net/forum?id=ryGs6iA5Km

[41] Menglin Yang, Ziqiao Meng, and Irwin King. 2020. FeatureNorm: L2 Feature Normalization for Dynamic Graph Embedding. In *2020 IEEE International Conference on Data Mining (ICDM)*. IEEE, 731–740.

[42] Menglin Yang, Min Zhou, Marcus Kalander, Zengfeng Huang, and Irwin King. 2021. Discrete-time Temporal Network Embedding via Implicit Hierarchical Learning in Hyperbolic Space. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 1975–1985.

[43] Hao Zhu and Piotr Koniusz. 2020. Simple spectral graph convolution. In *International Conference on Learning Representations*.

[44] Hao Zhu and Piotr Koniusz. 2021. Refine: Random range finder for network embedding. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 3682–3686.

[45] Hao Zhu, Ke Sun, and Peter Koniusz. 2021. Contrastive laplacian eigenmaps. *Advances in Neural Information Processing Systems* 34 (2021).

[46] Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. 2020. Beyond Homophily in Graph Neural Networks: Current Limitations and Effective Designs. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). https://proceedings.neurips.cc/paper/2020/hash/58ae23d878a47004366189884c2f8440-Abstract.html