

Empirical Analysis of Performance Bottlenecks in Graph Neural Network Training with GPUs

Zhaokang Wang, Yunpan Wang, Chunfeng Yuan, Rong Gu*, Yihua Huang*

*State Key Laboratory for Novel Software Technology,
Department of Computer Science and Technology, Nanjing University,
Nanjing 210023, China*

Abstract

The graph neural network (GNN) has become a popular research area for its state-of-the-art performance in many graph analysis tasks. Recently, various graph neural network libraries have emerged. They make the development of GNNs convenient, but their performance on large datasets is not satisfying. In this work, we analyze the performance bottleneck in training GNN with GPUs empirically. A GNN layer can be decomposed into two parts: the vertex and the edge calculation parts. According to their computational complexity, we select four representative GNNs (GCN, GGNN, GAT, GaAN) for evaluation. We breakdown their training time and memory usage, evaluate the effects of hyper-parameters, and assess the efficiency of the sampling techniques. The experimental evaluation indicates that the edge-related calculation is the performance bottleneck for most GNNs, dominating the training time and memory usage. Future optimization can focus on it. The sampling techniques are essential for training big graphs on GPUs, but their current implementations still have room for improvement.

Keywords: graph neural network, performance bottleneck analysis, empirical evaluation, machine learning system, GPU

1. Introduction

In recent years, the graph neural network (GNN) becomes a hot research topic in the field of artificial intelligence. Many GNNs [1, 2, 3, 4, 5, 6, 7] are proposed. They can learn the representation of vertices/edges in a graph from its topology and the original feature vectors in an *end-to-end* manner. The powerful expression ability makes GNNs achieve good accuracy in not only graph analytical tasks [8, 9, 10] (like node classification and link prediction) but also computer vision tasks (like human-object interaction [11], human parsing [12], and video object segmentation [13]).

To train GNNs easily, a series of GNN libraries/systems [14, 15, 16, 17, 18] are proposed. PyG [14], NeuGraph [16], PGL [18] and DGL [15] build upon the existing deep learning frameworks (PyG on PyTorch, NeuGraph on TensorFlow, PGL on PaddlePaddle, DGL on multiple backends). They provide users with a high-level programming models (the message-passing model for PyG/PGL/DGL and the SAGA-NN model for NeuGraph) to describe the structure of a GNN. They take advantage of the common tools provided by the underlying frameworks like the automatic differentiation to simplify the development. They utilize specially optimized CUDA kernels (like kernel fusion [15] [16]) and other implementation techniques (like 2D graph partitioning [16]) to improve the speed of GNN training on GPUs.

However, what is the real performance bottleneck in the GNN training is still in doubt. Yan et al. [19] and Zhang et al. [20] experimentally analyze the architectural characteristics of the GNN *inference*. They find that the GNN inference is more cache-friendly than the traditional graph analysis tasks (like PageRank) and is suitable for GPUs. They verify the effectiveness of the kernel fusion optimization in reducing the time

*Corresponding authors with equal contribution

Email address: {wangzhaokang, wangyp}@smail.nju.edu.cn, {cfyuan, gurong, yhuang}@nju.edu.cn (Zhaokang Wang, Yunpan Wang, Chunfeng Yuan, Rong Gu*, Yihua Huang*)

of inference. Nevertheless, they only analyze the inference stage, ignoring the effects of the backpropagation during the training.

To explore the essential performance bottleneck in the GNN training, we conduct a range of experimental analysis in deep in this work. We focus on the efficiency bottleneck of GNN training instead of accuracy. We model the GNNs with the message-passing framework that decomposes a GNN layer into two parts: the vertex calculation and the edge calculation. According to the time complexity of the two parts, we classify the typical GNNs into four quadrants($\{\text{high, low}\} \text{ complexity} \times \{\text{vertex, edge}\} \text{ calculation}$). We choose GCN [1], GGNN [4], GAT [6] and GaAN [7] as representative GNNs of the four quadrants.

We implement them with PyG and evaluate their efficiency with six real-world datasets on a GPU. We identify the most time-consuming stage in the GNN training by decomposing the training time per epoch from the layer level to the operator level. We also analyze the memory usage during the training to discover the main factor that limits the data scalability of GNN training on GPUs. Finally, we evaluate whether or not the sampling techniques affect the performance bottleneck. The key findings and insights are summarized below.

- **The training time and the memory usage of a GNN layer is mainly affected by the dimensions of the input/output hidden feature vectors.** Fixing other hyper-parameters, the training time and the memory usage of a GNN layer increase linearly as the input/output dimensions.
- **The edge-related calculation is the performance bottleneck for most GNNs.** For GNNs with high edge calculation complexity, most of the training time is spent on conducting the messaging function for every edge. For GNNs with low edge calculation complexity, the message collection and aggregation consumes most of the training time.
- **The high memory usage of the edge calculation is the main factor limiting the data scalability of GNN training.** The edge calculation generates and caches many intermediate results. They are an order of magnitude larger than the dataset itself. As GPUs have limited on-chip memory, the high memory consumption prevents us from training big graphs on GNNs.
- **The sampling techniques can significantly reduce training time and memory usage.** They are essential for training big graphs on GPUs. However, the existing implementation is still inefficient. Under big batch sizes, the time spent on the sampling may exceed the time spent on the training. Under small batch sizes, the sampled graphs are also small, wasting the computing power of GPUs.

Based on the insights, we provide several potential optimization directions:

- To reduce training time, **optimizations should focus on improving the efficiency of the edge calculation.** One may consider developing optimized operators for the messaging step that is the major source of computing costs in the edge calculation. Fusing operators of the collection step, messaging function and the aggregation step together is another way to reduce the overheads in the edge calculation.
- To reduce memory usage, **optimizations should focus on reducing the intermediate results in the edge calculation.** One may consider adopting the checkpoint mechanism to cache less intermediate results during the forward phase and re-calculate the needed data on the fly during the backpropagation.
- To improve the efficiency of the sampling techniques, **one may consider overlapping the sampling on the CPU side with the training on the GPU side.** Choosing a proper batch size automatically is another potential optimization.

We hope that our analysis can help the developers of the GNN libraries/systems have a better understanding of the characteristics of GNN training and propose more targeted optimizations.

Outline. We briefly survey the typical GNNs in Section 2. We introduce our experimental setting and targets in Section 3. The experimental results are presented and analyzed in Section 4. We summarize the key findings and give out potential optimization directions in Section 5. We introduce the related work in Section 6 and finally conclude our work in Section 7.

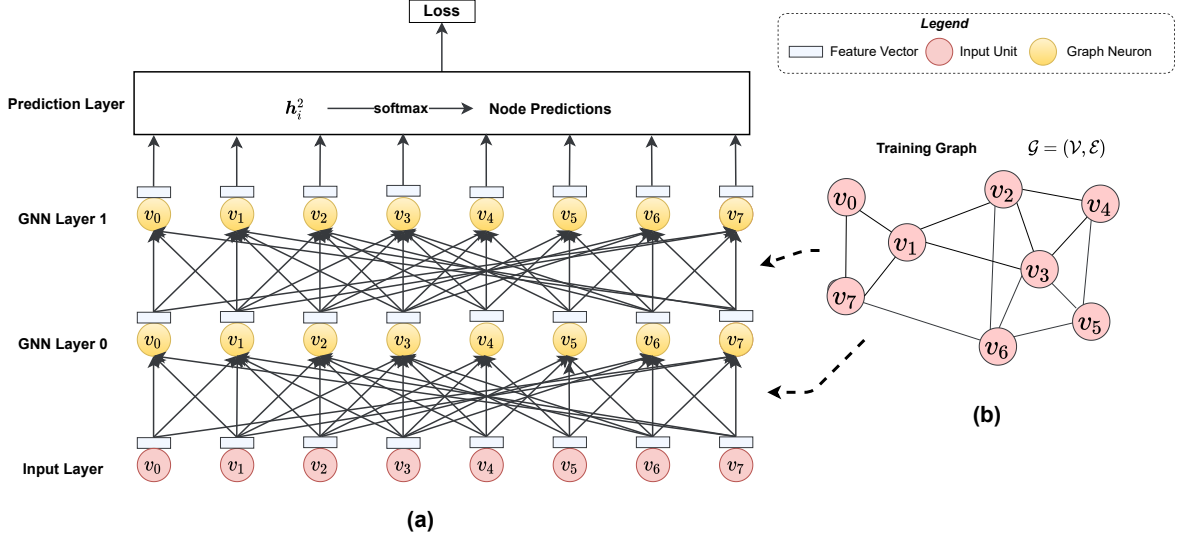


Figure 1: Structure of a typical graph neural network. (a) Demo GNN, (b) Demo graph. The target application is the node classification. The demo GNN has two GNN layers.

2. Review of Graph Neural Networks

In this section, we formally define the graph neural networks and briefly survey typical graph neural networks. We denote a simple graph \mathcal{G} as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} and \mathcal{E} are the vertex set and the edge set of \mathcal{G} , respectively. Let $n = |\mathcal{V}|$ and $m = |\mathcal{E}|$ as the number of vertices/edges. We use v_i ($0 \leq i < n$) to denote a vertex and $e_{i,j} = (v_i, v_j)$ to denote the edge pointing from v_i to v_j . The adjacency set of v_i is $\mathcal{N}(v_i) = \{v | (v_i, v) \in \mathcal{E}\}$. We denote a *vector* with a bold lower case letter like \mathbf{x} and a *matrix* with a bold upper case letter like \mathbf{X} .

2.1. Structure of Graph Neural Networks

As illustrated in Figure 1, a typical GNN can be decomposed into three parts: an input layer + several GNN layers + a prediction layer.

A GNN receives a graph \mathcal{G} as the input. Every vertex v_i in \mathcal{G} is attached with a feature vector \mathbf{x}_i to describe the properties of the vertex. The edges of \mathcal{G} may also be attached with feature vectors $\mathbf{e}_{i,j}$. The input layer of a GNN receives feature vectors from all vertices and passes them to GNN layers.

A GNN layer consists of n graph neurons, where n is the number of vertices in \mathcal{G} . Each graph neuron corresponds to a vertex in \mathcal{G} . In the first GNN layer (Layer 0), the graph neuron of the vertex v_i collects input feature vectors of itself and the vertices \mathbf{x}_j that are adjacent to v_i in \mathcal{G} (i.e., $v_j \in \mathcal{N}(v_i)$) from the input layer. The graph neuron aggregates input feature vectors, apply non-linear transformation, and outputs a hidden feature vector \mathbf{h}_i^1 for v_i . Take the demo GNN in Figure 1(a) as the example. Since $\mathcal{N}(v_3) = \{v_1, v_2, v_4, v_5, v_6\}$, the graph neuron of v_3 at the GNN layer 0 collects the feature vectors $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6\}$ from the input layer and outputs \mathbf{h}_3^1 . Different GNNs mainly differ in the graph neurons that they use. We elaborate on their details later.

Different from the traditional neural networks that layers are fully connected, the layers are sparsely connected in GNNs. The connection between the input layer and the first GNN layer is determined by the topology of \mathcal{G} . Two graph neurons are connected only if their corresponding vertices have an edge between them in \mathcal{G} . Most real-world graphs are very *sparse*, i.e. $|\mathcal{E}| \ll |\mathcal{V}|^2$.

In the next GNN layer (GNN Layer 1), the graph neuron of v_i collects the hidden feature vectors of itself \mathbf{h}_i^1 and its neighbors (\mathbf{h}_j^1 with $v_j \in \mathcal{N}(v_i)$) from the *previous* GNN layer. Based on the collected hidden vectors, the graph neuron in Layer 1 outputs a new hidden feature vector \mathbf{h}_i^2 for v_i . A GNN allows stacking more GNN layers to support deeper graph analysis. Assume there are L GNN layers in traditional. The last GNN layer (GNN Layer $L - 1$) outputs a hidden feature vector \mathbf{h}_i^L for every vertex v_i . As an embedding vector, \mathbf{h}_i^L encodes the knowledge learned from the input layer and all the previous GNN layers. Since \mathbf{h}_i^L is affected by v_i and the

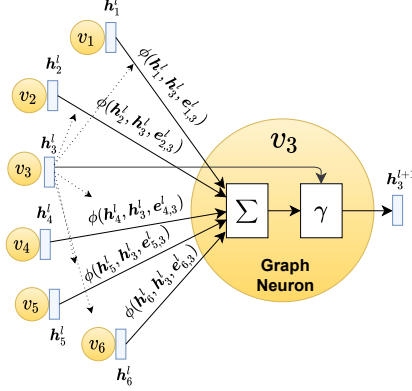


Figure 2: Graph neuron of v_3 at the GNN layer l with the graph \mathcal{G} in Figure 1b. $\phi/\Sigma/\gamma$ are the messaging/aggregation/vertex updating functions in the message-passing model, respectively.

vertices in the L -hop neighborhood of v_i , analyzing a graph with a *deeper* GNN means analyzing each vertex with a *wider* scope.

The hidden feature vectors \mathbf{h}_i^L of the last GNN layer are fed to the prediction layer to generate the output of the whole GNN. The prediction layer is a standard neural network. The structure of the prediction layer depends on the prediction task of the GNN. Take the node classification task in Figure 1 as the example. The node classification predicts a label for every vertex in \mathcal{G} . In this case, the prediction layer can be a simple softmax layer with \mathbf{h}_i^L as the input and a vector of probabilities as the output. If the prediction task is edge prediction, the hidden feature vectors of two vertices are concatenated and fed into a softmax layer. If we need to predict a label for the whole graph, a pooling (max/mean/...) layer is added to generate an embedding vector for the whole graph and the embedding vector is used to produce the final prediction.

Supporting end-to-end training is a prominent advantage of GNN, compared with traditional graph-based machine learning methods. The traditional methods need to construct feature vectors manually or from the embedding methods like DeepWalk [21] and node2vec [22]. The feature vector generation is independent from the downstream prediction task. Thus, the vectors have to be adjusted for every prediction task individually. In GNNs, the gradients are propagated from the prediction layer back to the previous GNN layers layer by layer. The model parameters in the GNN layers are updated based on the feedbacks from the downstream prediction task. In a fully parameterized way, the GNN automatically extracts an embedding vector for each vertex from its L -hop neighborhood, tuned according to the specific prediction task.

2.2. Graph Neuron and Message-passing Model

Graph neurons are building blocks of a GNN. A GNN layer consists of $|\mathcal{V}|$ graph neurons. Each vertex corresponds to a graph neuron. A graph neuron as shown in Figure 2 is a small neural network. The graph neuron of v_i at layer l receives hidden feature vectors \mathbf{h}_j^l from the graph neurons of v_i and its neighbors ($v_j \in \{v_i\} \cup \mathcal{N}(v_i)$) at the previous GNN layer¹. The graph neuron aggregates the received hidden feature vectors, applies non-linear transformations, and outputs a new hidden feature vector \mathbf{h}_i^{l+1} .

We follow the message-passing model [23] to formally define a graph neuron. The message-passing model is widely used in the cutting-edge GNN libraries like PyTorch Geometric (PyG) [14] and Deep Graph Library (DGL) [15]. Figure 2 shows the structure of a graph neuron in the message-passing model. Graph neurons at layer l are made of three *differentiable* functions: ϕ^l , Σ^l and γ^l . The graph neuron calculates the output hidden

¹For the GNN layer 0, graph neurons receive input feature vectors, i.e., $\mathbf{h}_i^0 = \mathbf{x}_i$

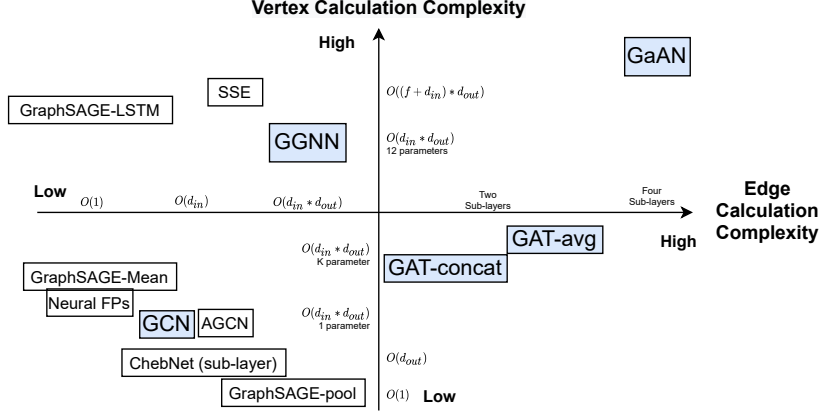


Figure 3: Complexity quadrants of typical GNNs. We compare the complexity according to the number of sub-layers, the Big-O notation and the number of parameters to train.

vector \mathbf{h}_i^{l+1} by

$$\mathbf{h}_i^{l+1} = \gamma^l(\mathbf{h}_i^l, \sum_{v_j \in \mathcal{N}(v_i)} \phi^l(\mathbf{h}_i^l, \mathbf{h}_j^l, \mathbf{e}_{j,i})).$$

ϕ^l is the *messaging* function. For every incident edge (v_j, v_i) of v_i , ϕ receives the output hidden feature vectors \mathbf{h}_i^l and \mathbf{h}_j^l from the previous GNN layer and the edge feature vector $\mathbf{e}_{j,i}$ as the input. ϕ^l emits a message vector $\mathbf{m}_{j,i}^l$ for every edge (v_j, v_i) at layer l , i.e., $\mathbf{m}_{j,i}^l = \phi^l(\mathbf{h}_i^l, \mathbf{h}_j^l, \mathbf{e}_{j,i})$. For every v_i , the message vectors $\mathbf{m}_{j,i}^l$ of its adjacent edges $(v_j \in \mathcal{N}(v_i))$ are aggregated by the *aggregation* function Σ^l to produce an aggregated vector \mathbf{s}_i^l , i.e., $\mathbf{s}_i^l = \Sigma_{v_j \in \mathcal{N}(v_i)}^l \mathbf{m}_{j,i}^l$. v_i 's aggregated vector \mathbf{s}_i^l and its hidden vector \mathbf{h}_i^l from the previous GNN layer are fed into the *vertex updating* function γ^l to calculate the output hidden vector \mathbf{h}_i^{l+1} of the current layer l , i.e., $\mathbf{h}_i^{l+1} = \gamma^l(\mathbf{h}_i^l, \mathbf{s}_i^l)$. The end-to-end training requires ϕ^l and γ^l (like multi-layer perceptrons and GRU) and Σ_l (like mean, sum, element-wise min/max) are *differentiable* to make the whole GNN differentiable.

Different GNNs have different definitions of the three functions. We regard ϕ and Σ as the *edge calculation* functions, since they are conducted over every edge in \mathcal{G} . We regard γ as the *vertex calculation* function, as it is conducted over every vertex in \mathcal{G} . Table 1 and Table 2 list the edge functions and the vertex functions of typical GNNs, respectively. For ChebNet, we report its GNN sub-layer in the tables. A ChebNet layer consists of K GNN sub-layers and a summation layer: $\mathbf{H}^{l+1} = \sum_{k=1}^K \mathbf{Z}^{(k)} \mathbf{W}^{(k)}$ with the GNN sub-layers $\mathbf{Z}^{(1)} = \mathbf{H}^l$, $\mathbf{Z}^{(2)} = \hat{\mathbf{L}} \mathbf{H}^l$, and $\mathbf{Z}^{(k)} = 2\hat{\mathbf{L}} \mathbf{Z}^{(k-1)} - \mathbf{Z}^{(k-2)}$, where \mathbf{H}^l is the matrix of output hidden feature vectors of the layer $l - 1$. For GAT, a GAT layer consists of two sub-layers and it conducts part of the vertex calculation before the two sub-layers. For GaAN, a GaAN layer consists of four sub-layers: the first sub-layer calculates the summation $\mathbf{m}_{j,i,(k)}^{l,0} = \exp((\mathbf{W}_{xa}^{l,k} \mathbf{h}_j^l + \mathbf{b}_{xa}^{l,k})^T (\mathbf{W}_{ya}^{l,k} \mathbf{h}_i^l + \mathbf{b}_{ya}^{l,k}))$, and the other three sub-layers calculate $\mathbf{m}_{j,i,1}^l / \mathbf{m}_{j,i,2}^l / \mathbf{m}_{j,i,3}^l$.

2.3. Classification of GNNs

Since we focus on analyzing the performance bottleneck in training GNNs, we classify the typical GNNs from the view of time complexity. We use $O_\phi/O_\Sigma/O_\gamma$ to denote the time complexity of the three functions in the message-passing model. The time complexity of a GNN layer is made up of two parts: the edge calculation complexity $m * (O_\phi + O_\Sigma)$ and the vertex calculation complexity $n * O_\gamma$.

In Table 1 and Table 2, we list the edge and vertex calculation complexity, respectively. The time complexity of a graph neuron is affected by the dimensions of the input/output hidden vectors d_{in} and d_{out} and the dimensions of the model parameters (like the number of heads K in GAT and the dimensions of the view vectors $d_a/d_v/d_m$ in GaAN).

We classify the typical GNNs into four quadrants based on their edge/vertex complexity as shown in Figure 3. We pick GCN, GGNN, GAT, and GaAN as the *representative GNNs* of the four quadrants.

GCN (low vertex & edge calculation complexity): Graph convolution network (GCN [1]) introduces the first-order approximation of the spectral-based graph convolutions. It has only one parameter to learn at each

GNN	Type	Σ	ϕ	Complexity
ChebNet [2]	Spectral	sum	$\mathbf{m}_{j,i}^{(k)} = e_{j,i} \mathbf{z}_j^{(k-1)}$	$O(d_{in})$
GCN [1]	Spectral	sum	$\mathbf{m}_{j,i}^l = e_{j,i} \mathbf{h}_j^l$	$O(d_{in})$
AGCN [3]	Spectral	sum	$\mathbf{m}_{j,i}^l = \tilde{e}_{j,i}^l \mathbf{h}_j^l$	$O(d_{in})$
GraphSAGE [5]	Non-spectral	mean/LSTM	$\mathbf{m}_{j,i}^l = \mathbf{h}_j^l$	$O(1)$
GraphSAGE-pool [5]	Non-spectral	max	$\mathbf{m}_{j,i}^l = \delta(\mathbf{W}_{pool}^l \mathbf{h}_j^l + \mathbf{b}^l)$	$O(d_{in} * d_{out})$
Neural FPs [24]	Non-spectral	sum	$\mathbf{m}_{j,i}^l = \mathbf{h}_j^l$	$O(1)$
SSE [25]	Recurrent	sum	$\mathbf{m}_{j,i}^l = [\mathbf{x}_j \parallel \mathbf{h}_j^l]$	$O(f + d_{in})$
GGNN [4]	Gated	sum	$\mathbf{m}_{j,i}^l = \mathbf{W}^l \mathbf{h}_j^l$	$O(d_{in} * d_{out})$
GAT [6]	Attention	sum	Sub-layer 0: $\mathbf{m}_{j,i,(k)}^{l,0} = \exp(\text{LeakyReLU}(\mathbf{a}^T [\hat{\mathbf{h}}_{i,(k)}^l \parallel \hat{\mathbf{h}}_{j,(k)}^l]))$	concat: $O(d_{out})$ average: $O(K * d_{out})$
			Sub-layer 1: $\alpha_{j,i,(k)}^l = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^T [\hat{\mathbf{h}}_{i,(k)}^l \parallel \hat{\mathbf{h}}_{j,(k)}^l]))}{\hat{\mathbf{a}}_{i,(k)}^{l,0}}$	Two sub-layers
			Multi-head concatenation : $\mathbf{m}_{j,i}^{l,1} = \parallel_{k=1}^K \alpha_{j,i,(k)}^l \hat{\mathbf{h}}_{j,(k)}^l$	
			Multi-head average : $\mathbf{m}_{j,i}^{l,1} = \frac{1}{K} \sum_{k=1}^K \alpha_{j,i,(k)}^l \hat{\mathbf{h}}_{j,(k)}^l$	
GaAN [7]	Attention	sum,max,mean	Sub-layer 0: $\mathbf{m}_{j,i,(k)}^{l,0} = \exp((\mathbf{W}_{xa}^{l,k} \mathbf{h}_j^l + \mathbf{b}_{xa}^{l,k})^T (\mathbf{W}_{ya}^{l,k} \mathbf{h}_i^l + \mathbf{b}_{ya}^{l,k}))$	$O(\max(d_a, d_m, d_v) * K * d_{in})$
			Sub-layer 1: $\alpha_{j,i,(k)}^l = \frac{\exp((\mathbf{W}_{xa}^{l,k} \mathbf{h}_j^l + \mathbf{b}_{xa}^{l,k})^T (\mathbf{W}_{ya}^{l,k} \mathbf{h}_i^l + \mathbf{b}_{ya}^{l,k}))}{\hat{\mathbf{a}}_{i,(k)}^{l,0}}$	Four sub-layers
			$\mathbf{m}_{j,i}^{l,1} = \parallel_{k=1}^K \alpha_{j,i,(k)}^l \text{LeakyReLU}(\mathbf{W}_v^{l,k} \mathbf{h}_j^l + \mathbf{b}_v^{l,k})$	
			Sub-layer 2: $\mathbf{m}_{j,i}^{l,2} = \mathbf{W}_m^l \mathbf{h}_j^l + \mathbf{b}_m^l$	
			Sub-layer 3: $\mathbf{m}_{j,i}^{l,3} = \mathbf{h}_j^l$	

Table 1: Typical graph neural networks and their edge calculation functions. d_{in} and d_{out} are dimensions of the input and output hidden feature vectors, respectively. Blue variables are model parameters to learn.

GNN	γ	Complexity
ChebNet [2]	$\mathbf{z}_i^{(k)} = 2\mathbf{s}_i^{(k)} - \mathbf{z}_i^{(k-2)}$	$O(d_{out})$
GCN [1]	$\mathbf{h}_i^{l+1} = \mathbf{W}^l \mathbf{s}_i^l$	$O(d_{in} * d_{out})$
AGCN [3]	$\mathbf{h}_i^{l+1} = \mathbf{W}^l \mathbf{s}_i^l$	$O(d_{in} * d_{out})$
GraphSAGE [5]	$\mathbf{h}_i^{l+1} = \delta(\mathbf{W}^l [\mathbf{s}_i^l \parallel \mathbf{h}_i^l])$	$O(d_{in} * d_{out})$
GraphSAGE-pool [5]	$\mathbf{h}_i^{l+1} = \mathbf{s}_i^l$	$O(1)$
Neural FPs [24]	$\mathbf{h}_i^{l+1} = \delta(\mathbf{W}^{l, \mathcal{N}(i) } (\mathbf{h}_i^l + \mathbf{s}_i^l))$	$O(d_{in} * d_{out})$
SSE [25]	$\mathbf{h}_i^{l+1} = (1 - \alpha)\mathbf{h}_i^l + \alpha\delta(\mathbf{W}_1^l \delta(\mathbf{W}_2^l [\mathbf{x}_i \parallel \mathbf{s}_i^l]))$	$O((f + d_{in}) * d_{out})$
GGNN [4]	$\mathbf{z}_i^l = \delta(\mathbf{W}^z \mathbf{s}_i^l + \mathbf{b}^{sz} + \mathbf{U}^z \mathbf{h}_i^l + \mathbf{b}^{hz})$ $\mathbf{r}_i^l = \delta(\mathbf{W}^r \mathbf{s}_i^l + \mathbf{b}^{sr} + \mathbf{U}^r \mathbf{h}_i^l + \mathbf{b}^{hr})$ $\mathbf{h}_i^{l+1} = \tanh(\mathbf{W} \mathbf{s}_i^l + \mathbf{b}^s + \mathbf{U}(\mathbf{r}_i^l \odot \mathbf{h}_i^l) + \mathbf{b}^h)$ $\mathbf{h}_i^{l+1} = (1 - \mathbf{z}_i^l) \odot \mathbf{h}_i^l + \mathbf{z}_i^l \odot \mathbf{h}_i^{l+1}$	$O(d_{in} * d_{out})$
GAT [6]	Preprocessing:	concat: $O(d_{in} * d_{out})$
	$\hat{\mathbf{h}}_{i,(k)}^l = \mathbf{W}^{l,k} \mathbf{h}_i^l$	average: $O(K * d_{in} * d_{out})$
	Sub-layer 0:	Two sub-layers
	$\hat{\mathbf{a}}_{i,(k)}^{l,0} = \mathbf{s}_{i,(k)}^{l,0}$	
GaAN [7]	Sub-layer 1:	
	$\mathbf{h}_i^{l+1} = \delta(\mathbf{s}_i^{l,1})$	
	Sub-layer 0:	$O(\max(K * d_v + d_{in}, 2 * d_{in} + d_m) * d_{out})$
	$\hat{\mathbf{a}}_{i,(k)}^{l,0} = \mathbf{s}_{i,(k)}^{l,0}$	Four sub-layers
	After sub-layer 0,1,2,3:	
	$\mathbf{g}_i^l = \mathbf{W}_g^l [\mathbf{h}_i^l \parallel \mathbf{s}_i^{l,2} \parallel \mathbf{s}_i^{l,3}] + \mathbf{b}_g^l$ $\mathbf{h}_i^{l+1} = \mathbf{W}_o^l [\mathbf{h}_i^l \parallel (\mathbf{g}_i^l \odot \mathbf{s}_i^{l,1})] + \mathbf{b}_o^l$	

Table 2: Typical graph neural networks and their vertex calculation functions. d_{in} and d_{out} are dimensions of the input and output hidden feature vectors, respectively. Blue variables are model parameters to learn. In Neural FPs, $\mathbf{W}^{l, |\mathcal{N}(i)|}$ is the weight matrix for vertices with degree $|\mathcal{N}(i)|$ at layer l .

layer, i.e. the weight matrix \mathbf{W}^l in γ . A GCN graph neuron can be expressed as $\mathbf{h}_i^{l+1} = \mathbf{W}^l \sum_{v_j \in \mathcal{N}(v_i)} e_{j,i} \mathbf{h}_j^l$, where $e_{j,i}$ is the normalized weight of the edge (v_j, v_i) . According to the associative law of the matrix multiplication, $\mathbf{h}_i^{l+1} = \sum_{v_j \in \mathcal{N}(v_i)} e_{j,i} \mathbf{W}^l \mathbf{h}_j^l$. Since the dimension of \mathbf{h}_i^{l+1} is usually smaller than \mathbf{h}_i^l in practical GCNs, the implementation of GCN in PyTorch Geometric chooses to first conduct the vertex calculation $\hat{\mathbf{h}}_j^l = \mathbf{W}^l \mathbf{h}_j^l$ for each vertex v_j and then conduct the edge calculation $\mathbf{h}_i^{l+1} = \sum_{v_j \in \mathcal{N}(v_i)} e_{j,i} \hat{\mathbf{h}}_j^l$. As $\hat{\mathbf{h}}_j^l$ has the same dimension as \mathbf{h}_i^{l+1} , the implementation significantly reduces the computation cost of the edge calculation.

GGNN (high vertex & low edge calculation complexity): GGNN [4] introduces the gated recurrent unit (GRU) into the graph neural networks. The vertex updating function ϕ of GGNN is a modified GRU unit that has 12 model parameters to learn, having high computational complexity. To lower the training cost, all GNN layers share the same group of parameters in GGNN. GGNN further requires the dimension of \mathbf{h}^{l+1} is equal to the dimension of \mathbf{h}^l . Since the messaging function ϕ only uses the hidden feature vector \mathbf{h}_j^l of the source vertex v_j of an edge (v_j, v_i) . In the implementation of PyG, GGNN conducts the pre-processing vertex calculation $\hat{\mathbf{h}}_i^l = \mathbf{W} \mathbf{h}_i^l$ for every vertex v_i before the message-passing. The messaging function ϕ directly uses $\hat{\mathbf{h}}_j^l$ as the message vector for the edge (v_j, v_i) . In this way, GGNN further reduces the time complexity of the edge calculation to $O(1)$ without increasing the time complexity of the vertex calculation.

GAT (low vertex & high edge calculation complexity): GAT [6] introduces the attention and multi-head mechanism into the graph neural networks. The K heads generate K independent views for an edge, where K is a hyper-parameter. The views of K heads can be merged by concatenating or by averaging. For concatenating, the dimension of the hidden feature vector of each head d_{head} is d_{out}/K . For averaging, d_{head} is d_{out} , multiplying the complexity by K . Each GAT layer consists of a preprocessing step and two sub-layers. In the preprocessing step, GAT calculates the attention vectors $\hat{\mathbf{h}}_{i,(k)}^l$ for every vertex v_i in every head k . The first sub-layer uses the attention vectors to calculate the attention weights of every edge in every head $\mathbf{m}_{j,i,(k)}^{l,0}$. After aggregation, the first sub-layer gets the summation of attention weights of every vertex in every head $\hat{\mathbf{a}}_{i,(k)}^{l,0}$. The second sub-layer calculates the normalized attention weights $\alpha_{j,i,(k)}$ for every edge in every head and aggregates the hidden feature vectors with the normalized weights in the message-passing. GAT uses the aggregated vector $\mathbf{s}_i^{l,1}$ of the second sub-layer as the output hidden vector directly.

GaAN (high vertex & edge calculation complexity): Based on the multi-head mechanism, GaAN [7] introduces a convolutional subnetwork to control the weight of each head. Each GaAN layer consists of four sub-layers. The first two sub-layers are similar to GAT, and the last two sub-layers build the convolutional subnetwork. The first sub-layer aims to get the summation of attention weights of every vertex in each head $\hat{\mathbf{a}}_{j,i,(k)}^{l,0}$. The second sub-layer aggregates the feature vectors with the normalized attention weights $\hat{\mathbf{a}}_{i,(k)}^{l,0}$. The third and fourth sub-layers aggregate feature vectors with element-wise max and element-wise mean operators separately. After the four sub-layers, the vertex updating function uses the aggregated vectors of last three sub-layers $\mathbf{s}_i^{l,1}, \mathbf{s}_i^{l,2}, \mathbf{s}_i^{l,3}$ to generate the output hidden vector \mathbf{h}_i^{l+1} .

2.4. Sampling Techniques

By default, GNNs are trained in a full-batch way, using the whole graph in each iteration. The full-batch gradient descent has two disadvantages [26]. It has to cache intermediate results of all vertices in the forward phase, which consumes lots of memory space. It updates the parameters only once for each epoch, slowing the convergence of gradient descent.

To train GNNs in a mini-batch way, the sampling techniques are proposed. In each mini-batch, they sample a small subgraph from the whole graph \mathcal{G} and uses the subgraph to update the model parameters. The sampling techniques only active the graph neurons and the connections that appear in the sampled subgraph between GNN layers, as shown in Figure 4. In active graph neurons and connections do not participate in the training of this mini-batch, saving lots of computation and storage costs. Moreover, it may reduce the risk of overfitting the training graph. The existing sampling techniques can be classified into two groups **neighbor sampling** and **graph sampling** based on whether different GNN layers sample different subgraphs.

The neighbor sampling techniques [5, 27, 28, 29, 30] sample nodes or edges layer by layer. The sampled subgraphs of different GNN layers may be *different*, as shown in Figure 4a. GraphSAGE [5] is the representative technique. The technique first samples several vertices from \mathcal{V} in the last GNN layer. Then it repeatedly samples the neighbors of those sampled vertices in the previous layer until the input layer. For every sampled vertex v_i in the GNN layer l , GraphSAGE samples at most S^l neighbors of v_i in the previous GNN layer. S^l is the

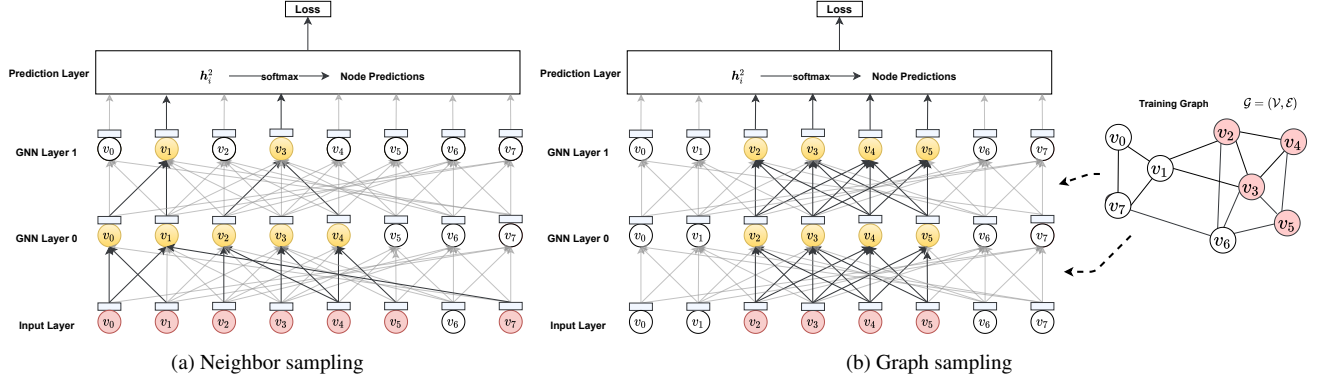


Figure 4: Training a GNN with sampling techniques. The faded graph neurons and their connections are inactivated.

hyper-parameter that is usually much smaller than $|\mathcal{V}|$. In this way, GraphSAGE limits the neighborhood sizes of the vertices in the sampled subgraph, especially high-degree vertices.

The graph sampling techniques [31, 26, 32] sample a subgraph for each mini-batch and use the same sampled graph for all GNN layers, as shown in Figure 4b. They differ in the methods to sample subgraphs. The cluster sampler technique [26] is the representative technique. Given a training graph \mathcal{G} , it partitions \mathcal{G} into dense clusters. For each mini-batch, it randomly pick K clusters to form the sampled subgraph, where K is the hyper-parameter.

3. Evaluation Design

We design a series of experiments to explore the performance bottleneck in training graph neural networks. We first introduce the experimental setting in Section 3.1 and then give out our experimental scheme in Section 3.2. The evaluation results are presented and analyzed later in Section 4.

3.1. Experimental Setting

Experimental Environment. All the experiments were conducted in a CentOS 7 Linux server with kernel version 3.10.0. The server had 40 cores and 90 GB main memory. The server was equipped with an NVIDIA Tesla T4 GPU with 16GB GDDR6 memory. For the software environment, we adopted Python 3.7.7, PyTorch 1.5.0, and CUDA 10.1. We implemented all the GNNs with PyTorch Geometric 1.5.0.

Dataset. We used six real-world graph datasets as listed in Table 3 that were popular in the GNN accuracy evaluation [33, 32, 34]. For directed graphs, PyG converts them into undirected ones during the data loading. Thus, the average degree of a directed graph $\bar{d} = \frac{2|\mathcal{E}|}{|\mathcal{V}|}$. For an undirected graph, \mathcal{E} already contains two-direction edges and $\bar{d} = \frac{|\mathcal{E}|}{|\mathcal{V}|}$. For the cam dataset, we generated random dense feature vectors. We also used random graphs generated by the R-MAT graph generator [35] in the experiments, to explore the effects of graph topological characteristics (like the average degree) on the performance bottleneck. Input feature vectors of random graphs were random dense vectors with a dimension of 32. Vertices of random graphs were classified into 10 classes randomly.

Learning Task. We used the node classification as the target task in GNNs due to its popularity in real-world applications. We trained GNNs with the semi-supervised learning setting. All vertices and their input feature vectors were used, but only parts of the vertices were attached with labels during the training and they were used to calculate the loss and gradients. The vertices with unseen labels were used in the evaluation phase to evaluate the accuracy of the current parameters.

Dataset	$ \mathcal{V} $	$ \mathcal{E} $	\bar{d}	$\dim(\mathbf{x})$	#Class	Directed
pubmed (pub) [33]	19,717	44,324	4.5	500	3	Yes
amazon-photo (amp) [34]	7,650	119,081	31.1	745	8	Yes
amazon-computers (amc) [34]	13,752	245,861	35.8	767	10	Yes
coauthor-physics (cph) [34]	34,493	247,962	14.4	8415	5	Yes
flickr (fli) [32]	89,250	899,756	10.1	500	7	No
com-amazon (cam) [36]	334,863	925,872	2.8	32	10	No

Table 3: Dataset overview. \bar{d} represents the average vertex degree. $\dim(\mathbf{x})$ is the dimension of the input feature vector.

GNN Implementation. We implemented the four typical GNNs (GCN, GGNN, GAT, GaAN) with PyTorch Geometric 1.5.0. To compare the performance characteristics of four GNNs side-by-side, we used a unified GNN structure for them: Input Layer \rightarrow GNN Layer 0 \rightarrow GNN Layer 1 \rightarrow Softmax Layer (to prediction). The structure was popular in the experimental evaluation of GCN [1], GAT [6] and GaAN [7]. Since a GGNN layer requires the input and output hidden feature vectors have the same dimension, we added two multi-layer perceptron (MLP) layers to transform the dimensions of the input/output feature vectors of the whole GNN: Input Layer \rightarrow MLP \rightarrow GGNN Layer 0 \rightarrow GGNN Layer 1 \rightarrow MLP \rightarrow Softmax Layer. We stored the dataset and the model parameters on the GPU side. All the training was conducted on the GPU.

Hyper-parameters. We use $\dim(v)$ to denote the dimension of a vector v . We picked the hyper-parameters of GNNs according to their popularity in the original papers. We used the same set of hyper-parameters for all the datasets. Some hyper-parameters like the dimensions of hidden feature vectors were common in the four GNNs and we set them to the same values. For GCN/GAT/GaAN, we set $\mathbf{h}_i^0 = \mathbf{x}_i$, $\dim(\mathbf{h}_i^1) = 64$, and $\dim(\mathbf{h}_i^2) = \#Classes$. For GAT, we set the hyper-parameters according to its paper[6]. The first GAT layer had 8 heads with the dimension of each head $d_{head}^0 = 8$. The first GAT layer merged the heads by concatenating. The second GAT layer used a single head with the dimension $d_{head}^1 = d_{out}^1 = \#Classes$. For GGNN, since it uses extra MLP layers to transform the dimensions of the input/output feature vectors, we set $\dim(\mathbf{h}_i^0) = \dim(\mathbf{h}_i^1) = \dim(\mathbf{h}_i^2) = 64$. We used 8 heads in the both GaAN layers with $d_a = d_v = 8$ and $d_m = 64$.

Sampling Techniques. We picked the neighbor sampler from GraphSAGE [5] and the cluster sampler from ClusterGCN [26] as the typical sampling techniques respectively. For the neighbor sampler, we set the neighborhood sample sizes as 25 for GNN Layer 1, and 10 for GNN Layer 0 and set the default batch size as 512 from its paper[5]. For the cluster sampler, we partitioned every input graph into 1500 partitions and used 20 partitions per batch according to the parameters used in its paper[26].

3.2. Experimental Scheme

To find out the performance bottleneck in GNN training, we conducted the experimental analysis with four questions. The answers to those questions will give us a more comprehensive view of the performance characteristics of GNN training.

Q1 *How do the hyper-parameters affect the training time and the memory usage of a GNN?* (Section 4.1)

Every GNN has a group of hyper-parameters like the number of GNN layers and the dimensions of hidden feature vectors. The hyper-parameters affect the training time per epoch and the peak memory usage during the training. To evaluate their effects, we measured how the training time per epoch and the peak memory usage (of the GPU) changed as we increased the values of the hyper-parameters. Through the experiments, we want to verify the validity of the computational complexity analysis in Table 1 and Table 2. If the complexity analysis is valid, we can analyze the bottleneck theoretically.

Q2 *Which stage is the most time-consuming stage in GNN training?* (Section 4.2)

We can decompose the training time on different levels: layer level, edge/vertex calculation level, and the basic operator level. On each level, we decomposed the training time of an epoch into several stages. The most time-consuming stage is the performance bottleneck. Optimizing its implementation can significantly reduce the training time.

Q3 Which consumes most of memory in GNN training? (Section 4.3)

The limited memory capacity of a GPU is the main factor preventing us from training GNNs on big graphs. We measured the peak memory usage of GNN training under different graph scales, input feature dimensions, and average vertex degrees. Based on the results, we analyzed which is the most memory-consuming component in a GNN. Reducing its memory usage will enable us to train GNNs on bigger graphs under the same memory capacity.

Q4 Can the sampling techniques remove the performance bottleneck in GNN training? (Section 4.4)

Theoretically, the sampling techniques can significantly reduce the number of graph neurons that participate in the training of a batch. Consequently, the training time and the memory usage should also decrease. To validate the effectiveness of the sampling techniques, we measured the training time and peak memory usage under different batch sizes. If the sampling techniques are effective, they are the keys to conduct GNN training on very big graphs. If they are not effective, we want to find out which impairs its efficiency.

4. Evaluation Results and Analysis

We answer the four questions in Section 3.2 one by one with experiments. Without otherwise mentioned, the reported training time per epoch is the average wall-clock training time of 50 epochs, excluding abnormal epochs².

4.1. Effects of Hyper-parameters on Performance

According to Table 1 and Table 2, the time complexities of ϕ and γ are linear to each hyper-parameter separately. If we increase one of the hyper-parameters and fix the others, the training time should increase linearly.

To verify the time complexity analysis in Table 1 and Table 2, we first compare the training time of the four GNNs on the real-world datasets in 3. The ranking of the training time is GaAN \gg GAT $>$ GGNN $>$ GCN in all cases. Since the real-world graphs have more edges than vertices ($m > n$), the time complexity of the edge calculation stage affects more than the vertex calculation stage. The ranking is consistent with the time complexity analysis.

To further evaluate the effects of the hyper-parameters, we measure the training time of each GNN with varying hyper-parameters in Figure 6.

Since $\dim(\mathbf{h}^0)$ and $\dim \mathbf{h}^1$ are determined by the dataset with $\dim(\mathbf{h}^0 = \dim(\mathbf{x}))$ and $\dim(\mathbf{h}^2) = \#Classes$, for GCN and GGNN, the only modifiable hyper-parameter is the hidden dimension $\dim(\mathbf{h}^1)$ with $\dim(\mathbf{h}^1) = d_{out}^0 = d_{in}^1$. According to the time complexity analysis, if we fix other hyper-parameters but only increase $\dim(\mathbf{h}^1)$, the computing costs of the GNN layer 0 and the GNN layer 1 both increase linearly with $\dim(\mathbf{h}^1)$, causing the training time of the whole GNN also increasing linearly. Figure 6a and Figure 6b show that the training time of GCN and GGNN increased linearly with $\dim(\mathbf{h}^1)$ when $\dim(\mathbf{h}^1)$ is big, consistent with the time complexity analysis.

For GAT, we modify the number of heads K and the dimension of each head d_{head} in the GAT layer 0. The dimension of the hidden feature vector is determined correspondingly as $d_{out}^0 = d_{in}^1 = \dim(\mathbf{h}^1) = K d_{head}$. Thus, the computing costs of the GAT layer 0 and the GAT layer 1 increase linearly with K and d_{head} separately. Figure 6c confirms the theoretical analysis.

For GaAN, it is also based on the multi-head mechanism. Its time complexity is affected by $\dim(\mathbf{h}^1)$ ($d_{out}^0 = d_{in}^1 = \dim(\mathbf{h}^1)$), d_a , d_v , d_m and the number of heads K . Figure 6d demonstrates that the training time increases linearly with the hyper-parameters, except for $\dim(\mathbf{h}^1)$. As $\dim(\mathbf{h}^1)$ increases, the training time increases first slightly and then linearly. We observe similar phenomena in GCN, GGNN, and GAT with low hyper-parameters. When the hyper-parameters is too low, the GNN training cannot make full use of the

²During the training of some epochs, there are extra profiling overheads from NVIDIA Nsight Systems and GC pauses from the Python interpreter that significantly increase the training time. Assume Q1 and Q3 are the 25% and 75% quantiles of the training time of 50 epochs, respectively. We regard the epochs with the training time outside the range of $[Q1 - 1.5 * (Q3 - Q1), Q3 + 1.5 * (Q3 - Q1)]$ as abnormal epochs.

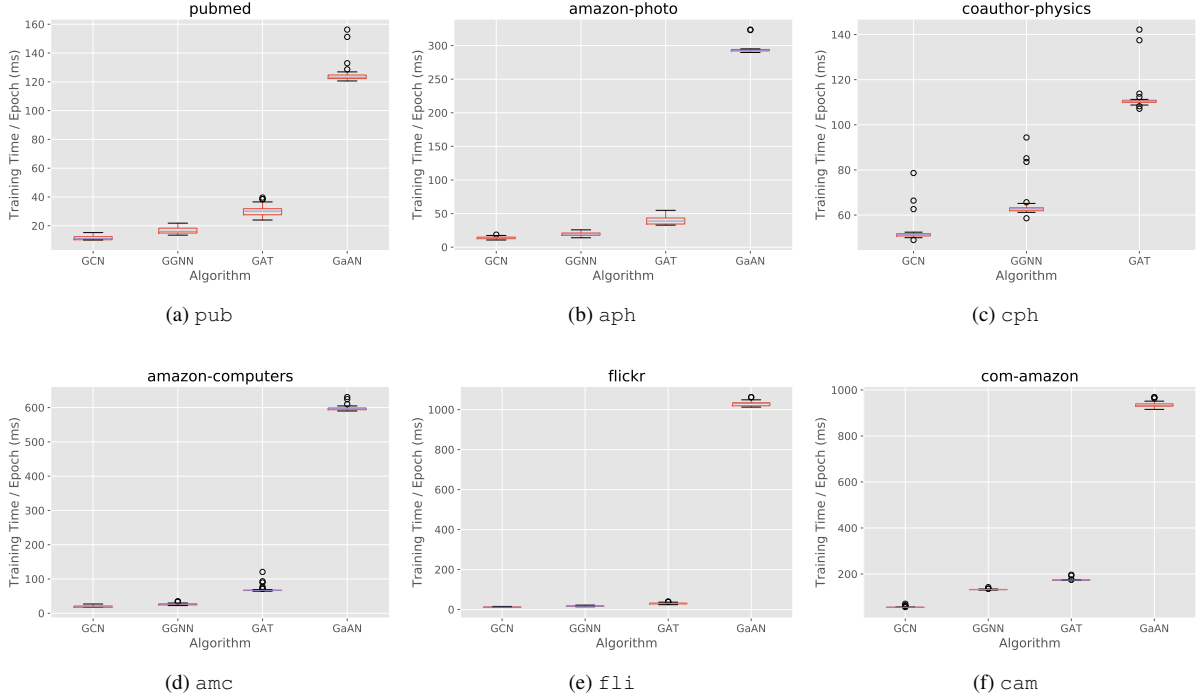


Figure 5: Distribution of the wall-clock training time of 50 epoches on different datasets. GaAN crashed due to out of memory exception on the `cph` dataset.

computing power of the GPU. When the hyper-parameter becomes high enough, the training time increases linearly, supporting the time complexity analysis.

We further measured the effects of the hyper-parameters on the peak GPU memory usage in Figure 7. The memory usage also increases linearly as the hyper-parameters increase for all GNNs, except for GaAN on $\dim(\mathbf{h}^1)$. As the hidden feature vectors \mathbf{h}^1 consume a small proportion of memory in GaAN, the growth in the memory usage is not noticeable until $\dim(\mathbf{h}^1)$ is large enough.

Summary. The complexity analysis in Table 1 and Table 2 is valid. Fixing other hyper-parameters, each hyper-parameter itself affects the training time and the memory usage of a GNN Layer **in a linear way**. Algorithm engineers can adjust hyper-parameters according to the time complexity to avoid explosive growth in the training time and memory usage.

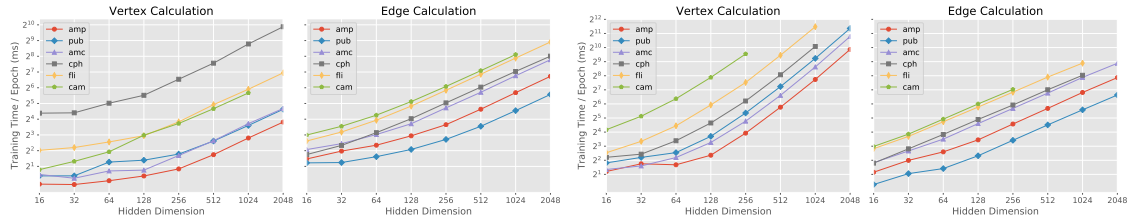
4.2. Training Time Breakdown

To find out which stage/step dominates the training time, we decompose the training time and analyze the performance bottleneck level by level.

4.2.1. Layer Level

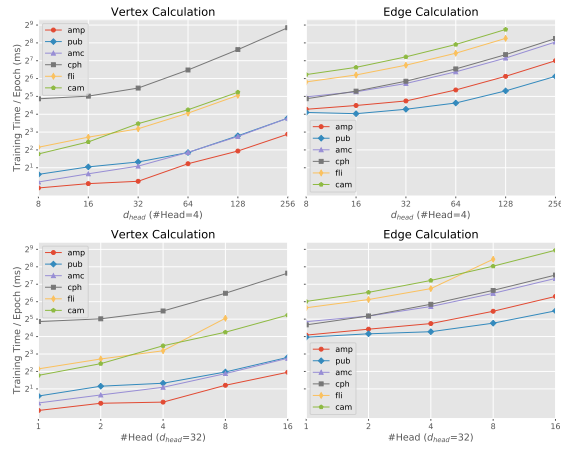
Figure 8 decomposes the training time of a GNN on the layer level. The training time of each layer is the summation of the time in the forward, backward, and evaluation phases. In GCN, GAT, and GaAN, the time spent on the layer 0 is much larger than the layer 1. In those GNNs, the dimensions of the input/output feature vectors in the layer 0 are much larger than the dimensions in the layer 1. $d_{in}^0 = \dim(\mathbf{x})$, $d_{out}^0 = d_{in}^1 = 64$ and $d_{out}^1 = \#Class$ and $\dim(\mathbf{x}) \gg \#Class$. For GaAN, since it requires the dimensions of the input/output feature vectors must be the same, the hyper-parameter are $d_{in}^0 = d_{out}^0 = d_{in}^1 = d_{out}^1 = 64$ and the training time of both layers is close.

Each GNN layer can be further divided into the vertex and the edge calculation stages. In Figure 8, GCN spends most of the training time on the edge calculation stage in most datasets. A special case is `cph` dataset. The dimension of the input feature vectors is very high in `cph`, making the vertex calculation stage of the GCN

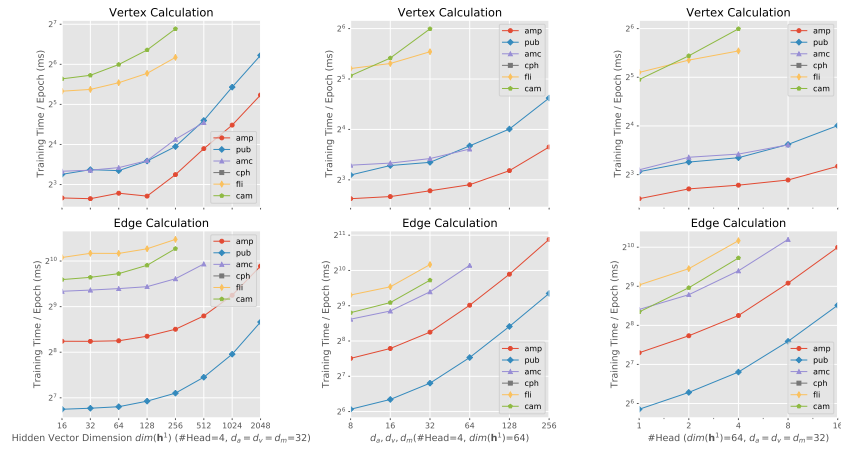


(a) GCN

(b) GGNN



(c) GAT



(d) GaAN

Figure 6: Effects of hyper-parameters on the edge/vertex calculation time.

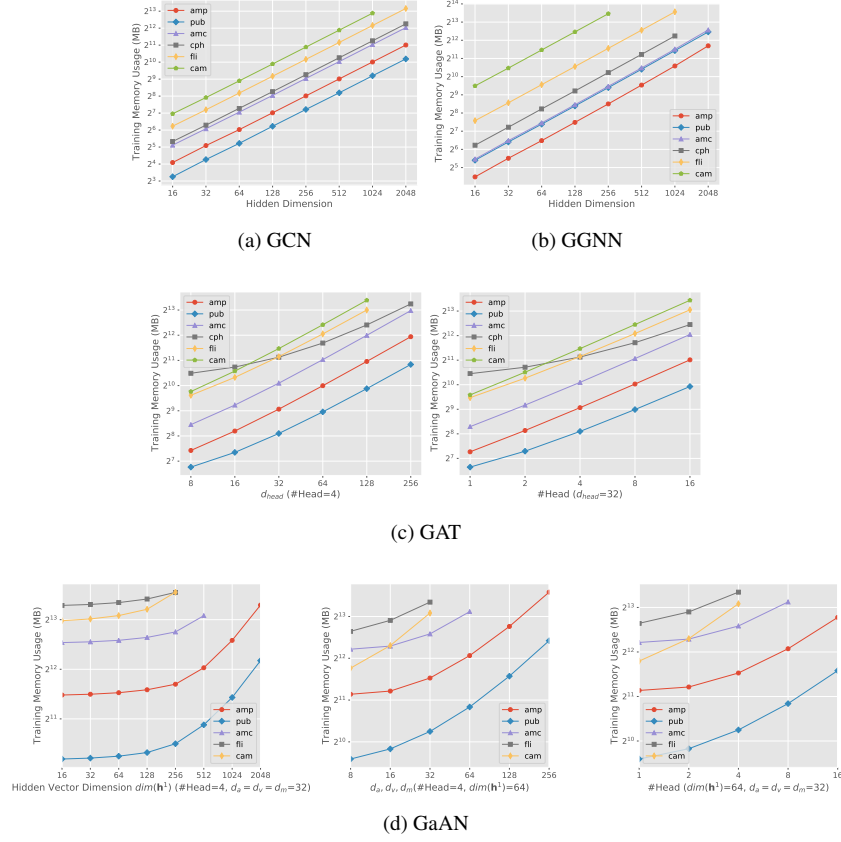


Figure 7: Effects of hyper-parameters on the peak GPU memory usage during the training, excluding the memory used by the dataset and the model parameters.

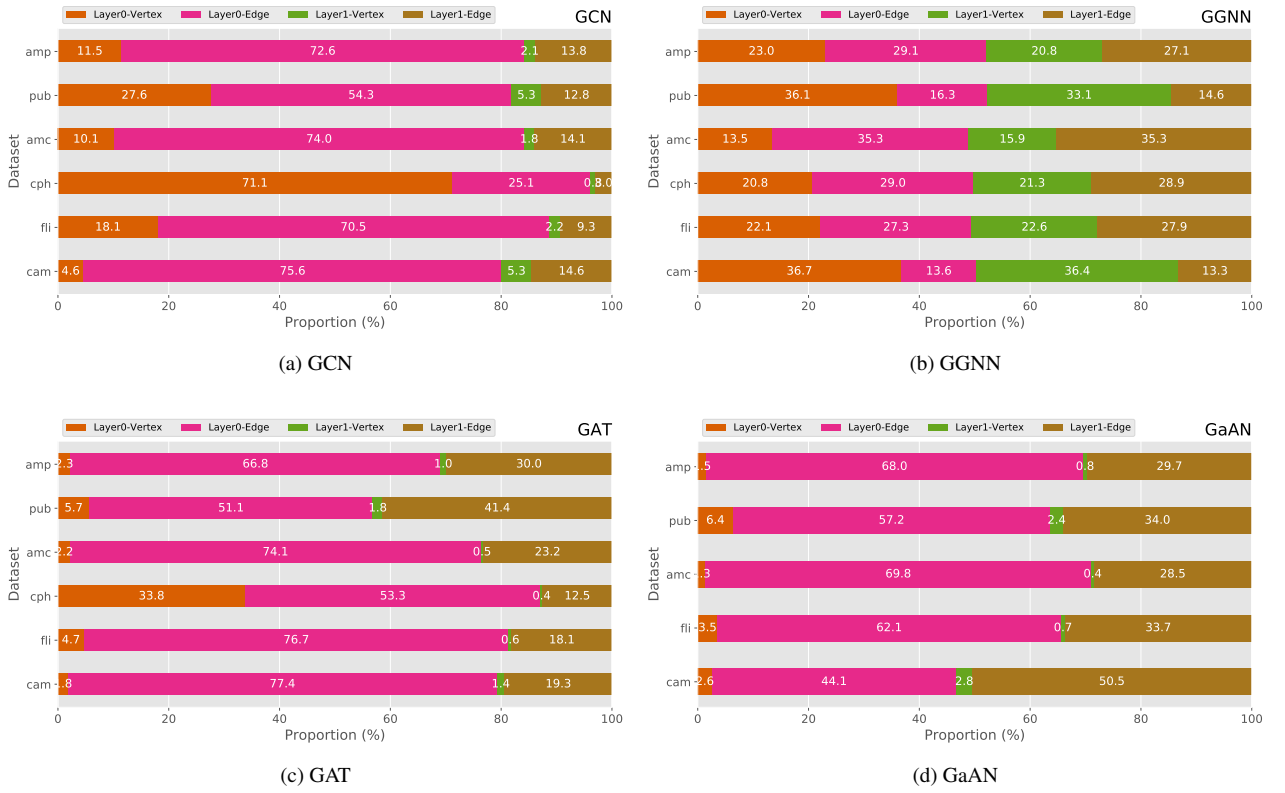


Figure 8: Training time breakdown on the layer level. The training time of each layer includes the time spent on the forward, backward and evaluation phases. Each layer is further decomposed into the vertex and the edge calculation stages.

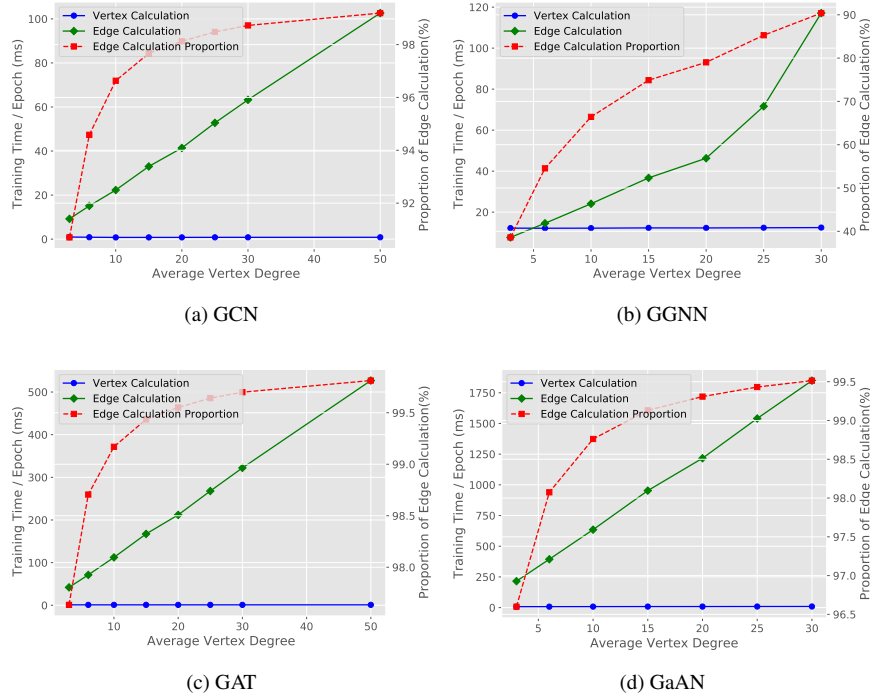


Figure 9: Effects of the average degree on the time proportion of the edge/vertex calculation. Graphs were generated with the R-MAT generator by fixing the number of vertices as 50,000.

Layer 0 spend considerable time. GGNN also spends the majority of its training time on the edge calculation stage. But the high time complexity of its vertex updating function γ makes the ratio of the vertex calculation in the total training time much higher than the other GNNs. For GAT and GaAN, due to their high edge calculation complexity, the edge calculation stage is the absolutely dominant stage. In summary, *the edge calculation stage is the most time-consuming stage in the GNN training.*

The experimental results also indicate that the average degree of the dataset affects the time-consuming proportion of the edge/vertex calculation time. For GaAN, the time spent on the vertex calculation stage exceeds the edge calculation stage on the pub and cam datasets, because the average degrees of the two datasets are low, making $|\mathcal{E}|$ and $|\mathcal{V}|$ much closer. To evaluate the effects of the average degree, we used the R-MAT generator to generate random graphs with 50k vertices and the average degrees ranging from 2 to 100. Figure 9 shows the training time of the four GNNs under different average degrees. As the average degree increases, the training time of the edge calculation stage grows *linearly*. For GCN, GAT, and GaAN, the edge calculation stage dominates the entire training time even under small average degrees. Only for GGNN that has high vertex and low edge calculation complexities, the training time of the vertex calculation stage exceeds the edge calculation stage under low average degrees (< 5). Therefore, *improving the efficiency of the edge calculation stage is the key to reduce the GNN training time.*

4.2.2. Step Level in Edge Calculation

In the implementation of PyG, the edge calculation stage can be decomposed into four steps: collection, messaging, aggregation, and updating, as shown in Figure 10. The edge index is a matrix with M rows and 2 columns that holds the edge set of the graph, where $M = |\mathcal{E}|$. The two columns of the matrix store the source vertex and the target vertex of each edge, respectively. The collection step copies the vertex feature vectors from the previous layer \mathbf{h}_i^l to the ends of each edge in the edge index, to form the parameters tensor $[\mathbf{h}_i^l, \mathbf{h}_j^l, \mathbf{e}_{i,j}^l]$ of the messaging function ϕ . This step only involves the data movement. The messaging step calls the messaging function ϕ to get message vectors of all edges $\mathbf{m}_{i,j}^l$. The aggregation step aggregates the message vectors with the same target vertex into an aggregated vector \mathbf{s}_i^l with the aggregation operation Σ . The updating step is optional. It performs an additional transformation on the aggregated vectors (for example, adding bias in GCN).

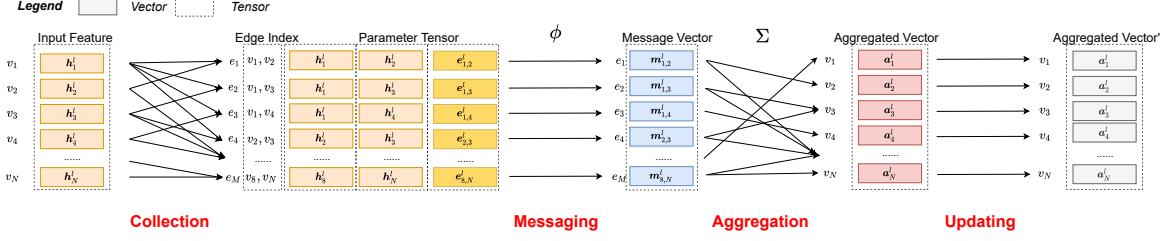


Figure 10: Step decomposition of the edge calculation in the GNN layer l .

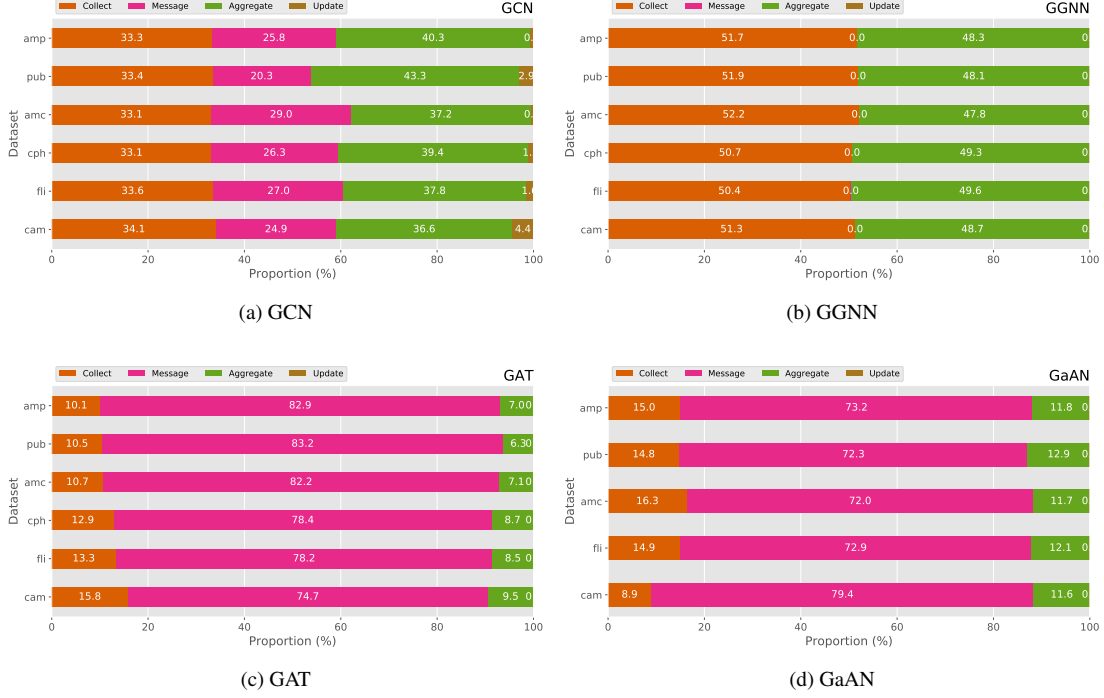


Figure 11: Training time breakdown of the edge calculation stage (including both GNN layers).

The aggregated vectors s_i^l (after the updating step) will be fed into the vertex updating function γ as one of the input parameters.

We decompose the execution time of the edge calculation stage in Figure 11. In each GNN, the proportions of the four steps are rather stable, rarely affected by datasets. For GAT and GaAN with the high edge calculation complexity, the messaging step consumes most of the training time. For GCN and GGNN with the low edge complexity, the proportions of the steps are close. Since the messaging function ϕ of GGNN uses the pre-computed \hat{h}_i^l as the message vector directly, the time spent on the messaging step of GGNN is negligible. Although the collecting step does not conduct any computation and only involves data movement, it occupies noticeable execution time in all the GNNs. The experiments show that *the performance bottleneck on the step level depends on the complexity of the messaging function ϕ* . For GNNs with the high complexity, the messaging function ϕ is the performance bottleneck. Optimizing the implementation of ϕ can significantly reduce the training time. For the other GNNs, optimization should focus on reducing the costs of the collection and the aggregation steps. Additionally, improving the efficiency of the collection step can benefit all GNNs.

4.2.3. Operator Level

The functions ϕ , Σ and γ in the vertex and edge calculation are made up of a series of basic operators implemented on the GPU, like the matrix multiplication `mm`, the elementwise multiplication `mul` and the index-based selection `index.select`. Figure 12 shows the five most time-consuming basic operators in each GNN,

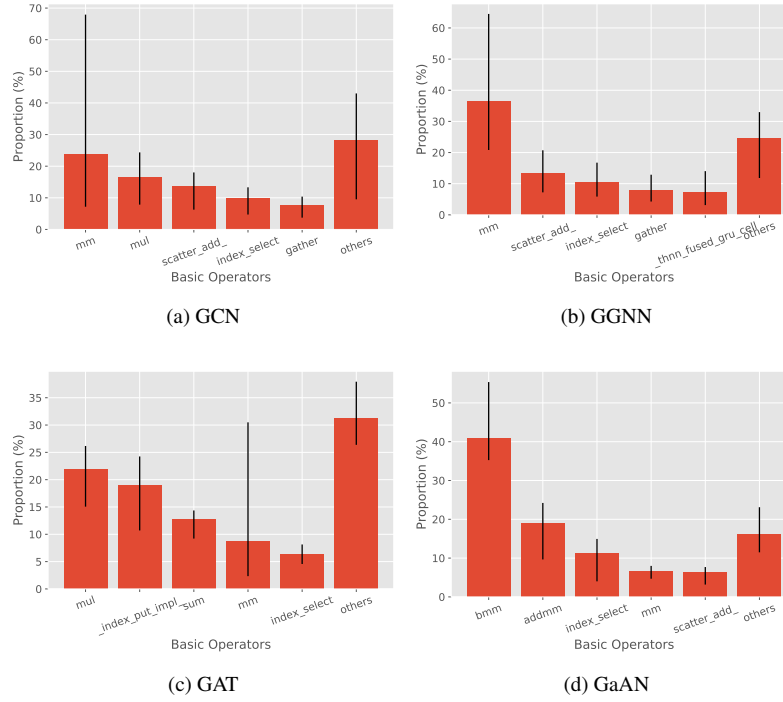


Figure 12: Top 5 time-consuming basic operators of typical GNNs. The time proportion of each basic operator is averaged over all graphs with the error bar indicating the maximal and the minimal.

averaged over all the real-world graphs in Table 3.

GCN. The most time-consuming basic operator is the matrix multiplication `mm` used in the vertex updating function γ . The elementwise multiplication `mul` used in the messaging function ϕ is also time-consuming. The other three operators are used in the edge calculation stage: `scatter_add` for the aggregation step in the forward phase, `gather` for the aggregation step in the backward phase, and `index_select` for the collection step. For GCN, the basic operators related to the edge calculation stage consume the majority of the training time.

GGNN. The top basic operator is `mm` used in the vertex updating function γ . Due to its high time complexity, the proportion of the `mm` is much higher than the other operators. The `thnn_fused_gru_cell` operator used in the backward phase of γ is also noticeable. The other three operators are used in the edge calculation stage.

GAT. All the top basic operators except for `mm` are related to the edge calculation stage. The `mm` operator is used in the vertex updating function γ .

GaAN. The top basic operator is `bmm` used in the messaging function ϕ . The `addmm` operator and the `mm` operator are used in both the vertex and the edge calculation stages, where the edge calculation stage is dominant.

In general, the most time-consuming operator in the four GNNs, is still the matrix multiplication `mm` and the elementwise multiplication `mul`, making GNN training suitable for GPUs. Although the aggregation step in the edge calculation stage is relatively simple (like sum and mean), the related operators `scatter_add` and `gather` still consume a certain amount of the time. The operators have to synchronize between hardware threads to avoid updating the same aggregated vector at the same time. They also conduct non-regular memory access with the access pattern determined by the edge set dynamically. For GPUs, they are less efficient than `mm`. The index-based selection `index_select` operator from the collection step consume around 10% of the time in all GNNs. Though GPUs have high on-chip memory bandwidth, improving the efficiency of `scatter_add`/`gather`/`index_select` can benefit the training of all kinds of GNNs.

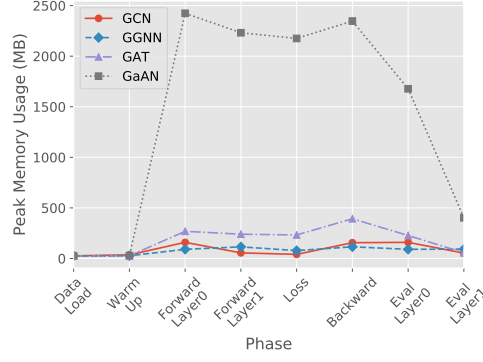


Figure 13: Memory usage of each phase. Dataset: amp.

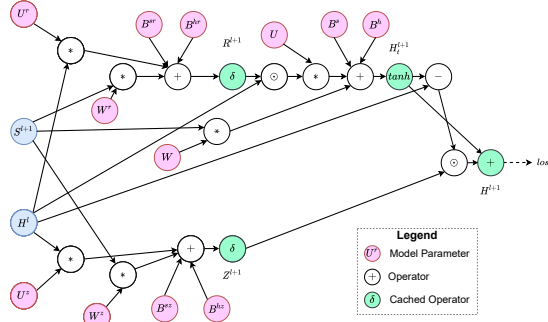


Figure 14: Computation graph of the vertex updating function γ of GGNN.

Summary. The GNN training is suitable for GPUs. **The edge calculation stage is the main performance bottleneck in most cases**, except for training GNNs with high vertex calculation complexity on low-average-degree graphs. The performance bottleneck in the edge calculation stage depends on the time complexity of the messaging function ϕ .

- If the time complexity of ϕ is **high**, the **efficiency of ϕ** limits the performance. Reducing its computation cost (via optimizing its implementation or modifying the algorithm) can significantly reduce the training time.
- If the time complexity of ϕ is **low**, the **collection step** and the **aggregation step** limit the performance. The collection step involves lots of data movement. The aggregation step suffers from data synchronization and non-regular data access. Optimizing their implementations can significantly reduce the training time.

4.3. Memory Usage Analysis

During the GNN training, all data (including datasets and intermediate results) are stored in the on-chip memory of the GPU. Compared with the main memory on the host side, the capacity of the GPU memory is very limited. *The GPU memory capacity limits the scales of the graphs that a GPU can train GNNs on.* For example, GaAN is unable to train on `cph` dataset due to the out of memory exception.

Figure 13 shows the peak memory usage of each phase in the GNN training on the `amp` dataset. The trend is similar on the other datasets. *The GNN training achieves its peak memory usage in the forward and the backward phases.* The forward phase generates lots of intermediate results. Some key intermediate results are cached, increasing memory usage. The cached results are used in the gradient calculation in the backward phase. Figure 14 shows the computation graph of the vertex updating function γ of GGNN. The computation graph has a large number of operators. Each operator generates an intermediate tensor. Some key intermediate tensors are cached. The cached tensors are the main source of memory usage in the loss phase. By the end of

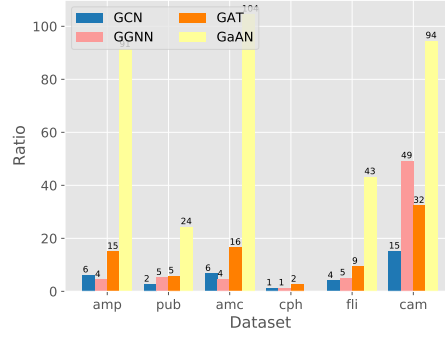


Figure 15: Memory expansion ratios of typical GNNs.

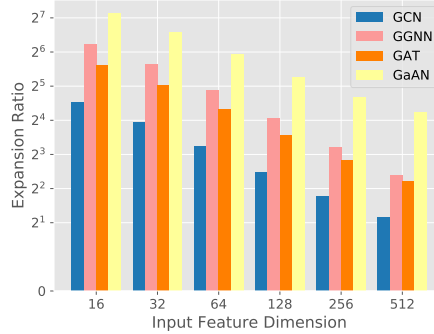


Figure 16: Memory expansion ratio under different dimensions of the input feature vectors. Dataset: cam.

the backward phase, the cached tensors are released. Since the evaluation phase does not need to calculate the gradients, it does not cache intermediate tensors. Its memory usage declines sharply.

The peak memory usage during the GNN training far exceeds the size of the dataset itself. We define the *memory expansion ratio* (MER) as the ratio of the peak memory usage during the training to the memory usage after loading the dataset. Figure 15 compares MER of different GNNs. GCN has the lowest MER (up to 15) while GaAN has the highest MER (up to 104). *The high MER limits the data scalability of GNNs*, making GPUs unable to handle big graphs.

Figure 15 also indicates that the same GNN has different MERs for different datasets. Two characteristics of a dataset affect the MER: the dimension of the input feature vector and the average degree.

Given the same graph, the scales of the intermediate results are mainly affected by the hyper-parameters of the GNN. If the dimension of the input feature vectors is high (like the `cph` dataset), the size of the dataset is large. The size may become comparable to the scales of the intermediate results, making the MER low. To find out how the dimension affects the MER, we generated random input feature vectors with different dimensions for the `cam` dataset and measured the MER in Figure 16. For a GNN under the same hyper-parameters, *the MER decreases as the dimension of the input feature vectors increases*.

The average degree of the graph also affects MER by influencing the relative scales of the intermediate results from the edge/vertex calculation stages. Fixing the number of vertices $|\mathcal{V}|$, we used the R-MAT generator to generate random graphs with different average degrees. Figure 17 shows how the memory usage changes according to the average degrees. As the average degree \bar{d} increases, the number of edges $|\mathcal{E}|$ increases and the peak memory usage increases *linearly* with \bar{d} . The edge calculation stage gradually dominates the memory usage and *the MER converges to a stable value*. The stable value is determined by the complexity of the edge calculation stage. Except for GGNN, the MER of the other GNNs increases as \bar{d} increases. As GGNN has high vertex calculation complexity, the MER related to the vertex calculation stage is much larger than the edge calculation stage. When the edge calculation stage dominates the memory usage, its MER becomes smaller.

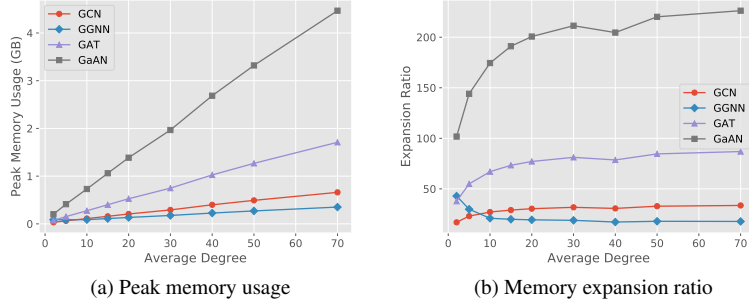


Figure 17: Memory usage under different average degrees of the graph. The graph was generated with the R-MAT generator fixing the number of vertices at 10K and the dimension of the input feature vectors at 32.

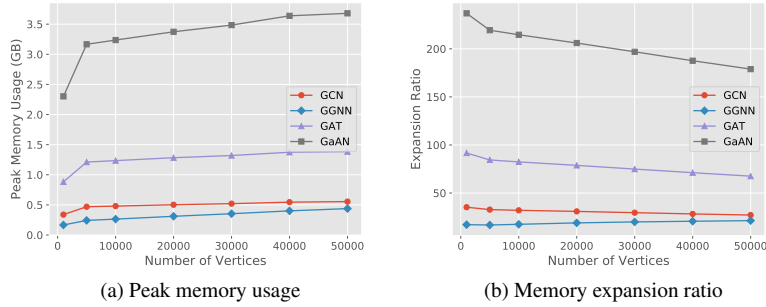


Figure 18: Memory usage under different numbers of vertices of the graph. The graph was generated with the R-MAT generator fixing the number of edges at 500K and the dimension of the input feature vectors at 32.

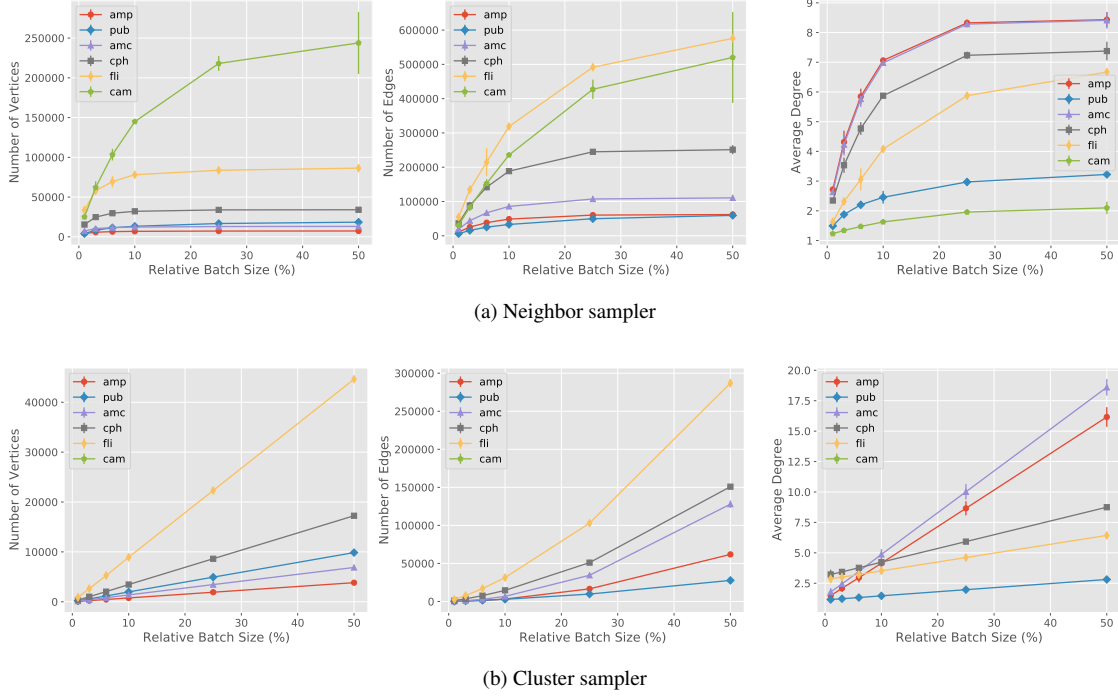


Figure 19: Sizes of the sampled subgraphs under different batch sizes. Each batch size was sampled 50 times and the average value was reported. The error bar indicates the standard deviation. The batch size is relative to the full graph.

We also fixed the number of edges $|\mathcal{E}|$ in the graphs and generated random graphs with different $|\mathcal{V}|$. Figure 18 shows how the memory usage changes according to $|\mathcal{V}|$. All GNNs are insensitive to the changes in $|\mathcal{V}|$ compared to $|\mathcal{E}|$. Except for GGNN, the MERs of the other GNNs decline as $|\mathcal{V}|$ increases because the scales of the datasets increase more quickly than the scales of the intermediate results. As GGNN has high vertex calculation complexity, the scales of the intermediate results are much more sensitive to $|\mathcal{V}|$. It indicates that *the intermediate results of the edge calculation stage dominated the memory usage during the GNN training.*

Summary. The *high* memory expansion ratio severely restricts the data scalability of the GNN training. The memory usage mainly comes from the intermediate results of the *edge calculation stage*. Fixing the number of vertices, the memory usage increases *linearly* along with the number of edges. Optimizing the memory usage of the edge calculation stage can significantly reduce the memory expansion ratio. Fixing the GNN structures and the hyper-parameters, increasing the dimension of the input feature vectors can also reduce the memory expansion ratio.

4.4. Effects of Sampling Techniques on Performance

With the sampling techniques, GNNs can be trained in a mini-batch manner. Each mini-batch updates the model parameters based on a small subgraph sampled from the original input graph. Thus, the training time per batch and the peak memory usage during the training should both decline significantly.

In the implementation in PyG, the GNN model and the dataset resides on the GPU side. To process each epoch, PyG samples the original dataset in the main memory and generates several batches. Each batch is a small subgraph of the dataset. To train on each batch, PyG sends the sampled subgraph to the GPU, calculates the gradients on the subgraph, and updates the model parameters directly on the GPU. With the sampling techniques, the model parameters are updated by a stochastic gradient descent optimizer. PyG conducts the evaluation phase every several epochs (either on the CPU side or the GPU side) to determine whether to stop the training. In this section, the experiments focus on the training phase of each batch.

Figure 19 shows how the size of the sampled subgraph changes with the batch size. For the neighbor sampler, the relative batch size is the proportion of the sampled vertices of the last GNN layer in $|\mathcal{V}|$. For the cluster sampler, the relative batch size is the proportion of the sampled partitions to all partitions of the graph.

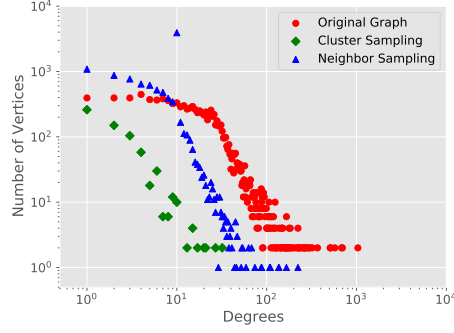


Figure 20: Vertex degree distribution of the sampled subgraph (relative batch size: 6%) and the original graph. Dataset:amp.

The neighbor sampler is very sensitive to the batch size. As the batch size increases, the size of the sampled subgraph first increases quickly and then stabilizes. The cluster sampler is much less sensitive compared to the neighbor sampler. The number of vertices and the average degree of the sampled subgraphs increases linearly with the batch size.

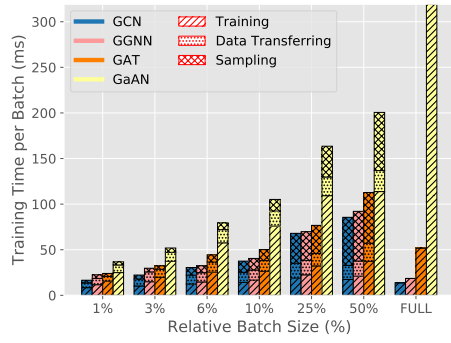
It is worth noting that the average degree of the sampled subgraph is *much lower* than the average degree of the whole graph, especially when the relative batch size is low. Taking the neighbor sampler with the relative batch size of 6% as an example, the average degree of the amp dataset is 31.1, but the average degree of the sampled subgraph is only 5.8. For the cluster sampler, the average degree is 3.0. Figure 20 compares the degree distribution of the sampled subgraphs with the original graph. The slopes of the curves are similar. It indicates that the sampled subgraphs still follow the power-law degree distribution. However, there are much less vertices in the sampled subgraphs, significantly lowering the average degrees. According to the experimental results in Section 4.2, if the average degree becomes lower, the proportion of the training time spent on the vertex calculation stage will become higher, especially for GGNN.

To find out the performance bottleneck with the sampling techniques, we decompose the training time per batch into three phases: *sampling* on the CPU, *transferring* sampled subgraphs from the CPU to the GPU and *training* with the subgraphs on the GPU. Figure 21 shows the time breakdown of the four GNNs under different relative batch sizes. For the neighbor sampler, the sampling technique reduces the training time per batch only when the batch size is very small. When the batch becomes bigger, the sampling and the data transferring phases introduce noticeable overheads, making the training time exceed the full-batch training. For the clustering sampler, the sampled subgraph is smaller than the neighbor sampler under the same relative batch size. The reduction in the training time is more obvious than the neighbor sampler. However, the overheads increase quickly as the relative batch size increase. The training time under the 25% relative batch size already exceeded the time of full-batch training. The experimental results indicate that the current implementation of the sampling techniques in PyG is inefficient. When the batch size is slightly big, more than 50% of the time has been spent on sampling and data transferring. *The sampling techniques are only efficient under small batch sizes.*

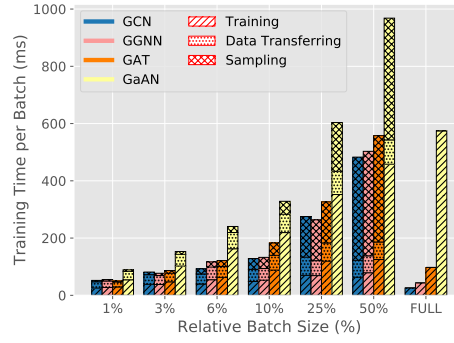
The main advantage of the sampling technique is *reducing the peak memory usage* during the training. Figure 22 shows the memory usage under different batch sizes. The peak memory usage declines significantly even under big batch sizes. The sampling techniques make training GNNs on big graphs possible for GPUs.

The disadvantage of the sampling technique is wasting GPU resources. As the sampling techniques are only effective under small batch sizes, the sampled subgraphs will be very small in those cases. They cannot make full use of the computing power of a GPU. To simulate the situation, we generated random graphs with few vertices and measured the training time per epoch in Figure 23. As the number of vertices increases, the training time is almost unchanged except for GaAN. The training time of GaAN increases only with $|\mathcal{V}| \geq 4000$.

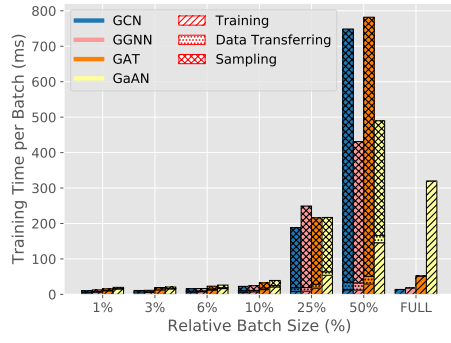
Summary. The sampled subgraphs have lower average degrees than the whole graph. With small batch sizes, the sampling techniques can significantly reduce the training time per batch and the peak memory usage during the training. However, small batch sizes cannot make full use of the computing power of a GPU. With big batch sizes, the current implementation of the sampling techniques in PyG is inefficient. The time spent on the sampling phase and the data transferring phase even exceeds the training phase.



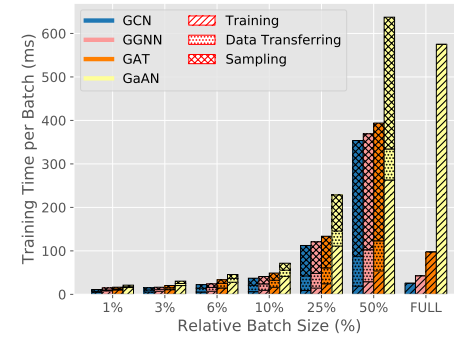
(a) Neighbor sampler on `amc`



(b) Neighbor sampler on `fli`



(c) Cluster sampler on `amc`



(d) Cluster sampler on `fli`

Figure 21: Training time per batch breakdown. FULL means that the full graph participates in the training.

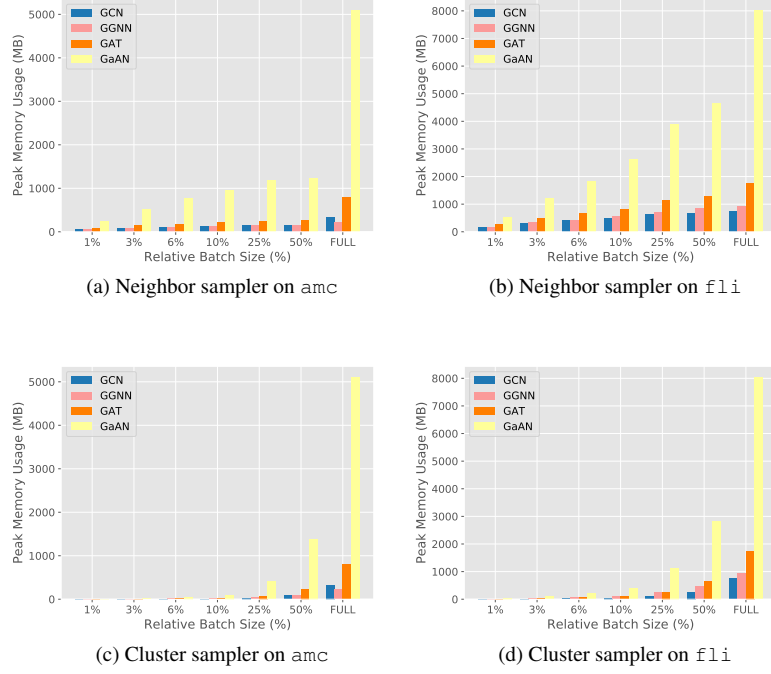


Figure 22: Peak memory usage under different batch sizes. FULL means the full graph participates in the training.

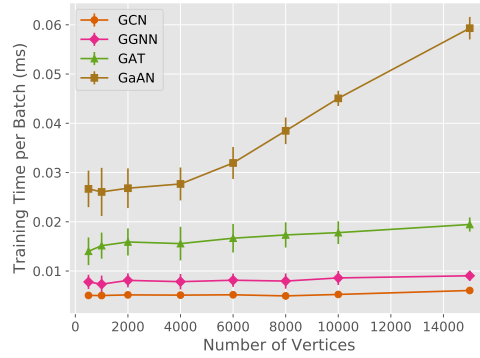


Figure 23: Training time per epoch on small random graphs. For each number of vertices, we generated 50 random R-MAT graphs with the average degree of 4.0 and reported the average training time per epoch (without the evaluation phase). The error bar indicates the standard deviation.

5. Insights

Through the extensive experiments, we propose the following key findings and suggestions for how to optimize the performance of the GNN training.

1. *The time complexity in Table 1 and Table 2 points out the performance bottleneck theoretically.* The experimental results validate the time complexity analysis. The time complexity points out where the bottleneck comes from. Optimization should focus on complex operations in the messaging function ϕ and the vertex updating function γ .
2. *The computational cost of a GNN layer is mainly affected by the dimensions of the input and the output hidden feature vectors.* Theoretically and empirically, the training time and the memory usage of a GNN layer both increase *linearly* with the dimensions of the input/output hidden feature vectors separately. GNNs are friendly to high-dimensional scenarios. Algorithm engineers can use high-dimensional feature vectors to improve the expressive power of a GNN without worrying exponential growth in the training time and memory usage.
3. *Performance optimizations should focus on improving the efficiency of the edge calculation stage.* The edge calculation stage is the most time-consuming stage in most GNNs.
 - If the complexity of the messaging function ϕ is high, the implementation of ϕ is critical to performance. Improving its efficiency can significantly reduce the training time. For example, the attention mechanism in GNNs (like GAT and GaAN) requires an extra sub-layer to calculate the attention weight of each edge. Implementing it with the specially optimized basic operators on the GPU is a potential optimization.
 - If the complexity of ϕ is low, the efficiency of the collection step and the aggregation step becomes critical. The existing GNN libraries [15, 14, 16] already introduce the *fused* operator to improve their efficiency. When the messaging function ϕ is an assignment or a scalar multiplication of the hidden feature vector of the source vertex, the libraries replace the collection, messaging, and aggregation steps with a single fused operator. The fused operator calculates the aggregated vectors directly from the input hidden feature vectors, minimizing the memory footprints and overlapping the memory accessing with computation. In this way, it significantly reduces the training time of GNNs with low edge calculation complexity (like GCN) [19, 20]. However, the applicable condition of the fused operator is very restricted. It does not work for ϕ with more complex operations like matrix multiplication. A potential optimization is to develop composite CUDA kernels that can read the input hidden feature vectors and aggregate message vectors on the fly, without materializing the parameter vectors and the message vectors.
4. *The high memory usage caused by the intermediate results of the edge calculation stage limits the data scalability of the GNN training.* The memory expansion ratios of the typical GNNs are very high, making GPUs unable to handle big graphs. One solution is to distribute the dataset among several GPUs and frequently swap parts of the dataset between GPUs and the main memory [16]. Another possible solution [37] comes from the deep neural network training. It only checkpoints key intermediate results during the forward propagation and re-calculates the missing results on demand during the backpropagation. Implementing the checkpoint mechanism in the GNN training is another potential optimization.
5. *Sampling techniques can significantly reduce the training time and memory usage, but its implementation is still inefficient.* The sampling techniques are effective under small batch sizes. Its current implementation brings considerable overheads when the batch size becomes large. Improving the efficiency of the sampling is a potential optimization. The sampled subgraphs are usually small. They cannot make full use of the computing power of a GPU. How to improve the GPU utilization under small batch sizes is another problem to solve. One possible solution is to train multiple batches asynchronously on the same GPU and use the asynchronous stochastic gradient descent to speed up the converge.

6. Related Work

Survey of GNNs. Zhou et al.[8], Zhang et al.[9] and Wu et al.[10] survey the existing graph neural networks and classify them from an algorithmic view. They summarize the similarities and differences between the architectures of different GNNs. The typical applications of GNNs are also briefly introduced. Those surveys focus on comparing the existing GNNs theoretically, not empirically.

Evaluation of GNNs. Shchur et al. [34] evaluate the accuracy of popular GNNs on the node classification task. Dwivedi et al. [38] further compare the accuracy of popular GNNs fairly in a controlled environment. Hu et al. [39] propose the open graph benchmark that provides a standard datasets and a standard evaluation workflow. The benchmark makes comparisons between GNNs easily and fairly. Those model evaluation efforts focus on evaluating the accuracy of different GNNs. They provide insightful suggestions to improve accuracy.

From the efficiency aspect, Yan et al. [19] compare the performance characteristics of graph convolutional networks, typical graph processing kernels (like PageRank), and the MLP-based neural networks on GPUs. They provide optimization guidelines for both the software and the hardware developers. Zhang et al. [20] analyze the architectural characteristics of the GNN inference on GPUs under SAGA-NN [16] model. They find that the GNN inference has no fixed performance bottleneck and all components deserve to optimize. These two efforts focus on the *inference* phase of GNNs and they investigate the potential optimizations mainly from an architectural view. In this work, our target is to find out the performance bottleneck in the training phase from a system view. We consider the performance bottleneck in both time and memory usage. We also evaluate the effects of the sampling techniques. Our work and the related evaluation [19, 20] form a complementary study on the efficiency issue of GNNs.

Libraries/Systems of GNNs. PyG [14] and DGL [15] both adopt the message-passing model as the underlying programming model for GNNs and support training big datasets with the sampling techniques. PyG [14] is built upon PyTorch and it uses optimized CUDA kernels for GNNs to achieve high performance. DGL [15] provides a group of high-level user APIs and supports training GNNs with a variety of backends (TensorFlow, MXNet, and PyTorch) transparently. It also supports LSTM as the aggregation functions. NeuGraph [16] proposes a new programming model *SAGA-NN* for GNNs. It focuses on training big datasets efficiently without sampling. It partitions the dataset sophisticatedly, schedules the training tasks among multiple GPUs, and swaps the data among GPUs and the host asynchronously. AliGraph [17] targets at training GNNs on big attributed heterogeneous graphs that are common in e-commerce platforms. The graphs are partitioned among multiple nodes in a cluster and AliGraph trains GNNs on the graphs in a distributed way with system optimizations. PGL [18] is another graph learning framework from Baidu based on the PaddlePaddle platform.

7. Conclusion

In this work, we systematically explore the performance bottleneck in graph neural network training. We model the existing GNNs with the message-passing framework. We classify the GNNs according to their edge and vertex calculation complexities to select four typical GNNs for evaluation. The experimental results validate our complexity analysis. Fixing other hyper-parameters, the training time and the memory usage increase linearly with each hyper-parameter of the four GNNs. To find out the performance bottleneck in the training time, we decompose the training time per epoch on different levels. The training time breakdown analysis indicates that the edge calculation stage and its related basic operators are the performance bottleneck for most GNNs. Moreover, the intermediate results produced by the edge calculation stage cause high memory usage, limiting the data scalability. Adopting sampling techniques can reduce the training time and the memory usage significantly. However, the current implementation of the sampling techniques in PyG brings considerable sampling overheads. The small sampled subgraphs cannot make full use of the computing power of a GPU card either. Our analysis indicates that the edge calculation stage should be the main target of optimizations. Reducing its memory usage and improving its efficiency can significantly improve the performance of the GNN training. Based on the analysis, we propose several potential optimizations for the GNN frameworks. We believe that our analysis can help developers to have a better understanding of the characteristics of GNN training.

Acknowledgements

This work is funded in part by National Key R&D Program of China [grant number 2019YFC1711000]; China NSF Grants [grant number U1811461]; Jiangsu Province Industry Support Program [grant number BE2017155]; Natural Science Foundation of Jiangsu Province [grant number BK20170651]; Collaborative Innovation Center of Novel Software Technology and Industrialization; and the program B for Outstanding PhD candidate of Nanjing University.

References

- [1] T. N. Kipf, M. Welling, Semi-supervised classification with graph convolutional networks, 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings, OpenReview.net, 2017.
- [2] M. Defferrard, X. Bresson, P. Vandergheynst, Convolutional neural networks on graphs with fast localized spectral filtering, 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.
- [3] R. Li, S. Wang, F. Zhu, J. Huang, Adaptive graph convolutional neural networks, Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018, 2018, pp. 3546–3553.
- [4] Y. Li, D. Tarlow, M. Brockschmidt, R. S. Zemel, Gated graph sequence neural networks, 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings, 2016.
- [5] W. L. Hamilton, Z. Ying, J. Leskovec, Inductive representation learning on large graphs, Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA, 2017, pp. 1024–1034.
- [6] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, Y. Bengio, Graph attention networks, 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, 2018.
- [7] J. Zhang, X. Shi, J. Xie, H. Ma, I. King, D. Yeung, Gaan: Gated attention networks for learning on large and spatiotemporal graphs, Proceedings of the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence, UAI 2018, Monterey, California, USA, August 6-10, 2018, 2018, pp. 339–349.
- [8] J. Zhou, G. Cui, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, M. Sun, Graph neural networks: A review of methods and applications, <https://arxiv.org/abs/1812.08434> (2018). arXiv:1812.08434.
- [9] Z. Zhang, P. Cui, W. Zhu, Deep learning on graphs: A survey, IEEE Transactions on Knowledge and Data Engineering (Early Access) (2020) 1–1. doi:10.1109/TKDE.2020.2981333.
- [10] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, P. S. Yu, A comprehensive survey on graph neural networks, IEEE Transactions on Neural Networks and Learning Systems (Early Access) (2020) 1–21. doi:10.1109/TNNLS.2020.2978386.
- [11] S. Qi, W. Wang, B. Jia, J. Shen, S. Zhu, Learning human-object interactions by graph parsing neural networks, in: V. Ferrari, M. Hebert, C. Sminchisescu, Y. Weiss (Eds.), Computer Vision – ECCV 2018, Vol. 11213 of Lecture Notes in Computer Science, Springer, 2018, pp. 407–423. doi:10.1007/978-3-030-01240-3_25.
- [12] W. Wang, H. Zhu, J. Dai, Y. Pang, J. Shen, L. Shao, Hierarchical human parsing with typed part-relation reasoning, in: Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, 2020, pp. 8926–8936. doi:10.1109/CVPR42600.2020.00895.
- [13] W. Wang, X. Lu, J. Shen, D. J. Crandall, L. Shao, Zero-shot video object segmentation via attentive graph neural networks, in: Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision (ICCV), IEEE, 2019, pp. 9235–9244. doi:10.1109/ICCV.2019.00933.
- [14] M. Fey, J. E. Lenssen, Fast graph representation learning with pytorch geometric, <https://pytorch-geometric.readthedocs.io/> (2019). arXiv:1903.02428.

- [15] M. Wang, L. Yu, D. Zheng, Q. Gan, Y. Gai, Z. Ye, M. Li, J. Zhou, Q. Huang, C. Ma, Z. Huang, Q. Guo, H. Zhang, H. Lin, J. Zhao, J. Li, A. J. Smola, Z. Zhang, Deep graph library: Towards efficient and scalable deep learning on graphs, <https://www.dgl.ai/> (2019). arXiv:1909.01315.
- [16] L. Ma, Z. Yang, Y. Miao, J. Xue, M. Wu, L. Zhou, Y. Dai, Neugraph: Parallel deep neural network computation on large graphs, 2019 USENIX Annual Technical Conference, USENIX ATC 2019, Renton, WA, USA, July 10-12, 2019, 2019, pp. 443–458.
- [17] R. Zhu, K. Zhao, H. Yang, W. Lin, C. Zhou, B. Ai, Y. Li, J. Zhou, Aligraph: A comprehensive graph neural network platform, *Proc. VLDB Endow.* 12 (12) (2019) 2094–2105. doi:10.14778/3352063.3352127.
- [18] Baidu, Paddle graph learning, <https://github.com/PaddlePaddle/PGL>.
- [19] M. Yan, Z. Chen, L. Deng, X. Ye, Z. Zhang, D. Fan, Y. Xie, Characterizing and understanding gcns on gpu, *IEEE Computer Architecture Letters* 19 (1) (2020) 22–25.
- [20] Z. Zhang, J. Leng, L. Ma, Y. Miao, C. Li, M. Guo, Architectural implications of graph neural networks, *IEEE Computer Architecture Letters* 19 (1) (2020) 59–62.
- [21] B. Perozzi, R. Al-Rfou, S. Skiena, Deepwalk: online learning of social representations, The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New York, NY, USA - August 24 - 27, 2014, 2014, pp. 701–710. doi:10.1145/2623330.2623732.
- [22] A. Grover, J. Leskovec, node2vec: Scalable feature learning for networks, Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016, 2016, pp. 855–864. doi:10.1145/2939672.2939754.
- [23] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, G. E. Dahl, Neural message passing for quantum chemistry, Vol. 70 of Proceedings of Machine Learning Research, Proceedings of the 34th International Conference on Machine Learning, PMLR, International Convention Centre, Sydney, Australia, 2017, pp. 1263–1272.
- [24] D. Duvenaud, D. Maclaurin, J. Aguilera-Iparraguirre, R. Gómez-Bombarelli, T. Hirzel, A. Aspuru-Guzik, R. P. Adams, Convolutional networks on graphs for learning molecular fingerprints, Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada, 2015, pp. 2224–2232.
- [25] H. Dai, Z. Kozareva, B. Dai, A. J. Smola, L. Song, Learning steady-states of iterative algorithms over graphs, Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018, 2018, pp. 1114–1122.
- [26] W.-L. Chiang, X. Liu, S. Si, Y. Li, S. Bengio, C.-J. Hsieh, Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks, Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, ACM, 2019, pp. 257–266. doi:10.1145/3292500.3330925.
- [27] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, J. Leskovec, Graph convolutional neural networks for web-scale recommender systems, Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018, 2018, pp. 974–983. doi:10.1145/3219819.3219890.
- [28] J. Chen, T. Ma, C. Xiao, Fastgcn: Fast learning with graph convolutional networks via importance sampling, 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, 2018.
- [29] J. Chen, J. Zhu, L. Song, Stochastic training of graph convolutional networks with variance reduction, Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018, 2018, pp. 941–949.

- [30] W. Huang, T. Zhang, Y. Rong, J. Huang, Adaptive sampling towards fast graph representation learning, Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada, 2018, pp. 4563–4572.
- [31] H. Zeng, H. Zhou, A. Srivastava, R. Kannan, V. K. Prasanna, Accurate, efficient and scalable graph embedding, 2019 IEEE International Parallel and Distributed Processing Symposium, IPDPS 2019, Rio de Janeiro, Brazil, May 20-24, 2019, 2019, pp. 462–471. doi:10.1109/IPDPS.2019.00056.
- [32] H. Zeng, H. Zhou, A. Srivastava, R. Kannan, V. K. Prasanna, Graphsaint: Graph sampling based inductive learning method, 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020, OpenReview.net, 2020.
- [33] Z. Yang, W. W. Cohen, R. Salakhutdinov, Revisiting semi-supervised learning with graph embeddings, Vol. 48 of JMLR Workshop and Conference Proceedings, Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016, JMLR.org, 2016, pp. 40–48.
- [34] O. Shchur, M. Mumme, A. Bojchevski, S. Günnemann, Pitfalls of graph neural network evaluation, CoRR abs/1811.05868. arXiv:1811.05868.
- [35] D. Chakrabarti, Y. Zhan, C. Faloutsos, R-mat: A recursive model for graph mining, Proceedings of the 2004 SIAM International Conference on Data Mining, pp. 442–446. doi:10.1137/1.9781611972740.43.
- [36] J. Yang, J. Leskovec, Defining and evaluating network communities based on ground-truth, 12th IEEE International Conference on Data Mining, ICDM 2012, Brussels, Belgium, December 10-13, 2012, IEEE Computer Society, 2012, pp. 745–754. doi:10.1109/ICDM.2012.138.
- [37] T. Chen, B. Xu, C. Zhang, C. Guestrin, Training deep nets with sublinear memory cost, https://arxiv.org/abs/1604.06174 (2016). arXiv:1604.06174.
- [38] V. P. Dwivedi, C. K. Joshi, T. Laurent, Y. Bengio, X. Bresson, Benchmarking graph neural networks (2020). arXiv:2003.00982.
- [39] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, J. Leskovec, Open graph benchmark: Datasets for machine learning on graphs, https://arxiv.org/abs/2005.00687 (2020). arXiv:2005.00687.