

Architectural Implication of Graph Neural Networks

Zhihui Zhang¹, Jingwen Leng^{1*}, Lingxiao Ma², Youshan Miao³, Chao Li,¹ Minyi Guo^{1*}

Abstract—Graph neural networks (GNN) represent an emerging line of deep learning models that operate on graph structures. It is becoming more and more popular due to its high accuracy achieved in many graph-related tasks. However, GNN is not as well understood in the system and architecture community as its counterparts such as multi-layer perceptrons and convolutional neural networks. This work tries to introduce the GNN to our community. In contrast to prior work that only presents characterizations of GCNs, our work covers a large portion of the varieties for GNN workloads based on a general GNN description framework. By constructing the models on top of two widely-used libraries, we characterize the GNN computation at inference stage concerning general-purpose and application-specific architectures and hope our work can foster more system and architecture research for GNNs.

Index Terms—Graph neural networks, computation analysis, deep learning, characterization.

1 INTRODUCTION

GRAPH neural networks (GNN) start to gain momentum as researchers are considering important tasks involving the *graph* structure such as social media. Deep learning (DL) has achieved great success in domains with *grid* data structure, e.g., images and sequences, which, however, only represents a small portion of the real-world data. In contrast, *graph* structure reflects the vast majority of real-world data such as molecular structure and knowledge graph.

Graph representation learning is one of the most important graph-related problems [1]. It converts the irregular graph structure to embedding vectors, which are the compressed representations of vertices (i.e. vertex embedding) and the entire graph (i.e. graph embedding). A downstream task such as molecular property prediction can then take in the regular embeddings rather than the raw graph for efficient processing. As such, the quality of the embeddings directly determines the accuracy of downstream tasks.

Traditional graph representation methods like DeepWalk [2] and node2vec [3] mostly rely on hand-crafted or intuition-based algorithms. In contrast, GNNs extend the graph analytics with DL’s end-to-end learning capability, which has led to better accuracies in a variety of domains including molecular science, recommendation, and transportation. To realize the full potential of GNNs, we should adapt the existing software and hardware platforms to GNNs’ unique characteristics.

The combination of DL and graph analytics makes GNN a new computation paradigm, which is quite different from their counterparts such as multi-layer perceptrons (MLP) and convolutional neural networks (CNN). Figure 1 compares the graphics processing unit (GPU) kernel distribution for ResNet-50 and three popular GNNs. It is well understood that CNNs are dominated by convolutional layers, which are implemented through general matrix multiplication (GEMM) on GPU. In contrast, the computation-intensive GEMM kernel is not the hotspot in the three GNNs, which also demonstrate model-specific patterns.

In this work, we aim to introduce GNN to our community. In contrast to prior work [4] that only presents

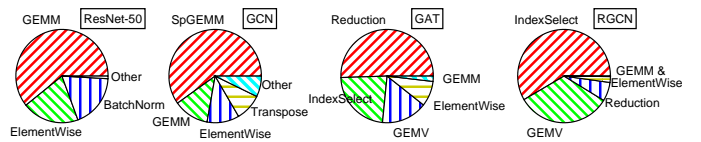


Fig. 1: Kernel breakdown of CNN (ResNet-50) and three GNN models.

characterizations of GCN [5], we select five representative GNN models that cover a large portion of the varieties for GNN workloads on the basis of a general GNN description framework and our model review. By constructing the models on top of two widely-used libraries, we characterize the efficiency of the GNN computation at inference stage on the existing GPU architecture and suggest directions for GNN-specific accelerators. We hope that our analysis can help architects and system designers have a better understanding of GNN computation and foster more future work.

2 GNN BENCHMARK SUITE CONSTRUCTION

This section describes our methodology of constructing a representative GNN benchmark suite. We use a general description framework to perform a detailed review of the recently published GNN models. The result identifies a few common patterns across the surveyed models, which lets us choose five models to cover almost all patterns.

Model Survey. A multi-layered GNN model is designed to learn the embeddings (i.e., vectors) for each vertex and edge in a graph. The input for a GNN layer is the graph structure in the form of adjacency matrix, together with the vertex (edge) embedding matrix $vertex^l$ ($edge^l$). The layer generates the transformed embedding matrix $vertex^{l+1}$ and

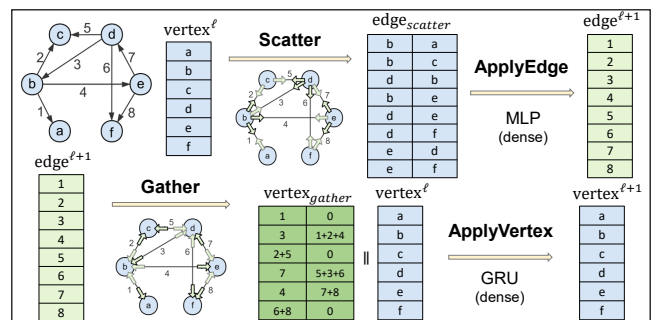


Fig. 2: The computation stages in a GGNN [6] layer.

• ¹Shanghai Jiao Tong University, ²Peking University, ³Microsoft Research Asia. *Corresponding authors: Jingwen Leng and Minyi Guo.

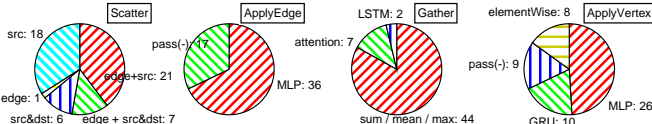


Fig. 3: The proportion of operations used in each stage of SAGA-NN.

TABLE 1: The SAGA-NN stage breakdown of the selected models.

Model	Scatter	ApplyEdge	Gather	ApplyVertex
GCN [5]	-	-	sum(in_edges)	MLP
GAT [10]	src&dst	MLP	attention(in_edges)	MLP
GGNN [6]	src&dst	MLP	sum(all_edges), vertex	GRU
R-GCN [11]	edge+src	MLP	sum(in_edges), vertex	sum
GraphSAGE [12]	-	-	LSTM(in_edges), vertex	MLP
Coverage	85%	100%	100%	100%

$edge^{l+1}$ for the next layer, as shown in Figure 2. We reviewed 53 GNN models published in recent top conferences and journals, spanning across molecular science, recommendation, and transportation domains. Since GNNs mix the graph and DL computation, the surveyed models demonstrate significant variability, which poses a great challenge for the analysis of their computational characteristics.

Model Decomposition. To overcome the diversity challenge, we adopt the recent GNN description framework SAGA-NN [7], which defines four stages in a GNN layer (Scatter, ApplyEdge, Gather, and ApplyVertex). The current popular GNN libraries DGL [8] and PyG [9] also implicitly follow a similar framework. We express and decompose the surveyed models into different stages in the SAGA-NN framework. We then categorize each stage’s operations to mine common patterns and simplify the analysis.

We use GGNN [6] in Figure 2 as an example to illustrate the different stages. In the first Scatter stage, each edge concatenates its source and destination vertex embeddings as $edge_{scatter}$. In the second ApplyEdge stage, each edge transforms $edge_{scatter}$ with a MLP operation to produce the new edge embedding $edge^{l+1}$. In the third Gather stage, each vertex first sums the $edge^{l+1}$ of all its inedges (incoming edges) to output $vertex_{gather}$, and then concatenates it with the vertex’s existing embedding $vertex^l$ as the input for the next stage. In the fourth ApplyVertex stage, each vertex transforms this input with a GRU (Gated Recurrent Unit) operation to the new vertex embedding $vertex^{l+1}$.

Survey Result. For the surveyed 53 models, we decompose them into different SAGA-NN stages. Figure 3 summarizes the operation distribution in each stage. The results show that *although GNNs have a tremendous design space (diversity), there exist a few commonly used operations in each stage.* GNNs can adopt a mix of source/destination vertex embedding, and edge embedding in the Scatter stage. For the Gather stage, some models use the simple sum/mean/max operation while others use more complex attention/LSTM (Long short-term memory) operations. Some models bypass the ApplyEdge or ApplyVertex, while most models use MLP operations. The survey insight suggests that we can cover the large GNN space by carefully selecting a few models.

Studied Models. In the end, we select five representative GNN models in Table 1, which cover almost 100% of

TABLE 2: The evaluated datasets for the vertex classification task.

Name	Cora	Citeseer	Pubmed	AIFB	MUTAG	BGS
Vertex#	2,708	3,327	19,717	8,285	23,644	333,845
Edge#	5,429	4,732	44,338	29,043	74,227	916,199
Graph Type	Citation Network	Citation Network	Citation Network	Semantic Network	Molecular Structure	Geological Graph

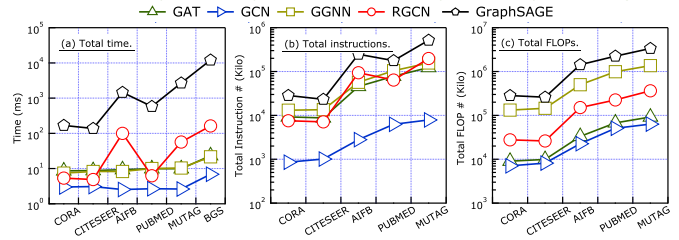


Fig. 4: Execution statistics of selected models with the different datasets.

the operations used in the 53 surveyed models. In other words, the 53 models can be viewed as recombination of different stages in the five models. We study and analyze those models on existing GNN execution library DGL [8] and PyG [9]. However, the two libraries do not explicitly implement SAGA-NN so we refactor the codes following the SAGA-NN description framework for a more systematic analysis. We focus more on DGL analysis because the PyG currently does not support GraphSAGE. We also select 6 commonly used datasets from the literature, with the number of vertices ranging from 5 K to about 1 M in Table 2.

3 GNN COMPUTATION ANALYSIS ON GPU

We study the computational characteristics of the selected GNN models on the contemporary GPU. We conduct an end-to-end analysis to examine the model difference and the impact of the dataset. We then analyze the stage-level characteristics in detail, which lets us identify the possible bottleneck and suggest new optimizations.

3.1 End-to-end Analysis

Figure 4 plots the execution statistics of our selected models with the different datasets (Table 2). The experiments are carried on a machine with two Intel Xeon 4110 CPUs and one NVIDIA RTX 2080Ti GPU. The machine runs Ubuntu 16.04 with CUDA 10 and cuDNN 7. We use GNN libraries DGL 0.4 [8] and PyG 1.3 [9], both with PyTorch 1.4 as backend. We summarize the key insights as follows.

Execution Time. Figure 4(a) plots the inference execution time with different graph structures that are ranked by their edge/vertex count. The inference time for GraphSAGE and RGCN generally increases with graph vertex/edge count except for the AIFB dataset. The reason for RGCN is that this dataset has additional edge types that lead to more computation in the two models, while the reason for GraphSAGE is its bottleneck stage Gather that we describe later. In contrast, the inference time for the rest three models do not vary except the largest graph BGS. The reason is that the CPU time dominates the execution when the graph is small.

Instruction & FLOP. Figure 4(b) and (c) compare the total instructions and FLOPs (floating-point operations), which increase as the graph size increases. When the graph is too small, the CPU time dominates the entire inference time so the increase of instructions and FLOPs are not reflected in Figure 4(a). Unlike the CNNs whose input is usually the fixed-size image, the size of graphs varies significantly across different domains and problem settings. As such, the design of GNN acceleration architecture must be aware of the graph size (also embedding vector size) of the targeted problems to balance resource utilization.

We also observe that GCN is the simplest model with

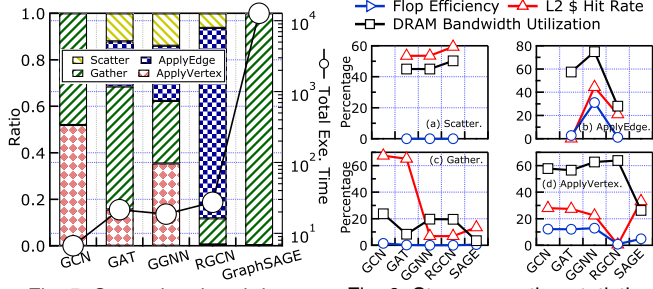


Fig. 5: Stage time breakdown. Fig. 6: Stage execution statistics.

less number of instructions and FLOPs. The rest four models have similar instruction counts for the same graph but their FLOPs are quite different. This variability can be attributed to their different model complexity: GraphSAGE is more complicated with the LSTM-based Gather stage while GGNN uses GRU cell in the ApplyVertex stage. Next, we perform a detailed stage-level analysis for a deeper understanding of these models.

3.2 Stage-level Analysis

Figure 5 shows the stage-level time breakdown and Figure 6 shows each stage’s execution statistics. Both results are obtained on the largest dataset BGS. The stage time distribution varies greatly among models, which confirms the diversity of our selected models, and also suggests that there is no fixed hotspot in GNNs. Therefore, we must equally and jointly consider all the stages. We present our detailed stage-by-stage analysis as follows.

Scatter. This stage prepares data for its following edge transformation stage ApplyEdge so that it only involves data movement. As such, Figure 6(a) shows that this stage has no floating-point operation and intensively uses DRAM bandwidth (note that GCN and GraphSAGE bypass this stage). Through the kernel trace profiling (not shown owing to space limits), we find that this stage uses the CUDA kernel `indexSelectLargeIndex` to implement the data movement, which puts together the embedding vector of each edge and the embedding vector of its two connected vertices into a new edge embedding vector (Figure 2). This stage also has a high L2 cache hit rate while graph processing is well known as low locality. The reason is that GNNs use a typical size of 32 or more for vertex/edge embedding while graph processing like BFS/PageRank uses one.

ApplyEdge. This stage performs edge embedding matrix transformation with MLPs. In general, it is possible to batch all edges for a parallel process so that this stage should have high computation efficiency. However, only GGNN has high FLOP efficiency while both GAT and RGCN have low-efficiency values. The size of edge embedding in GAT and GGNN is 1 (as attention value) and 32, respectively. As such, the ApplyEdge stage can be implemented by GEMV (matrix-vector multiplication) and GEMM in GAT and GGNN, respectively. This explains why the FLOP efficiency in GGNN is higher than the GAT. Meanwhile, the RGCN assigns another edge type attribute to different edges. As a result, its ApplyEdge stage first needs to select edges with the same type and then batches them to GEMM operation, which leads to the overall low FLOP efficiency.

Gather. This stage gathers edge embeddings for the following vertex transformation stage ApplyVertex. Since different vertices have different edge counts, this stage uses

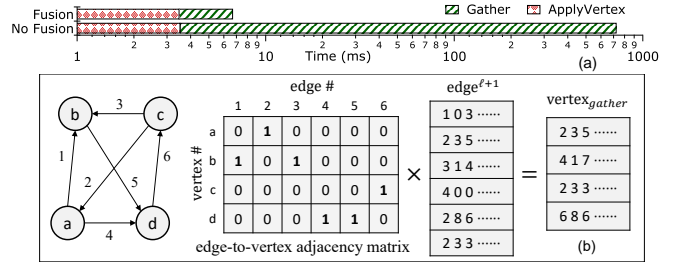


Fig. 7: (a) GCN time with/without fusion. (b) The fused Gather can be implemented by multiplying the graph adjacency matrix and edge embedding matrix.

a Reduction function that reduces the varying edge count to a fixed size embedding. Through profiling, we find that the Gather stage uses an important fusion optimization. To emphasize its importance, we run the GCN model with and without the fusion that DGL exposes. Figure 7-a shows that the stage fusion leads to about 200× speedup.

We use an example with the sum reduction function to illustrate how the fusion works. Assume an edge-to-vertex adjacency matrix (Figure 7-b middle), where a nonzero element represents an inedge (column) of the corresponding vertex (row). With the sum reduction function, the Gather stage calculates the new embedding of each vertex by summing the embeddings of all inedges. For each vertex, the summation can be implemented by multiplying its corresponding row in the adjacency matrix and the entire edge embedding matrix, where each row is an edge embedding vector. As such, the fused operation equals the multiplication between the adjacency matrix and the embedding matrix. Through the kernel trace profiling, we find that *this multiplication uses sparse GEMM for great performance and efficiency as the adjacency matrix is highly sparse*.

The exception is GraphSAGE which adopts LSTM as the reduction function, which treats the inedges of a vertex as a sequence and outputs a new vector. The LSTM requires a serialized input of a vertex’s inedges, which cannot leverage the fusion optimization and becomes the bottleneck (Figure 5). We find that DGL implements the LSTM-based Gather in a *degree traversal* technique. It batches the vertices with the same degree for parallel computation on the GPU. Figure 8-b shows that the number of kernel invocations equals the number of vertex degrees in the graph.

The degree traversal approach causes a severe GPU under-utilization as the idle time dominates in the stage (Figure 8-c). Figure 8-a shows the degree distribution in the BGS dataset, which obeys the power-law distribution: the number of vertices with the same degree decreases exponentially when the degree increases. As such, this approach is close to the sequential vertex-by-vertex computation for large degrees. Through profiling, we find that the interval time between two adjacent vertex computations can be 20× larger than the invoked kernel duration. As a result, the Gather becomes the major bottleneck in GraphSAGE.

ApplyVertex. This stage usually performs vertex embedding matrix transformation with MLPs, which is similar to the ApplyEdge stage. All vertices can be batched so that this stage can be implemented by GEMM, which leads to the highest FLOP efficiency among all stages (Figure 6d). The only exception of RGCN is due to the simple sum defined for the stage (Table 1). On the other hand, ApplyEdge and ApplyVertex often play a key role in the model accuracy, and more and more GNN models are being developed

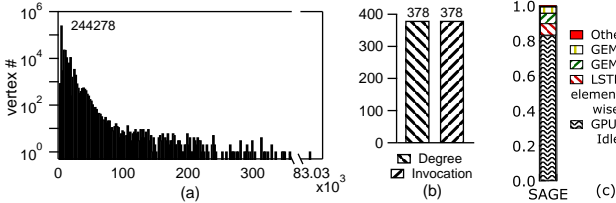


Fig. 8: (a) Degree distributions of BGS. (b) The number of degrees in the graph and invocations in the Gather stage. (c) Time distribution of kernels in the Gather stage of GraphSAGE.

with complex functions in those two stages to capture more information on the vertex, edge, and graph.

Library. We observe a similar result on PyG, especially for `ApplyEdge/ApplyVertex` that both leverage PyTorch’s existing features for efficient computation. Their greatest difference lies in `Scatter` and `Gather` stages. For example, PyG lets users customize the data movement in `Scatter` while DGL moves all data (embeddings of the edge and its vertices) by default. As such, PyG performs better when part of the data is fetched (e.g., RGCN). Moreover, the PyG does not support the LSTM-based `Gather` stage.

Summary. Table 3 summarizes the characteristics of each stage. It shows that although GNNs have a large design space, the stage-level characteristics are relatively stable across different models. In other words, our stage-level characterization can lay the foundation for future hardware/software acceleration of GNNs.

4 IMPLICATIONS FOR HARDWARE ACCELERATION

We now study the feasibility and challenges for the GNN accelerator, which can further improve the performance or energy efficiency of GNNs, on the basis of our previous analysis. Instead of building an accelerator from scratch, we take the existing DL accelerator TPU, and study its performance of running the GNN models. Based on the detailed stage-level analysis, we shed light on the efficient GNN accelerator design. We use TensorFlow 1.12 on top of one Google Cloud TPU v2 for this part of the experiments.

Projection Methodology. Because both DGL and PyG only support CPU and GPU platforms, we adopt a micro-benchmark based methodology to project the performance of running a typical GNN model on the TPU. We extract the parameters for the regular computation and run them on the TPU. For the irregular data movement, we use the time of our local CPU due to the lack of native TPU support. Figure 9 shows the result of two datasets. We report both the dense and sparse performance on the CPU/GPU, and report only the dense performance of TPU because it does not support the sparse GEMM.

Result Analysis. We show a counter-intuitive result in Figure 9. The TPU does not achieve the best performance on both datasets (Total bar), although it does achieve the best dense GEMM performance (ApplyV bar). On the small dataset Cora, TPU performs worse than the GPU because the TPU has to rely on CPU for the data movement. On the

TABLE 3: The characteristics of different GNN computation stages.

Stage	Description	Kernel
Scatter	Vertex/edge embedding movement	IndexSelection
ApplyEdge	DL-based edge embedding transformation	GEMM/GEMV
Gather (fused)	Edge embedding reduction	Sparse GEMM
Gather (non-fused)	Edge embedding movement Complex reduction (e.g., LSTM)	IndexSelection GEMM/GEMV
ApplyVertex	DL based vertex embedding transformation	GEMM/GEMV

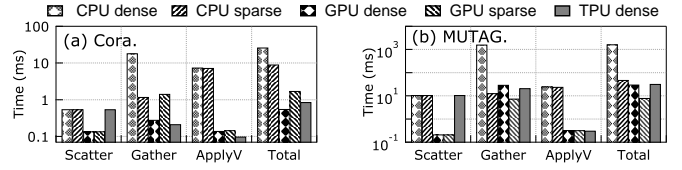


Fig. 9: The inference time (ms) on difference architectures.

large dataset MUTAG, the GPU outperforms TPU with the sparse stage fusion due to the higher graph sparsity.

Forward Looking. In summary, we think it is important for GNN accelerators to support both sparse and dense matrix operations for efficient GNN acceleration. Efficient data movement for the graph structure is also indispensable to avoid unnecessary data copy between processors. As such, an ideal GNN accelerator in our outlook consists of three key components: data movement component, dense operation component, and sparse operation component. Meanwhile, the current execution paradigm for GNN is stage-by-stage and layer-by-layer. We believe that there is a huge design space to break this sequential paradigm by, for example, by fine-grained pipelining the different stages. We leave those for our future work.

5 CONCLUSION

We systematically study the computation characteristics of graph neural networks. We first construct a representative GNN benchmark based on the extensive model review and a general GNN description framework. We then analyze their computational efficiency and microarchitectural characteristics on the existing GPU architecture. Our analysis suggests that the GNN is a unique workload with the mixed features from graph analytics and DL computation, which warrants more future research.

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their constructive feedback. This work was supported by National Key R&D Program of China (2018YFB1305900), the National Natural Science Foundation of China Grant (61702328, 61832006, and 61972247), Microsoft Research Asia Collaborative Research Grant. Any opinions, findings, and conclusions in this paper are those of the authors only.

REFERENCES

- [1] W. L. Hamilton *et al.*, “Representation learning on graphs: Methods and applications,” *IEEE Data Eng. Bull.*, 2017.
- [2] B. Perozzi *et al.*, “DeepWalk: online learning of social representations,” in *KDD*, 2014.
- [3] A. Grover *et al.*, “node2vec: Scalable feature learning for networks,” in *KDD*, 2016.
- [4] M. Yan *et al.*, “Characterizing and understanding GCNs on GPU,” in *IEEE Computer Architecture Letters*, 2020.
- [5] T. N. Kipf *et al.*, “Semi-supervised classification with graph convolutional networks,” in *ICLR*, 2017.
- [6] Y. Li *et al.*, “Gated graph sequence neural networks,” in *ICLR*, 2016.
- [7] L. Ma *et al.*, “NeuGraph: Parallel deep neural network computation on large graphs,” in *USENIX ATC*, 2019.
- [8] M. Wang *et al.*, “Deep graph library: Towards efficient and scalable deep learning on graphs,” 2019.
- [9] M. Fey *et al.*, “Fast graph representation learning with PyTorch Geometric,” in *ICLR Workshop*, 2019.
- [10] P. Velickovic *et al.*, “Graph attention networks,” in *ICLR*, 2018.
- [11] M. S. Schlichtkrull *et al.*, “Modeling relational data with graph convolutional networks,” in *ESWC*, 2018.
- [12] W. L. Hamilton *et al.*, “Inductive representation learning on large graphs,” in *NeurIPS*, 2017.