

GNNerator: A Hardware/Software Framework for Accelerating Graph Neural Networks

Jacob R. Stevens¹, Dipankar Das², Sasikanth Avancha², Bharat Kaul², Anand Raghunathan¹

Purdue University, West Lafayette¹

Intel²

{steven69, raghunathan}@purdue.edu

{dipankar.das, sasikanth.avancha, bharat.kaul}@intel.com

Abstract—Graph Neural Networks (GNNs) apply deep learning to inputs represented as graphs. They use fully-connected layers to extract features from the nodes/edges of a graph and aggregate these features using message passing between nodes, thereby combining two distinct computational patterns: dense, regular computations and sparse, irregular computations. To address the computational challenges posed by GNNs, we propose GNNERATOR, an accelerator with heterogeneous compute engines optimized for these two patterns. Further, we propose feature-blocking, a novel GNN dataflow that beneficially trades off irregular memory accesses during aggregation for regular memory accesses during feature extraction. We show that GNNERATOR achieves speedups of 5.7–37x over an NVIDIA RTX 2080-Ti, and 2.3x–3.8x over HyGCN, a state-of-the-art GNN accelerator.

Index Terms—neural network accelerators, graph neural networks

I. INTRODUCTION

Recently, graph neural networks (GNNs) have seen great success in achieving state-of-the-art results in a variety of applications, such as physics modeling [1], chemical synthesis [2], and electronic design automation (EDA) [3]. These GNNs consist of two main stages: the *feature extraction* stage and the *aggregation* stage. The feature extraction stage is typically a fully-connected network whose weights are shared across all nodes. In the aggregation stage, each node in the graph aggregates features from its neighbors into a new feature representation, resulting in sparse, random memory accesses. Since each node is performing this aggregation independently, there is ample *inter-node* parallelism. Further, in contrast to traditional graph-processing workloads such as PageRank where node features tend to be a few bytes, the features in GNNs can be thousands of bytes long. Due to this high dimensionality, there is also ample *intra-node* parallelism, since the aggregation of each dimension can be done in parallel. Finally, the feature extraction and aggregation stages have a producer-consumer relationship, resulting in *inter-stage* parallelism.

As a result of GNNs’ unique computational characteristics, traditional DNN accelerators are not well suited for GNNs. For example, DNN accelerators such as Eyeriss [4] or Google’s TPU [5] are not optimized for the memory patterns present in

the graph-based aggregation step, leading to poor utilization of on-chip resources. Similarly, graph analytics accelerators such as GraphP and Tesseract [6], [7] are also ill-suited for GNNs, as these accelerators lack the compute resources for the dense, matrix-based operations found in the feature extraction stage.

Due to the aforementioned challenges, there have been a few recent efforts to design GNN-specific accelerators. In particular, HyGCN [8] uses heterogeneous compute units for the feature extraction and aggregation stages, which are then tightly coupled together. This architecture allows HyGCN to address both major kernels described above. However, HyGCN has some drawbacks. Specifically, HyGCN fails to fully exploit inter-node parallelism, since it processes a single node at a time. Further, HyGCN only supports a dataflow in which the aggregation stage is the producer and the feature extractor is the consumer. This limits its applicability to workloads such as GraphsagePool, where this is not the case.

To overcome the aforementioned limitations, we propose GNNERATOR a Graph Neural Network accelerator that exploits all three forms of parallelism found in GNNs: inter-node, intra-node, and inter-stage. GNNERATOR provides flexible, fine-grained control of how data flows between its dense engine and graph engine, supporting a wider variety of GNN topologies. We also propose a novel feature dimension-blocking dataflow. This dataflow improves upon the dataflows used in previous GNN accelerators, reducing the overhead associated with processing the irregular, feature aggregation stage by exploiting the independent nature of the feature dimensions. We evaluate GNNERATOR on a suite of nine GNN benchmarks and demonstrate that it outperforms a GPU baseline, as well as HyGCN.

In summary, we make the following contributions:

- We design GNNERATOR, a programmable graph neural network accelerator consisting of a Dense Engine and a Graph Engine, exploiting the inter-stage parallelism inherent in GNNs while allowing for flexibility in the producer-consumer relationship between them.
- We provision the Dense and Graph Engines to exploit both the inter- and intra-node parallelism abundant in GNNs.
- We propose a novel, feature-dimension blocking dataflow for GNNs and provide hardware support for this dataflow in GNNERATOR.

This work was supported in part by C-BRIC, one of six centers in JUMP, a Semiconductor Research Corporation (SRC) program sponsored by DARPA

- We develop a simulation framework that implements the above proposals and use this framework to demonstrate the benefits of GNNERATOR. Our experiments indicate that GNNERATOR achieves an average speedup of 8x over an NVIDIA 2080-Ti GPU and an average speedup of 3.15x over HyGCN, a recent state-of-the-art GNN accelerator.

II. BACKGROUND

A. Graph Neural Networks

The term Graph neural networks (GNNs) is used to collectively refer to a diverse family of networks. Many popular GNNs [9], [10] consist of two distinct stages: a *feature extraction* stage, and an *aggregation* stage. In the feature extraction stage, the node features are passed through one or more fully-connected linear layers. In the aggregation stage, a node aggregates feature vectors from its neighbors and then uses this aggregated feature vector to update its own feature.

These two stages are combined to make up a single GNN layer. Either stage may precede the other. A full GNN can then be constructed by stacking these layers. By stacking layers, a GNN incorporates information from nodes that are increasingly far away from the original node. For example, a single layer GNN only considers a node’s neighbors, while a two layer GNN will consider nodes in the two-hop neighborhood.

We provide a brief overview of two popular GNNs below in order to further illustrate this concept.

Graphsage. The Graphsage network [11] is an extremely popular example of GNNs. Graphsage applies the local weight sharing found in traditional convolutional neural networks to graphs such that every node (edge) feature shares the same weights. As shown in Equation (1), Graphsage first applies an aggregation stage, where the mean of the features of a node’s neighbors is calculated. The aggregated node feature \bar{z} is then passed through a linear layer to obtain the updated node feature, h'_u .

$$\bar{z} = \frac{1}{|N(u)|} \sum \{h_v | v \in N(u) \cup u\} \quad (1)$$

$$h' = \sigma(\mathbf{W} \cdot (\bar{z} \cup \mathbf{h}))$$

GraphsagePool. The GraphsagePool variant of Graphsage replaces the mean-based aggregator computation of \bar{z} with a symmetric, trainable aggregator, as described in Equation (2). Specifically, each node’s feature is fed through a linear layer, defined by W_{pool} . The resulting features are then aggregated using an element-wise pooling operation (typically, max). Note that in this case, the feature extraction for z is consumed by the aggregation for \bar{z} .

$$z = \sigma(\mathbf{W}_{pool} \cdot \mathbf{h}) \quad (2)$$

$$\bar{z} = \max(\{z'_v | v \in N(u) \cup u\})$$

B. Graph Sharding

Many real-world graphs are too large to fit into a given level of the memory hierarchy, severely degrading the performance

of graph processing algorithms, as it makes it extremely difficult to exploit any potential reuse. This is especially true when high-dimensional features are associated with the nodes/edges. To address this challenge, graphs are often broken into smaller pieces, such that each subgraph can fit in a given level of the memory hierarchy. This process is referred to as *graph sharding* [12]. Similar to [12], we adopt a two-dimensional sharding paradigm, as depicted in Figure 1. In this paradigm, a graph’s edge list is divided into shards such that each shard contains a maximum of n^2 edges, where n is a tunable parameter that limits the number of source/destination nodes per shard.

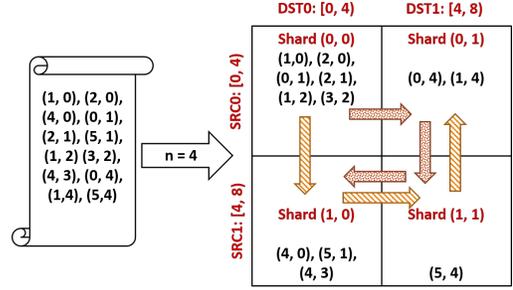


Fig. 1: A graph sharding algorithm divides an edge list into shards (subgraphs) which can then be processed in either a source-stationary (dotted arrow) or destination-stationary (slanted arrow) manner.

III. GNNERATOR ARCHITECTURE

The GNNERATOR architecture, presented in Figure , consists of two different compute engines, the Dense Engine and the Graph Engine. The Dense Engine is used for performing the dense, regular feature extraction steps of a GNN, while the Graph Engine performs the graph-based aggregation steps. The interoperation of these engines is controlled by the GNNERATOR Controller.

A. Dense Engine Overview

The Dense Engine consists of a two-dimensional systolic array-based matrix multiplication unit, an activation unit, and

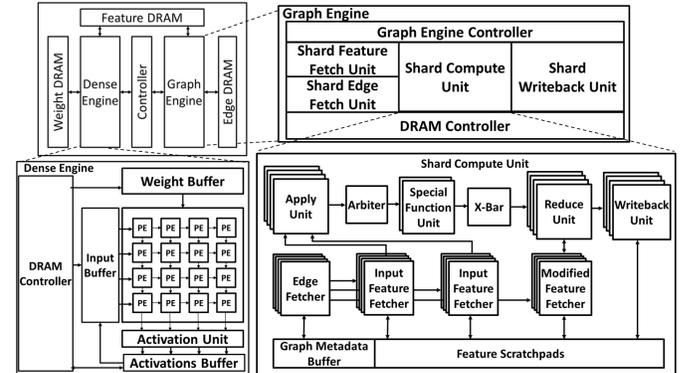


Fig. 2: GNNERATOR (top left) consists of two heterogeneous compute engines, a Dense Engine and a Graph Engine, that share feature storage.

input, weight, and activation on-chip double-buffered scratchpads, much like traditional DNN accelerators. The input and weight buffers feed the systolic array. The output of the two-dimensional systolic array is connected to a one-dimensional activation unit, which performs any required activation operations (*e.g.*, ReLu). The results from the activation unit are stored in the output buffer, where they can either be transferred out to DRAM or to the input buffer to be reused as input to the systolic array. We note that, unlike HyGCN’s combination engine, the Dense Engine has its own memory controller. This is necessary to support GNNs where the Dense Engine needs to act as the producer. It also enables reloading of partial sums, which in turn facilitates the novel dataflow proposed in the following section.

In order to support fine-grain pipelining of the feature extraction and aggregation stages, the Dense Engine also contains connections with the Graph Engine that are used to communicate the current state of the respective computing engines.

B. Graph Engine Overview

The Graph Engine is tailored for the irregular nature of graph processing. It consists of four major units— Shard Feature Fetch, Shard Edge Fetch, Shard Compute, and Shard Writeback Units— as well as a Graph Engine Controller that orchestrates the processing pipeline between the four units.

The Shard Edge Fetch and Feature Fetch Units work in parallel to load the edge data and feature data required to process a given graph shard from main memory to the on-chip scratchpads. After a shard is loaded, the Shard Compute Unit steps through its associated edges and performs the necessary computations. Finally, the Shard Writeback Unit stores output data from the on-chip scratchpads to main memory. As in the Dense Engine, all of the on-chip buffers in the Graph Engine are also double-buffered, enabling the pipelining of the above computations, such that the next shard is being prefetched while the current shard is being executed.

The Shard Compute Unit contains an Edge Fetcher, which steps through the edges associated with the current graph shard and distributes the required edge information (*e.g.* source node ID) to the Feature Fetcher Units, as well as the Writeback Unit. These units use the edge information in order to generate read (write) accesses to the on-chip scratchpad. The fetch units feed the actual compute units: the Apply Unit, which performs binary operations, and the Reduce Unit, which performs an aggregation operation. Each of these compute units are vectorized in a single-instruction, multiple data (SIMD) manner in order to exploit *intra-node parallelism* across the different dimensions of a given node’s feature, which are computed independently.

In order to exploit *inter-node parallelism*, the Shard Compute Unit contains multiple copies of the set of units described above, referred to collectively as a Graph Processing Element (GPE). Each GPE is assigned to a subset of the edges for a given graph shard, thereby processing multiple nodes at one time.

C. GNNerator Controller

The GNNerator Controller coordinates the interaction between the Dense Engine and the Graph Engine. Crucially, our controller allows each engine to be either the producer or the consumer. This is important to efficiently support various GNN configurations, as in some networks the computation is a feature extraction followed by an aggregation and sometimes it is the other way around.

Dense first. If the feature extraction is first, then the Dense Engine must run ahead of the Graph Engine. Hence, the GNNerator Controller reads the state of the Dense Engine and stalls the Graph Engine until the source nodes for the current shard of the Graph Engine have been processed by the Dense Engine.

Graph first. If aggregation is first, then the Graph Engine must run ahead of the Dense Engine. Hence, the GNNerator Controller reads the state of the Graph Engine and stalls the Dense Engine until the Graph Engine is done processing a set of destination nodes (*i.e.*, the Graph Engine has completed a full column of the shard grid).

IV. GNN DATAFLOWS

In this section, we first describe the conventional dataflow used for executing graph neural networks, and then propose a novel feature dimension-blocking GNN dataflow, which we implement in GNNERATOR. The conventional GNN dataflow can be expressed using our proposed dataflow in Algorithm 1, with B set to the length of the node feature, D , thus eliminating the second loop (line 2).

A. Conventional GNN Dataflow

In order to process a given sharded graph, GNNERATOR must step through the two-dimensional grid depicted in Figure 1. This can be done in a source-major or destination-major manner. This corresponds to the third and fourth loops in Algorithm 1 (lines 3-4); in this case, we are processing in a destination-major manner.

When traversing in a source-major fashion (*i.e.*, across a row of the shard grid), a set of *source* vertices and their corresponding feature(s) are loaded on-chip and remain on-chip for the entire row. The destination vertices, however, must be written back and reloaded as we move from shard to shard. Conversely, when traversing in a destination-major fashion, a set of *destination* vertices and their corresponding feature(s) are loaded on-chip and remain on-chip until they are done aggregating, while the source features must be reloaded as we move from shard to shard. The loading/storing of these features are performed by the Shard Feature Fetch Unit and the Shard Writeback Unit, respectively.

Assuming an S-pattern, we show the read and write costs associated with the two different orders in Table I, where S is the number of shards and I is the maximum number of input features required to be on-chip at one time. With these costs, assuming equal costs for read and write transactions, we can analytically determine the best ordering.

Algorithm 1: Dimension-blocking Algorithm

Input: Sharded Graph G ; Width/Height of Square Shard Grid S , Hidden Dimension Size D , Features h , Layers L

```
1 for  $l$  in range ( $L$ ) do
2   for  $blockD$  in range ( $D/B$ ) do
3     for  $dst$  in range ( $S$ ) do
4       for  $src$  in range ( $S$ ) do
5          $Shard = G.Shards(src, dst)$ 
6         for  $v$  in range ( $Shard(src, dst).V$ ) do
7           for  $u$  in range ( $v.U$ ) do
8             for  $d$  in range ( $B$ ) do
9                $dim = f(d, blockD)$ 
10               $h_{agg}[v][dim] =$ 
11                 $Aggregate(h_u[dim], h_v[dim])$ 
12               $h'[dst][:] =$ 
13                 $FeatureExtract(h_{agg}[dst][blockD * B :$ 
14                   $(blockD + 1) * B], h'[dst][:])$ 
15             $h = h'$ 
```

TABLE I: Analytical description of shard dataflows

	Read Cost	Write Cost
SRC Stationary	$S * I + (S - 1) * S - S + 1$	$S^2 - S + 1$
DST Stationary	$(S^2 - S + 1) * I$	S

Within each sub-graph shard, the edge list that defines the sub-graph must also be processed. The Edge Fetcher within the Shard Compute Unit steps through this edge list (corresponding to lines 6-7 in Algorithm 1), and uses the list to program the Input Feature Fetcher Units, which perform the requisite loads. Finally, the nodes are aggregated (line 10) by the Apply and Reduce Units. Note that the Modified Feature Fetcher is used when reloading partial computed accumulations is required. Finally, after aggregation, feature extraction (line 12) is performed by the Dense Engine. We note that in the traditional dataflow, wherein $B = D$, this is done in one step without the need to reload partial sums.

B. Feature Dimension-blocking

We first note that the feature dimensions are treated *independently* during the graph processing phase of GNNs. With this insight in mind, we propose, and provide hardware support for, a novel feature-blocking dataflow (Algorithm 1) wherein the graph processing computations are broken into two loops, such that only a block of dimensions is on-chip and processed at one time. A block is processed for every subshard before moving on to the next block of dimensions.

There are two main implications of the proposed dataflow. First, as shown in Table I, the cost associated with read and writing the subshards is dependent on S , the number of shards. Thus, to minimize data transfers, we would like to maximize the number of nodes that can be held on-chip at one time – that is, we would like the shards to be as large as possible. However, large shards are difficult to fit on-chip in the context

of the traditional GNN dataflow, since each node or edge is associated with one or more features, each of which can be of a high dimension and must be held on-chip in the traditional dataflow. By contrast, in our proposed dataflow, only a subset of the dimensions are kept on chip at any time. Since each node requires less on-chip storage, more nodes can be stored on-chip, reducing S in Table I. Note, however, that the edge list is processed multiple times (*i.e.*, lines 3-4), thus increasing the number of on-chip memory accesses. This is an overhead associated with feature-blocking, but it is favorably offset by the reduction in off-chip data transfers.

Second, since the traditional dataflow computes on entire features at once, the Dense Engine must load the weights for the entire feature, which are then shared across a relatively smaller number of nodes, for the reasons described above. In our dataflow, by contrast, the Dense Engine computes on more nodes, with fewer features per node at a time. This increases reuse for the Dense Engine, improving efficiency and helping to mitigate the overheads of the new need to reload partial sums during feature extraction (line 12).

In effect, when comparing the conventional and proposed dimension-blocked dataflows, the irregular off-chip accesses of feature aggregation are reduced at the cost of additional regular accesses for feature extraction and additional on-chip accesses for processing the edges. As borne out by our results, this is a beneficial tradeoff. Moreover, the additional feature extraction accesses for partial sums are mitigated by the increased reuse enable by dimension-blocking. Finally, feature extraction is performed on only destination shards whereas aggregation depends on both source and destination shards, improving the benefits of dimension-blocking.

V. METHODOLOGY

Simulation infrastructure. We developed a cycle-level simulator as well as a prototype compiler and runtime for GNNERATOR. The cycle-level simulator was developed using the PyMTL3 framework [13]. We implemented cycle-level models of all of the Graph Engine and GNNerator Controller components shown in Figure 2 and integrated the cycle-accurate SCALE-Sim simulator [14] for the Dense Engine.

TABLE II: Summary of Graph Datasets

Dataset	Vertices	Edges	Feature Dim.	Size
CORA	2708	10556	1433	15.6 MB
CITeseer	3327	9104	3703	49 MB
PUBMED	19717	88648	500	40.5 MB

TABLE III: Summary of Graph Neural Networks

Network	Hidden Layers	Hidden Dimension
GCN [15]	1	16
Graphsage [11]	1	16
GraphsagePool [11]	1	16

Benchmarks. In Table II, we summarize the input graph datasets used in our experiments. These datasets represent

standard graph datasets used in GNNs. Note that most of the datasets cannot fit on-chip due to the large feature dimension sizes. We run each of these input graph datasets on the three graph neural network architectures outlined in Table III using the Deep Graph Library (DGL) [16] with the PyTorch backend.

Platforms. Table IV presents the configuration of GNNERATOR used in our evaluations, along with the relevant parameters for the baseline platforms: an NVIDIA RTX 2080 Ti GPU and HyGCN, a recently proposed state-of-the-art GNN accelerator. HyGCN was chosen as the baseline as it is the most similar design to GNNERATOR in terms of hardware resources (FLOPs and bandwidth).

VI. EVALUATION

A. Performance

Figure 3 shows the normalized speedup with respect to the 2080-Ti GPU. We consider two variants of GNNERATOR: (i) the standard baseline GNNERATOR that uses the proposed novel dimension-blocking scheme described in Section IV with $blockD$ set equal to the width of the Dense Engine (*i.e.*, 64) and (ii) a variant that does not use dimension blocking. Without feature dimension-blocking, GNNERATOR demonstrates an average speed up of $4.2\times$ over the GPU baseline, while with feature dimension-blocking it has an average speedup of $8.0\times$. Without feature blocking, GNNERATOR’s speedup comes from a variety of sources. First, as discussed in Section III, GNNERATOR fully exploits all potential sources of parallelism in GNNs: inter-stage, inter-node, and intra-node parallelism. Further, unlike a GPU, the memory hierarchy of GNNERATOR’s Graph Engine is specialized for the unique needs of graph processing. For example, the memory width of the edge memories are sized such that no bandwidth is wasted due to mismatches between edge data size and memory width. GNNERATOR’s additional performance improvement through the use of dimension-blocking stems from two main sources. First, dimension-blocking allows for more nodes’ features to be held on-chip, reducing the memory bottleneck of transferring these features on- and off-chip. Second, dimension-

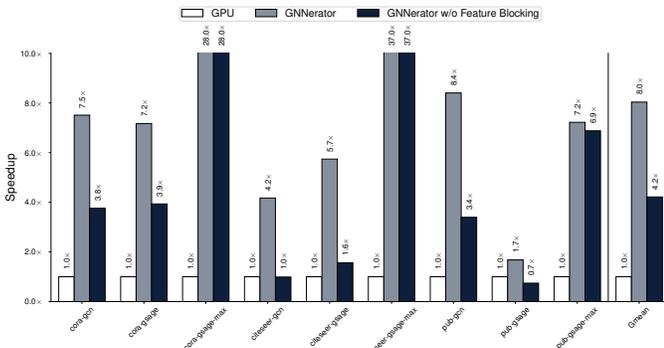


Fig. 3: GNNerator achieves an $8\times$ speed up over the 2080-Ti baseline. Roughly half of this speedup results from the specialized architecture and the other half comes from the novel dimension blocking dataflow.

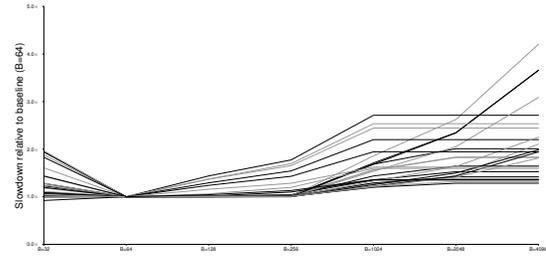


Fig. 4: Generally, a smaller feature block size B is desirable. However, performance suffers if B is less than the width of the systolic array of the Dense Engine (*i.e.*, $B = 32$)

blocking reduces the amount of time the Dense Engine must wait for the Graph Engine to finish aggregating a node’s neighborhood, as the Graph Engine only has to aggregate a small fraction of the dimensions before the Dense Engine can begin.

HyGCN Comparison. Table V compares GNNERATOR to the recently proposed state-of-the-art GNN accelerator HyGCN using their results described in [8]. We chose HyGCN as the comparison platform, as it is closest to GNNERATOR in terms of hardware resources. As demonstrated in the table, without feature dimension-blocking, GNNERATOR and HyGCN are quite similar in performance, with comparable performance improvements over the GPU baseline. We note that HyGCN uses a sparsity elimination optimization which is particularly effective for the Citeseer dataset (roughly, $1.1\times$ for Cora/Pubmed, $3\times$ for Citeseer). This optimization is orthogonal to our work and can be added to GNNERATOR. When utilizing dimension-blocking, however, GNNERATOR consistently and considerably outperforms HyGCN, with an average speedup of $3.15\times$.

Feature-blocking Ablation. A key question regarding our novel feature-blocking is the optimal size of the feature block, B —that is, how many dimensions of the feature will be kept on-chip at once. To explore this, we performed a sweep of this parameter, while running a large number of various networks and datasets. As seen in Figure 4, a smaller B is desirable. However, there is a lower limit to this. If the block size is set to the smallest possible value (*i.e.*, the width of the Graph Engine lanes), the performance suffers. This is because the block size is then less than the width of the Dense Engine systolic array of sixty-four, leading to under-utilization of the engine’s resources.

B. Scaling GNNerator

Finally, we consider the scenario of scaling GNNERATOR to larger chip sizes and examine the question of where to invest the additional hardware resources in order to maximize the performance return on that investment. We present three possible versions of a “next-generation” GNNERATOR. One version doubles the amount of on-chip memory in the Graph Engine, allowing for more larger shards to be held on-chip. Another version doubles both the height and width of the Dense Engine, increasing the compute available for the linear

TABLE IV: Summary of Compute Platforms

	RTX 2080 Ti	GNNerator	HyGCN
Peak Compute	13 TFLOPs	10 TFLOPs (2 for Graph, 8 for Dense)	9TFLOPs (1 for Graph, 8 for Dense)
On-chip Memory	29.5 MiB	30 MiB (24 MiB Graph, 6 MiB Dense)	24 MiB
Off-chip Memory	616 GB/s	256 GB/s	256 GB/s
Area	775 mm ²	14.5 mm ²	7.8 mm ²

TABLE V: Speedups of GNNERATOR over HyGCN for GCN

	Cora	Citeseer	Pubmed
GNNerator w/o blocking	1.8x	0.8x	1.0x
GNNerator	3.8x	3.2x	2.3x

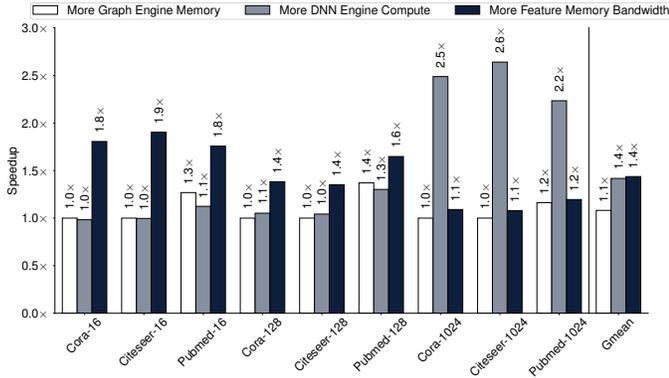


Fig. 5: Adding feature memory bandwidth tend to improve performance for networks with smaller hidden dimension size while more DNN Engine compute results in the largest speedups for networks with large hidden size.

layers. The final version doubles the bandwidth available for the shared feature memory DRAM.

We find that the best investment depends on the target networks. For networks with smaller hidden dimension sizes, increasing the feature memory bandwidth provides the highest return. At larger hidden sizes, on the other hand, increasing the size of the Dense Engine results in the largest speedups. This is driven primarily by large speedups for the larger hidden dimension sizes, as seen in Figure 5.

VII. DISCUSSION AND RELATED WORK

Hardware acceleration of neural networks has been a very active topic in the past decade. However, very few efforts directly target graph neural networks. The closest effort to GNNERATOR is HyGCN [8], which proposes two heterogeneous compute engines, one optimized for the feature extraction stage and one for the aggregation stage. Our work differs from HyGCN in a few major ways. First, unlike in our work, HyGCN only exploits intra-node parallelism, processing a single node’s entire feature across all cores before moving on to the next node. In contrast, GNNERATOR exploits both intra-node and inter-node parallelism. Second, GNNERATOR’s dual engine architecture is more flexible. For example, in GNNERATOR, the Graph Engine can be either the producer or the consumer— in HyGCN, the Aggregation Engine is the only producer. This directly limits the applicability of HyGCN in workloads where the Dense Engine is the producer, such as GraphsagePool. Finally, HyGCN does not contain the architectural support necessary for feature blocking.

Due to the wide variety in hardware resources used and benchmarks evaluated, it is difficult to compare GNNERATOR to other recent GNN proposals such as GNNa [17] and EnGN [18]. However, these architectures rely on the traditional GNN dataflow and can therefore benefit from our proposed feature blocking-based dataflow with suitable enhancements.

VIII. CONCLUSION

GNNs are a promising new area of machine learning that aim to bring the success of deep learning from Euclidean domains to graph-based inputs. In this work, we detail the limitations of current hardware in efficiently realizing GNNs and propose GNNERATOR, a specialized hardware accelerator for GNNs that is able to exploit the abundant intra- and inter-node parallelism inherent in GNNs. GNNERATOR utilizes feature dimension-blocking, a novel GNN dataflow that allows for processing more nodes in a graph on-chip at one time. We evaluate GNNERATOR on a suite of benchmarks and demonstrate significant performance benefits over GPUs, as well as a recently proposed GNN accelerator.

REFERENCES

- [1] P. Battaglia, R. Pascanu, M. Lai *et al.*, “Interaction networks for learning about objects, relations and physics,” in *Proc. of NeurIPS*, 2016.
- [2] Y. Li, O. Vinyals, C. Dyer *et al.*, “Learning deep generative models of graphs,” *arXiv preprint arXiv:1803.03324*, 2018.
- [3] G. Zhang, H. He, and D. Katabi, “Circuit-gnn: Graph neural networks for distributed circuit design,” in *Proc. of ICML*, 2019, pp. 7364–7373.
- [4] Y.-H. Chen, T. Krishna, J. S. Emer, and V. Sze, “Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks,” *IEEE Journal of Solid-State Circuits*, 2017.
- [5] N. P. Jouppi *et al.*, “In-datacenter performance analysis of a tensor processing unit,” in *Proc. of ISCA*, 2017.
- [6] M. Zhang, Y. Zhuo, C. Wang *et al.*, “Graphp: Reducing communication for pim-based graph processing with efficient data partition,” in *Proc. of HPCA*. IEEE, 2018, pp. 544–557.
- [7] J. Ahn, S. Hong, S. Yoo, O. Mutlu, and K. Choi, “A scalable processing-in-memory accelerator for parallel graph processing,” in *Proc. of ISCA*, 2015, pp. 105–117.
- [8] M. Yan *et al.*, “HygcN: A gcN accelerator with hybrid architecture,” in *Proc. of HPCA*. IEEE, 2020, pp. 15–29.
- [9] R. Ying *et al.*, “Graph convolutional neural networks for web-scale recommender systems,” in *Proc. of ICKDDM*. ACM, 2018.
- [10] P. Veličković, G. Cucurull, A. Casanova *et al.*, “Graph attention networks,” *arXiv preprint arXiv:1710.10903*, 2017.
- [11] W. Hamilton, Z. Ying, and J. Leskovec, “Inductive representation learning on large graphs,” in *Proc. of NeurIPS*, 2017, pp. 1024–1034.
- [12] X. Zhu, W. Han, and W. Chen, “Gridgraph: Large-scale graph processing on a single machine using 2-level hierarchical partitioning,” in *USENIX*, 2015, pp. 375–386.
- [13] S. Jiang, B. Ilbeyi, and C. Batten, “Mamba: closing the performance gap in productive hardware development frameworks,” in *Proc. of DAC*. IEEE, 2018, pp. 1–6.
- [14] A. Samajdar, Y. Zhu, P. Whatmough, M. Mattina, and T. Krishna, “Scale-sim: Systolic cnn accelerator simulator,” *arXiv preprint arXiv:1811.02883*, 2018.
- [15] T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” *arXiv preprint arXiv:1609.02907*, 2016.

- [16] M. Wang, L. Yu, D. Zheng *et al.*, “Deep graph library: Towards efficient and scalable deep learning on graphs,” *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [17] A. Auten, M. Tomei, and R. Kumar, “Hardware acceleration of graph neural networks,” in *Proc. of DAC*. IEEE, 2020, pp. 1–6.
- [18] S. Liang *et al.*, “Engn: A high-throughput and energy-efficient accelerator for large graph neural networks,” *IEEE Transactions on Computers*, 2020.