

# SPA-GCN: Efficient and Flexible GCN Accelerator with an Application for Graph Similarity Computation

Atefeh Sohrabizadeh, Yuze Chi, and Jason Cong

Computer Science Department, University of California, Los Angeles, USA  
{atefehsz,chiyuze,cong}@cs.ucla.edu

## Abstract

While there have been many studies on hardware acceleration for deep learning on images, there has been a rather limited focus on accelerating deep learning applications involving graphs. The unique characteristics of graphs, such as the irregular memory access and dynamic parallelism, impose several challenges when the algorithm is mapped to a CPU or GPU. To address these challenges while exploiting all the available sparsity, we propose a flexible architecture called SPA-GCN for accelerating Graph Convolutional Networks (GCN), the core computation unit in deep learning algorithms on graphs. The architecture is specialized for dealing with many small graphs since the graph size has a significant impact on design considerations. In this context, we use SimGNN, a neural-network-based graph matching algorithm, as a case study to demonstrate the effectiveness of our architecture. The experimental results demonstrate that SPA-GCN can deliver a high speedup compared to a multi-core CPU implementation and a GPU implementation, showing the efficiency of our design.

## 1 Introduction

Graphs are the core data structure used in datacenters and have a wide application in different domains such as recommender systems, social networks, and the World Wide Web. Although they are widely used, they are mainly unstructured and have a high dimensionality, making them computationally expensive to process. In fact, many graph algorithms, such as Graph Edit Distance [46] and Maximum Common Subgraphs [9] are known to be NP-complete [29, 75]. This problem has motivated researchers to apply deep learning on graphs with the goal of extracting structured, low-dimensional features from it (e.g., [31, 72]). More specifically, the purpose of this line of research is to assign a feature vector to each node of the graph that can show the role of the node in the graph. Such feature vectors are called the *node embeddings*.

In this context, Graph Convolutional Networks (GCN) [31] are widely used to extract the node embeddings in a graph. GCNs follow the same behavior as Convolutional Neural Networks (CNN) in learning. They consist of multiple layers in which the features of the nodes are propagated within them until rich information of the input graph is derived. Like in a CNN, in each layer, the GCN updates the node features by gathering the neighbors' features and passing their summation through a filter. GCNs have shown to be successful in many domains including traffic prediction [80], facilitating web-scale recommender systems [72], molecular footprint calculation [19], logic optimization for EDA tools [22], etc.

While some graph data tend to scale rapidly, there are also many graph data that are naturally limited in size, for example, chemical compounds and molecules [6, 7, 10, 40, 54] that have a wide

application in different domains including drug development, quantum mechanics, physical chemistry, biophysics, etc [10, 64]; the LINUX dataset containing operating system program dependence graph [57]; the GREC database consisting of graphs representing symbols from architectural and electronic drawings [18]; the Fingerprint database [61], etc [44, 64]. The average number of nodes for the graphs of these databases ranges from 5 to 50.

Because of the vast application of small graphs, numerous algorithms have been proposed to obtain their information by extracting their features or learning how similar two graphs are [3, 11, 26, 30, 33, 36, 43], especially using GCNs [3, 26, 33, 36, 43]. In particular, SimGNN [3] proposed a GCN-based approach to learn a similarity score for such graphs. It demonstrates that a GCN-based approach is able to approximate the GED with high accuracy; hence, it expedites the graph similarity computation significantly for many applications. SimGNN targets searching/matching graphs from real-world graph databases, such as AIDS [40], LINUX [57], and IMDB [70], which consists of antivirus chemical compounds, program dependency graphs, and actor/actress ego-networks, respectively. The generated target graphs are relatively small, with 10 nodes on average, but the database contains millions of graph pairs, creating a large number of graph matching queries. Although the CPU implementation can finish each SimGNN query for such graphs in milliseconds (5.8ms to be exact), processing millions of queries can take several hours, which is not acceptable and requires customized acceleration. Such a workload of graph searching/mining is increasing in importance. For example, searching for antivirus chemical compounds is an important step in drug repurposing for COVID-19.

Despite the popularity and effectiveness of *graph neural network* (GNN) approaches, there has been limited research on developing an accelerator for them. Besides, the differences between an image and a graph structure, as explained below, make the countless CNN accelerators proposed in the literature (e.g., [12, 28, 37, 39, 47–51, 62, 63, 73, 77, 79]) incompatible. Furthermore, the characteristics of the graphs impose several challenges when the algorithms are executed using CPUs or GPUs. More specifically, the unique features of GNN impose the following challenges in designing an accelerator:

- **Irregular memory access and low data reuse:** As opposed to images in which the neighbors of a pixel are stored either in contiguous locations or have a fixed distance to one another, the neighbors of a node in a graph may be stored in any location in memory. This will result in many irregular memory accesses to all levels of the memory hierarchy. To make matters worse, GNNs have much lower data reuse compared to CNNs. These characteristics make the application memory-bounded. Compared to the traditional graph algorithms such as breadth first search (BFS), single-source shortest path (SSSP), PageRank (PR),

**Table 1: Our approach compared to state-of-the-art GCN accelerators.**

Work	Graph Size	Layer Customization	Sparse Engine for Feature Transformation	Read Each Element Only Once	Parallelization			Batch
					Inter-layer	Feature-level (Sparse Part)	Node-level (Sparse Part)	
HyGCN [69]	Large	✗	✗	✗	✗	✓	✓	✗
GraphACT [74]	Small	✗	✗	✗	✗	✓	✗	✗
ASAP'20 [76]	Large	✗	✗	✗	✗	✓	Limited	✗
AWB-GCN [21]	Large	✓	✓	✗	✓	✗	✓	✗
<b>Ours (SPA-GCN)</b>	Small	✓	✓	✓	✓	✓	✓	✓

etc., the nodes have long feature vectors instead of a single scalar value. As a result, although they both have irregular memory access, the access pattern is different. Also, this characteristic introduces new kinds of parallelism and data reuse. Now that each node deals with a long vector (often length of 256 or higher) rather than a scalar, we can exploit intra-node parallelism. Furthermore, all the vectors associated with different nodes share a weight matrix which brings in data reuse opportunities. These differences make most graph-based accelerators (e.g., [2, 5, 13–16, 23, 34, 38, 58, 59, 71, 81, 82]) ill-suited for GNNs.

- **Computation pattern disparity:** Different steps of the GCN algorithm deal with different sparsity rates as will be discussed in Section 2.1. In addition, a GNN may utilize other types of computation patterns, such as neural tensor network computation in SimGNN (see Section 4.1 for details) to make an end-to-end application. Such computation pattern variations call for a customized processing unit for each of these steps.
- **Dynamic workload and parallelism:** Apart from the random memory location of the neighbors of a node, the *number* of neighbors also varies across different nodes. This will result in load-imbalance between the graph’s nodes that can degrade the performance significantly.

In addition to the challenges mentioned above, dealing with small graphs (which is particularly challenging for GPUs as shown in Section 5.4.2) requires special design considerations as we will explain in Section 3. To solve these challenges, in this paper, we present SPA-GCN as an efficient and flexible **GCN** accelerator for small graphs that exploits all the available *sparsity*. Then, we apply it to accelerate the entire processing pipeline of SimGNN as an end-to-end application. Since we are facing a *memory-bounded* application, our goal is to have the least number of global memory transactions and to read each element only once. Because of the *computation pattern disparity*, we analyze the requirements for each step of the computation pipeline and, accordingly, develop a dedicated architecture for each of them. To deal with *the irregular memory access*, we utilize a scratchpad memory to store the matrices that need random access. This is doable as we are targeting small graphs. We further propose an efficient workload distribution mechanism to alleviate the *load imbalance* problem.

Concisely, we fuse all the stages together and employ a very deep pipeline with four different levels of nested parallelization and customize the computation units for each of the stages based on its workload as listed in Table 1. These design decisions distinguish SPA-GCN from recently proposed GCN accelerators [21, 69, 74, 76] as summarized in Table 1. Section 6 compares our approach to prior works in more details. While we use SimGNN for illustrating our

approach, the same optimizations can be applied to other GCN-based networks dealing with small graphs such as [4, 26, 43] as well. We implement SPA-GCN on three different FPGAs showing the flexibility and adaptivity of SPA-GCN to different platforms with different global memory bandwidth. Experimental results suggest that SPA-GCN mapped to an HBM FPGA outperforms a multi-core CPU and a GPU implementation by 18.2× and 26.9×, respectively.

In summary, the key contributions of this paper are:

- We design and develop SPA-GCN, a flexible architecture for accelerating GCN specialized for processing many small graphs.
- We adopt SPA-GCN to accelerate SimGNN as an end-to-end application, resulting in an efficient architecture with a very deep pipeline and four levels of parallelization. To our knowledge, this is the first hardware accelerator for GCN-based graph matching.
- We demonstrate the flexibility of our architecture by mapping and customizing it to three different FPGAs with different capacities and memory systems.
- Experimental results suggest that our accelerator can outperform multi-core CPU by 18.2× and GPU by 26.9×, demonstrating the efficiency of our design.

## 2 Background

### 2.1 Graph Convolutional Network (GCN)

Inspired by the success of Convolutional Neural Networks (CNN) on images, GCN [31] was developed to apply a series of convolutional layers on graphs. Layer  $l$  of a GCN takes an undirected graph  $G(V, E, H^l)$  as the input, where  $V$  and  $E$  denote the nodes and edges of the graph, respectively.  $H^l \in \mathbb{R}^{|V| \times f_l}$  is the matrix of the *input node embeddings* for this layer, with each row containing the embedding (also called node features in the literature) of one of the nodes where  $f_l$  indicates the number of features of each node at layer  $l$ . The core computation of a GCN layer to produce the *output node embeddings*,  $H^{l+1} \in \mathbb{R}^{|V| \times f_{l+1}}$  is as follows:

$$H^{l+1} = \sigma_{activation} \left( A' \cdot H^l \cdot W^l \right) \quad (1)$$

where  $\sigma_{activation}(\cdot)$  is an activation function which typically is a ReLU and  $W^l$  is a layer-specific *trainable weight matrix*.  $A'$  is the *normalized adjacency matrix with added self-connections* and can be calculated as follows:

$$\tilde{A} = A + I_N, \quad \tilde{D}_{ii} = \sum_j \tilde{A}_{ij}, \quad A' = \tilde{D}^{-\frac{1}{2}} \cdot \tilde{A} \cdot \tilde{D}^{-\frac{1}{2}} \quad (2)$$

here,  $A$  and  $I_N$  are the *adjacency* and the *identity matrix*, respectively.  $\tilde{D}$  is a *diagonal matrix* where  $\tilde{D}_{ii}$  is the degree of node  $i$  after the self-connection edges are added.

Fig. 1 depicts the computation in Eq. 1 on a simple toy graph assuming that  $A'$  is given, the activation function is ReLU, and each node embedding is a vector of size 4. The first step in the computation gathers the information from the neighbors of each of the nodes which is denoted in Eq. 1 as  $A' \cdot H^l$ . Since  $A'$  is a normalized adjacency matrix, the computation here is a weighted aggregation. After the Aggregation step, the node embeddings are transformed by applying a pre-trained set of weights. Note that the output embeddings may have a different vector size. In the end, the embeddings are passed through a ReLU layer to introduce non-linearity to the model. The time complexity for each GCN layer is  $O(|E|f_{in}f_{out})$ , where  $|E|$  denotes the number of edges including the self-connection ones, since  $A' \cdot H^l$  can be implemented as a sparse-dense matrix multiplication [31]. As we will explain in Section 3.4, because of the ReLU unit at the end of each GCN layer, the node embeddings will also be a sparse matrix. However, the rate of sparsity is different as the adjacency matrix is, in fact, an *ultra* sparse matrix [21].

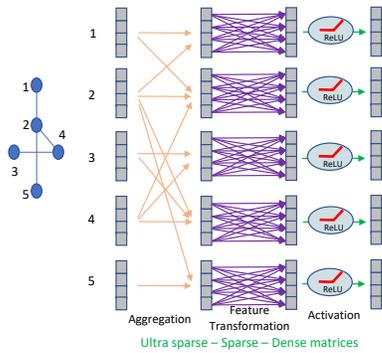


Figure 1: Overview of the core computation in a GCN Layer

### 3 SPA-GCN Architecture

The first step in designing an architecture for GCN is to determine the order of the matrix multiplications. More specifically, we can either compute Eq. 1 as  $(A' \times H^l) \times W^l$  or  $A' \times (H^l \times W^l)$ . We have chosen the latter since it results in a fewer number of operations. Intuitively, this is because both matrices  $A'$  and  $H^l$  are sparse, but their multiplication creates a dense matrix. As a result, in the former, we end up doing a dense-dense multiplication for the second multiplication. However, if we go with the latter, both multiplications are sparse-dense that as shown in AWB-GCN [21], reduces the number of operations. Figure 4 illustrates the high-level view of GCN architecture in SPA-GCN. In this section, we employ a bottom-up approach to highlight the optimization opportunities when GCN is applied to small graphs and how we used them to build the GCN accelerator as demonstrated in Figure 4.

#### 3.1 SPA-GCN Design Principles

As our application is memory-bounded and we are targeting small graphs, SPA-GCN is designed:

- To reduce the number of times we access the global memory to the least amount possible. In our final architecture, each input element is read only once and there is no need to store any of the intermediate results in the global memory.
- To exploit all the available sparsity.

- To employ a deep pipeline with varying levels and degrees of parallelization for matching the workload of different stages and maximizing the overall performance.
- To customize the architecture to efficiently handle small graphs.

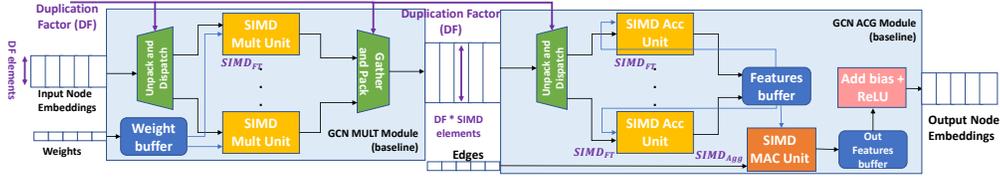
### 3.2 Baseline Architecture

In this section, we describe the basic optimizations that can be applied for processing GCNs. We explain the optimized way of scheduling the operations of a GCN when targeting small graphs and justify the use of an outer-product-based multiplication [8, 55] for it. Furthermore, we demonstrate how we support sparse computation in the Aggregation step. Although these optimizations are necessary, they are not enough when dealing with many small graphs. We cover the rest of the optimizations that can be applied here in the subsequent sections.

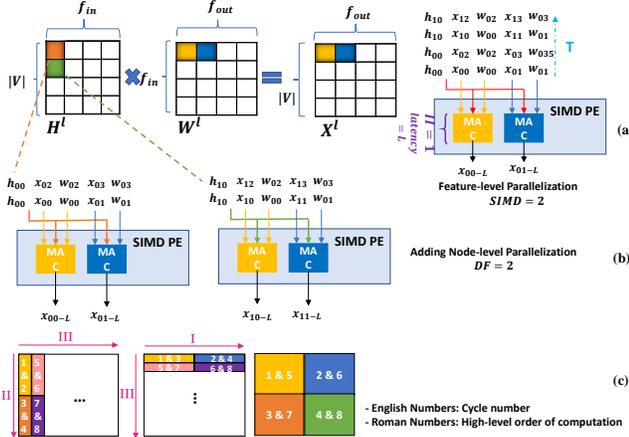
**3.2.1 Feature Transformation:** In this step, one has to multiply matrices  $H^l \in \mathbb{R}^{|V| \times f_{in}}$  and  $W^l \in \mathbb{R}^{f_{in} \times f_{out}}$ , where  $f_{in}$  and  $f_{out}$  denote the number of input and output features, respectively. Since the matrix multiplication has  $|V| \times f_{in} \times f_{out}$  operations, the minimum latency for doing this computation is:  $\frac{|V| \times f_{in} \times f_{out}}{\#MAC_{used}}$  where  $\#MAC_{used}$  is the number of used *multiply and accumulate* (MAC) units. To achieve this latency, we should be able to schedule a new set of operations in each clock cycle. However, adopting an inner-product-based matrix multiplication results in updating the same output feature in the consecutive iterations which introduces read-after-write (RAW) dependency between them. As a result, our pipeline cannot achieve an initiation interval (II) of one which degrades the efficiency of our accelerator.

**Optimized Scheduling:** *Read the Weight Matrix Row-wise; Stream the Embeddings Matrix Column-wise.* We can solve the aforementioned problem by changing the order of computations. To break the dependency, we should update different output locations by taking an element from one of the input matrices (read as a stream) and broadcasting that to parallel MAC units while each MAC unit reads different elements from the other matrix. We choose to read  $H^l$  as a stream and prefetch and cache  $W^l$  since it will be reused by  $H^l$ . With this scheduling we can support sparsity for this step more efficiently as discussed in Section 3.4.

Fig. 3 (a) depicts the modified scheduling for cycle  $l$ . As the figure suggests, we divide the workload within a PE by parallelizing *SIMD* operations across the output features. To read each element only once, for each fetched element of  $H^l$ , (i.e.,  $h_{00}$ ), we schedule all the operations it is involved with, before its eviction. In addition to reusing the  $h_{00}$ , this scheduling increases the number of cycles before a RAW dependency happens. This is because by reading each element of the input node embedding, we can ensure that we will be updating different output locations in the next  $\frac{f_{out}}{SIMD}$  cycles.  $x_{ij-L}$  elements in the figure represent the result of the MAC operations that were scheduled  $L$  cycles ( $L$  being the latency of a MAC unit) before the current inputs:  $x_{ij}$  elements. As long as  $\forall l \in [1 .. L-1] x_{ij-l} \neq x_{ij}$ , there is no dependency; hence, we can schedule a new operation every cycle and achieve  $II=1$ . To make sure that  $II=1$  is possible, we read matrix  $H^l$  in column order. Note that if we read it rather in row order, we update the same location every  $\frac{f_{out}}{SIMD}$  iterations instead of every  $|V| \times \frac{f_{out}}{SIMD}$  iterations.



**Figure 2: SPA-GCN baseline architecture: intra-layer pipelining between MULT module (multiplication unit for Feature Transformation step) and ACG Module (accumulation unit for Feature Transformation step + Aggregation step)**



**Figure 3: Scheduling order for Feature Transformation: (a) Feature-level parallelization. (b) Node-level parallelization. (c) The overall computation order.**

We add a second level of parallelization by duplicating the SIMD PE by a *duplication factor* ( $DF$ ) which parallelizes the node dimension. As a result, we can make use of memory coalescing by packing and reading  $DF$  elements in each cycle. Fig. 3(c) illustrates the final execution order of this step. The arrows denote the high-level ordering of traversing different dimensions and the numbers show the elements that are accessed at their respective cycles. It is important to traverse the input feature dimension ( $f_{in}$ ) last (arrow III) since it is the dimension causing the dependencies. In the baseline architecture, to ensure that the dependency distance is larger than  $L$  (to get  $\Pi=1$ ), we pad the matrix  $H^l$  with zeros by adding rows until  $\frac{(|V|+padding)}{DF} \times \frac{f_{out}}{SIMD} \geq L$ . In Section 3.4, however, we insert bubbles in case of a dependency since the location of non-zero elements is dynamic.

**3.2.2 Aggregation:** In this step, we must multiply matrices  $A' \in \mathbb{R}^{|V| \times |V|}$  and  $X^l \in \mathbb{R}^{|V| \times f_{out}}$  where  $X^l$  is the result of the Feature Transformation step and  $A'$  denotes the normalized adjacency matrix with added self-connections. Due to the highly irregular access to the matrix  $X^l$  to aggregate features of the neighbors, we allocate a scratchpad memory to it. Matrix  $A'$ , is often ultra sparse [21]. To reduce the number of both transferred elements and operations, we prune this matrix and only pass its non-zero elements, which represent edges, to FPGA. For most graph databases, like the ones we target, this step is not needed since the graphs are already stored as a data structure that contains a list of the vertices and edges. Instead of dedicating an on-chip memory for storing the edges, we

read them as a stream and update all the features of the destination node, before retiring the edge.

We further pre-process the adjacency matrix to not only pre-compute  $A'$ , but also re-organize the edges to prevent having RAW dependencies. More specifically, before sending the edges to the FPGA, they are re-arranged offline so that the ones with the same destination node are at least  $L$  locations apart to make sure there is not more than one update to the same node within the window of  $L$  cycles. Without this re-ordering, we either have to increase the  $\Pi$  for the whole operation or add a control unit to insert bubbles in the pipeline of this computation engine in case there is a RAW dependency to ensure correctness. We chose to pre-process the edges since as we are dealing with small graphs, the time for this pre-processing is negligible. We only make use of feature-level parallelism to distribute the workload in this step. Edge-level parallelism can result in bank conflicts since they update random nodes.

**3.2.3 Intra-layer Pipelining:** In the last two sections, we explained the basic optimizations that should be applied to the major computation steps of GCN, Feature Transformation and Aggregation. The last step (ReLU) has a very low area and latency overhead. We only utilize a  $max(0, \cdot)$  unit at the end of the aggregation module for it. To further boost the performance, we add intra-layer pipelining by exploiting a dataflow architecture for connecting the modules (see Fig. 2). As a result, the overall latency will be close to the latency of the slowest module (here, the ACG module). In addition, we can avoid off-chip memory accesses in between these modules.

The computation in the Feature Transformation step can be divided into two separate tasks. Task 1, the MULT unit, is responsible for doing all the multiplications and Task 2, the ACC unit, does the additions. We pipeline these tasks by implementing them as separate modules that are connected by FIFOs. To avoid deadlocks, the throughput (i.e.  $\Pi$ ) of these units should match. The MULT module, as depicted in Fig. 2, has a local buffer to store the weights and streams the elements of  $H^l$  from input FIFO. Each entry of this FIFO is a concatenation of  $DF$  elements (see Section 3.2.1). Both the SIMD factor and  $DF$  can be customized based on the target network configuration. Once the multiplication results are ready, they are packed and sent to the ACG module.

The ACC unit of the Feature Transformation step and the Aggregation step share the matrix  $X^l$ ; thus, to save memory resources, we merge their computation into one module as shown in ACG module in Fig. 2. After fetching the output of the MULT module, the ACG module unpacks the data based on the same  $DF$ , and dispatches SIMD elements to each *SIMD Acc Unit* with the same SIMD factor. Once the additions are done, it will store the partial results to the local buffer *features buffer*. After all updates are committed to the

features buffer, the matrix  $X^l$  is computed and the Aggregation step can start. The SIMD factor of this step is higher than the one in Feature Transformation step since we only exploit feature-level parallelization here. After this step is finished, the elements of the *out features buffer* are added with a bias, passed through a ReLU unit, and stored into off-chip memory. Note that in the baseline architecture, we reuse the same modules for all the GCN layers.

### 3.3 Adding Inter-layer Pipelining

As it is commonly practiced ([35, 69, 74, 76]), in the baseline architecture, we only exploit intra-layer pipelining and reuse the modules for all the GCN layers. However, this is not sufficient when we are dealing with small graphs. The off-chip communication is a serious burden for this application since it deals with small-sized inputs. To alleviate this problem, we intend to reduce the number of accesses to the off-chip memory as much as possible. The baseline architecture is inefficient with this regard since, at the end of each layer, the output should be stored to the off-chip memory and read back again for the next layer. To avoid these redundant accesses, we propose to extend the dataflow architecture described in Section 3.2.3 to all layers of GCN. To realize this, we instantiate new modules for each layer and connect them with FIFOs as depicted in Fig. 4.

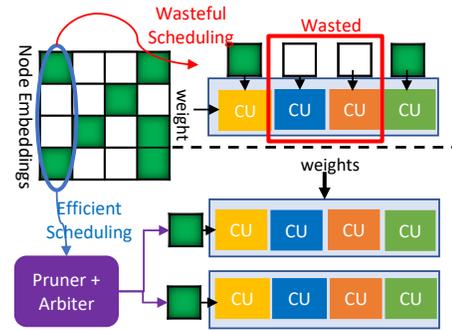
The various (LUT, BRAM, and URAM) and large on-chip memory resources of FPGA and the rather small matrices used for our target application enable us to dedicate separate modules for each of the layers. Fusing the computation for all the layers by enabling dataflow architecture has several benefits such as: 1) we can avoid writing the intermediate results to the global memory by forwarding them to the next layer through FIFOs. 2) The operations will be dynamically scheduled since each module can perform its operation whenever it has a data available. 3) Since we are instantiating different modules for each layer, we can customize the parallelization factors of each module based on the available workload of their respective GCN layer. 4) As the adjacency matrix of a graph does not change across different layers, we can read the edges from the global memory only once for the first layer and reuse them for the subsequent ones by transferring them through the on-chip FIFOs.

### 3.4 Extending the Support for Sparse Computation

The input node embedding to the first layer of GCN usually contains many zero elements since it often uses one-hot encoding for assigning initial vectors to nodes. Furthermore, as shown in Fig. 1, at the end of each GCN layer, there is a ReLU unit. As a result, the matrix generated by each layer, which is the input to the next layer, is sparse. In fact, we saw 52% and 47% sparsity on average for the input to the second and the third layers of GCN in SimGNN for randomly drawn graphs from our target dataset. As a result, the Feature Transformation step also needs to have the support for sparse computation. To reduce the number of operations, we prune the zero elements and only pass the non-zero ones to the next layer. As the updates to the output buffer may come in random cycles, it is necessary to store the buffer containing the partial results on-chip to enable random access. For the same reason, we pack the node features with their address which includes their row and column ID. Packing the elements with their address helps to make the dispatch

unit simpler since each SIMD PE is free to work with any data and knows which partial result should be updated; hence, there is no need to take special considerations to navigate the data to the correct PE. We only need to make sure that at all the times each SIMD PE is working with a different memory bank. We employ an arbiter for this matter, as explained below.

As mentioned in Section 3.2, to reduce the number of RAW dependencies, we chose to stream the node embedding matrix and broadcast it to different computation units (CU) which read the weight matrix as a batch. Since the node embedding is a sparse matrix, reading it as a stream facilitates the pruning mechanism we employ and enables us to distribute the workload more efficiently. Fig. 5 demonstrates a toy example illustrating the benefit of this scheduling. The colored squares show the non-zero elements of the node embedding matrix. By mapping the weights, which are non-zero, to the SIMD dimension, all the CUs in the PE would execute useful operations and we can skip all the operations involving a zero node embedding.



**Figure 5: The benefit of streaming the node embeddings and mapping the weights to the SIMD dimension for supporting sparse computation.**

When skipping the zero node embeddings, the dependency distance for output elements may change dynamically since the number of non-zero inputs between the updates to the same location can be different. Even though the scheduling discussed in Section 3.2 increases the dependency distance as much as possible by doing all the MAC operations when a non-zero input is encountered (each non-zero element would fill  $\frac{f_{out}}{SIMD_{FT}}$  cycles of the dependency window), there still may be some cases where the dependency distance is less than  $L$  (the latency of our function unit (FU) causing the dependency) after this optimization. Instead of setting the  $\Pi$  to  $L$  to ensure correctness, we first insert  $L$  registers to store the partial results of FU at the end of each of its pipeline stages; hence, we can schedule a new set of operations at each clock cycle ( $\Pi=1$ ). There may be cases where the new scheduled operations want to update a location whose old value is still in the registers and have not updated the buffer. To ensure the correctness, we add a control unit which keeps track of the last cycle that each of the output locations was updated. If the number of cycles between two updates to the same location is less than  $L$ , the control unit will insert bubbles into the pipeline until the previous update is committed.

We insert a unit for pruning zeros at the end of the Aggregation step in *ACG* module. As Fig. 6 demonstrates, at each cycle, we

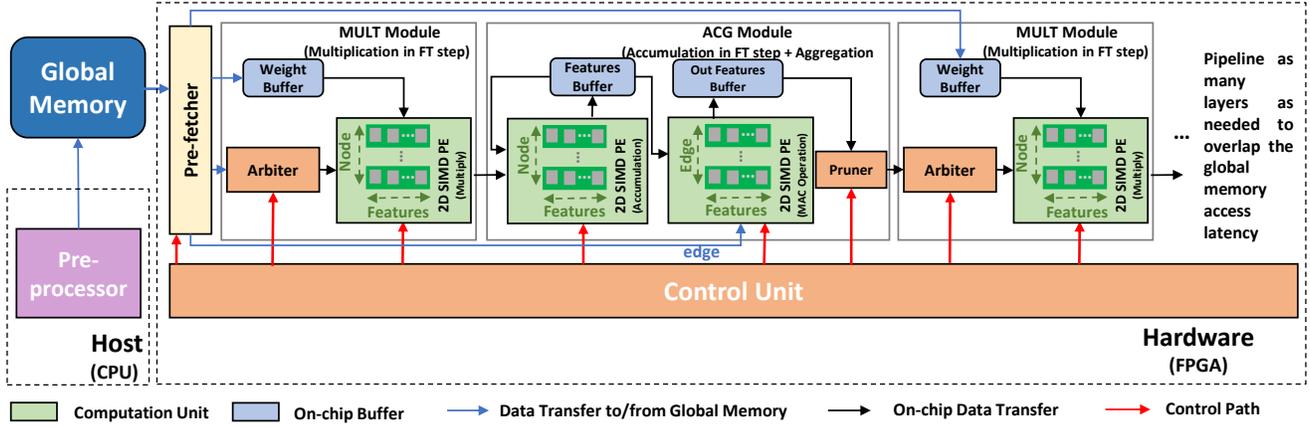


Figure 4: High-level overview of GCN accelerator architecture in SPA-GCN

evaluate  $P$  elements of the node embeddings and pass each to a FIFO if it is not zero. To prevent deadlocks in the system, we use non-blocking read or writes for these FIFOs and only write (read) elements when the FIFO is not full (empty). The *MULT* module of the Feature Transformation step of the next layer takes the  $P$  FIFOs as the input and uses an arbiter to fetch, at most,  $DF$  of them ( $DF \leq P$ ) for passing to  $DF$  SIMD PEs. An arbiter keeps track of the FIFO whose turn it is to be read first in the next cycle. It then uses a round-robin ordering for dispatching the elements from the non-empty FIFOs. After dispatching the inputs, it checks for the RAW dependency by scanning the *prev iter* buffer which contains the last cycle when each element was seen as the input. If the distance was less than  $L$ , it will insert bubbles in the pipeline until the previous input has committed its update to avoid the conflict on IL. If there is no dependency, for each memory bank no more than one element from the dispatched inputs will be issued to a SIMD PE and the current cycle number will be stored in *prev iter* buffer for that input.

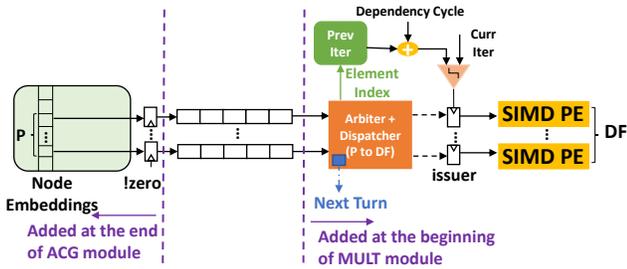


Figure 6: Architecture support for sparse computation in Feature Transformation step.

Table 2: Summary of architecture parameters for the accelerator of each GCN layer in SPA-GCN.

Design Parameter	Explanation
$SIMD_{FT}$	SIMD factor of the FT step
$SIMD_{Agg}$	SIMD factor of the Aggregation step
$DF$	Duplication factor of the PEs in FT step
$P$	Number of input FIFOs to the arbiter

FT: Feature Transformation

To recap, the SPA-GCN architecture provides a flexibility in choosing the parallelization rates. Table 2 lists the parameters that can be tuned for each GCN layer based on its workload. The SIMD factors correspond to feature-level parallelization, while  $DF$  and  $P$  map to node dimension.

## 4 SPA-GCN Application to Graph Matching

In Section 3, we proposed an architecture for GCN specialized for small graphs. In this section, we extend our architecture to accelerate an end-to-end application, SimGNN. An end-to-end application introduces new computation patterns beyond GCN which requires customized processing units. We introduce a new level of parallelization in Section 5.4.3.

### 4.1 SimGNN

Bai et al. [3] proposed a neural-network-based approach to assign a similarity score to two graphs. Its computation pipeline consists of four major stages as depicted in Fig. 7. The first stage is made up of three layers of GCN to extract the *node embeddings*  $H \in \mathbb{R}^{|V| \times F}$  where  $F$  is the number of features of the last layer. In the second stage, it uses a *Global Context-Aware Attention layer (Att)* to combine the node embeddings and generate a single embedding per graph  $h_G \in \mathbb{R}^F$ . For this matter, at first, the Att stage computes the global context by taking an average of the node embeddings followed by a non-linear transformation:  $c = \tanh(\frac{1}{|V|} W_{Att} \sum_{n=1}^{|V|} h_n)$  where  $W_{Att} \in \mathbb{R}^{F \times F}$  is a learnable *weight matrix*. Then, it calculates the importance of each node by constructing an *attention weight*  $a_n$  for them to denote their similarity score to the global context by computing:  $a_n = \frac{1}{1 + \exp(h_n^T \cdot c)}$ . Finally, the graph embedding can be calculated by taking a weighted sum of the node embeddings using the attention weights. The computation in this stage can be summarized as follows:

$$h_G = \sum_{n=1}^{|V|} \sigma(h_n^T \cdot \tanh(\frac{1}{|V|} W_{Att} \sum_{n=1}^{|V|} h_n)) \cdot h_n \quad (3)$$

where  $\sigma(\cdot)$  denotes the sigmoid function to produce  $a_n$ . The time complexity of this stage can be seen to be  $O(|V|F)$ .

The third stage is a *Neural Tensor Network (NTN)* that calculates a vector of similarity score between the two graphs:

$$s(h_{G_1}, h_{G_2}) = \sigma(h_{G_1}^T W_{NTN}^{[1:K]} h_{G_2} + V \cdot \text{concat}(h_{G_1}, h_{G_2}) + b) \quad (4)$$

where  $W_{NTN}^{[1:K]} \in \mathbb{R}^{F \times F \times K}$ ,  $V \in \mathbb{R}^{K \times 2F}$ , and  $b \in \mathbb{R}^K$  are learnable *weight tensor*, *weight matrix*, and *bias vector*, respectively.  $K$  is a hyper-parameter that controls the number of similarity scores. The time complexity of this stage is  $O(F^2K)$  where  $F$  denotes the dimension of the graph level embedding. In the last stage, a standard fully connected network is used to gradually reduce the similarity vector to only one score.

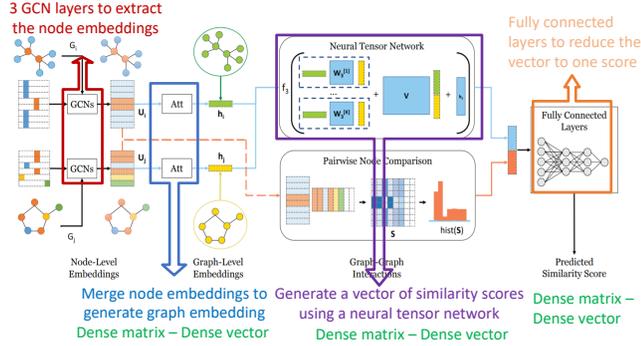


Figure 7: The SimGNN pipeline as depicted in [3].

Since the core computation stage of SimGNN consists of GCN layers, the challenges mentioned in previous sections still exist here. Furthermore, the rest of the stages make use of exponential and hyperbolic tangent functions which are expensive to have on FPGA that can limit the rate of parallelism for them. By comparing the computation complexity of these stages, with GCN (Section 2.1), we can see that the GCN step is the most computation-intensive step; hence, when pipelining all the stages together, the accelerator will be bottlenecked by the GCN step. Therefore, we do not aggressively parallelize the rest of the steps. We rather focus on reducing their resource utilization.

## 4.2 Att Architecture

The SimGNN pipeline applies the GCN stage to two graphs for each comparison query. Instead of duplicating the architecture in Fig. 4 twice, we process the graphs serially and reuse the GCN module for the two input graphs in the query. Reusing the GCN module enables us to map the design to smaller FPGAs as well. Note that the GCN stage is the most resource-hungry one of them all. We improve the performance for processing one query by overlapping the GCN computation of one graph with the *Att* computation of the other one. Thus, the total performance will be bottlenecked with the performance of GCN. As a result, we focus on reducing the area and reusing the resources for *Att*. Let  $H \in \mathbb{R}^{F \times |V|}$  be the transposed result of the GCN stage of SimGNN, then, its  $n$ -th column ( $x_n$ ) will be the vector  $h_n$  in Eq. 3. In computing  $v = W_{Att} \sum_{n=1}^{|V|} h_n$ , we first have to add  $h_n$  vectors and then do a matrix-vector multiplication (MVM). Instead of instantiating separate adders for the first additions and the ones in MVM, we rewrite the equation as follows to

reuse the adders:

$$v = W_{Att} \sum_{n=1}^{|V|} h_n = \sum_{n=1}^{|V|} W_{Att} h_n = \text{sum}(W_{Att} \cdot H, 2) \quad (5)$$

where  $\text{sum}(W_{Att} \cdot H, 2)$  denotes the reduction of the result matrix across its second dimension (columns), meaning that all the multiplications associated with a column of  $H$  should be added together.

Fig. 8 demonstrates an overview of the *Att* module. As in the GCN stage, we divide the matrix multiplication to two different modules, one responsible for multiplications and the other for additions. Again, we use SIMD PEs to implement these modules. However, the SIMD factor here can be set to a different value compared to the GCN stage since they have different computation complexities. The *Repack* module is responsible for adjusting the output of GCN with the SIMD factor of this stage. We use the implementation of *tanh* and *exp* functions available in *Xilinx HLS Math* library since they are already optimized. Note that the last summation in Eq. 3 can be seen as  $H \times a$  where  $a \in \mathbb{R}^{|V|}$  contains the sigmoid results. Hence, we use a matrix vector multiply (MVM) unit at the end.

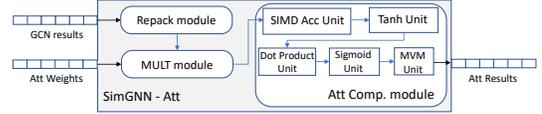


Figure 8: Architecture overview of the second stage of SimGNN in SPA-GCN: Att

## 4.3 NTN + Fully-connected Network Architecture

The computation in the NTN stage is rather simple since it is a series of fixed-size MVMs followed by a bias addition and an activation function. Furthermore, the layers of the fully-connected network (FCN) in the last stage either need an MVM unit or a reduction tree to lower a vector to a scalar. Like the previous stages, we implement all the sub-modules of these two stages in a dataflow-manner. Fig. 9 depicts the architecture of these two steps.

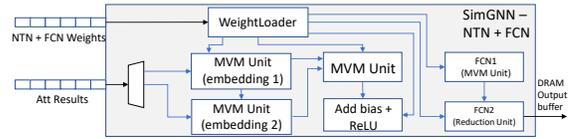


Figure 9: Architecture overview of the last two stages of SimGNN in SPA-GCN: NTN and FCN.

## 4.4 Putting It All Together

Finally, we develop the SimGNN architecture by connecting the modules described in the previous sections. The whole computation pipeline is implemented as a three-level dataflow architecture. The first two levels resemble an inter-stage pipelining while the last one is for intra-stage pipelining. The first level enables a task-level parallelization by grouping the graph-related steps, the *GCN* (Section 3) and *Att* (Section 4.2) modules, and overlapping them with the rest, *NTN\_FC* module (Section 4.3). The second level of dataflow architecture overlaps the *GCN* stage with the *Att*. Finally, the last level applies dataflow architecture to each of the *GCN*, *Att*, and *NTN\_FC* modules as shown in Fig. 4, 8, and 9, respectively. The

*Pre-fetcher* of the GCN module reads the weights and bias of the rest of the stages as well and distributes them to the right place.

We apply three optimizations for reducing the off-chip communication latency: 1) each input buffer can be mapped to a different DRAM bank or HBM channel to enable parallel access to them, 2) the available global memory bandwidth is fully utilized by applying memory coalescing. Memory burst is also applied to amortize the initialization overhead, 3) the modules that access the global memory are overlapped by the computation modules by implementing the accelerator as a dataflow architecture.

## 5 Experimental Results

In this section, we first review our target dataset in Section 5.1 and present the experimental setup in Section 5.2. We then evaluate the effects of our optimizations to the GCN architecture (Section 3) step-by-step in Section 5.3. The whole pipeline of SPA-GCN is then evaluated in Section 5.4 on three different FPGAs that differ in the amount of resources available and the type of global memory they employ. Finally, we compare the performance of SPA-GCN to CPU and GPU implementation and demonstrate the superiority of our design. Related work on general GCN acceleration, sparse CNN, and sparse matrix multiplication will be discussed in Section 6.

### 5.1 Benchmark

We consider a publicly available dataset, AIDS [40], for benchmarking our design. AIDS is a real-life graph dataset containing 42,687 antivirus chemical compounds gathered by the Developmental Therapeutics Program at NCI/NIH. The graphs in AIDS have 25.6 nodes and 27.6 edges on average. We randomly select 10,000 pairs to form 10,000 queries. The kernel time and end-to-end (E2E) time reported in this section are the average of all queries.

### 5.2 Experimental Setup

The SPA-GCN architecture is described using Vivado HLS C++ [68]. The design is synthesized and implemented using Xilinx Vitis 2019.2 [67] on three different target platforms: Xilinx Alveo U50, Xilinx Alveo U280, and Xilinx Kintex UltraScale+ KU15P. The first two are equipped with HBM2 [27] and, ideally, can achieve a bandwidth of 316 GB/s (460 GB/s) with a TDP of 75W (225W) [56, 66]; while the last one utilizes DDR4 as the global memory. Table 3 compares the hardware resources of these boards. For comparison to CPU and GPU, the PyTorch-based implementation of SimGNN from [45] is used that is built using the state-of-the-art PyTorch Geometric (PyG) library [20] which is commonly used as a baseline by previous works [21, 35, 69]. For the Aggregation step, PyG exploits sparsity and edge-level parallelism by adapting the PyTorch Scatter library [1]. For the Feature Transformation step, it uses Intel MKL [25] and NVIDIA cuBLAS library [41] for CPU and GPU respectively, making it a reasonable and optimized baseline. The target CPU in our experiments is Intel(R) Xeon(R) CPU E5-2699 v4 running at 2.2 GHz. For testing on GPU, we use an AWS p3.2xlarge instance which has an NVIDIA V100 GPU.

### 5.3 Impact of GCN Architecture Optimizations

*5.3.1 Inter-Layer Pipelining:* Table 4 shows the resource usage and performance of the SPA-GCN architecture when accelerating three

**Table 3: Properties of the FPGAs used in this paper.**

Platform	BRAM (Mb)	LUT (K)	FF (K)	DSP	URAM (Mb)	Max BW (GB/s)
KU15P	34.6	523	1045	1968	36	19.2
U50	47.3	872	1743	5952	180	316
U280	70.9	1304	2607	9024	270	420

GCN layers of SimGNN on the U280 FPGA. The baseline uses the same hardware for all GCN layers. With inter-layer pipeline added, all 3 GCN layers run in parallel as a coarse-grained pipeline. Since each layer utilizes different pieces of hardware, we can customize the design parameters to match the throughput of each layer. As a result, the 3 layers require 2.4× more DSPs compared with the baseline. Although BRAM usage appears to be smaller, more URAM is used with inter-layer pipelining since we distribute the required resources to obtain a better frequency. The total storage usage is increased in order to store all the buffers for different GCN layers and the FIFOs for the intermediate results between them. The GCN kernel time is reduced by 36% with inter-layer pipelining added to the baseline. However, if we look at the latency-area product metric, i.e., Kernel×DSP, we can see that the performance improvement does not catch up with the computation units (DSP) increment, suggesting potential for further optimizations.

*5.3.2 Extending Sparsity to Feature Transformation:* Indeed, we can extend sparsity support to the Feature Transformation step in GCN to improve both the kernel time and the latency-area product metric, as explained in Section 3.4. Although using  $P$  queue helps the arbiter fetch non-zero elements more frequently, it may still not be enough to dispatch data to all the  $DF$  PEs. Furthermore, by increasing the  $DF$ , we may need to insert more bubbles in the pipeline to avoid RAW dependency since it reduces the number of cycles between the updates to the same location. As a result, there is a trade-off in choosing the right  $DF$  for each layer. Since it is a highly workload-dependent decision, we employ profiling results for setting each of the parallelization factors.

The best parallelization factors are summarized in Table 4. When  $DF$  is set to 1, we no longer need to have separate banks in the row dimension of the buffers which can lessen the number of needed memory blocks. This makes it more efficient to use dense memory blocks (BRAM and URAM) as opposed to LUTs for the buffers, thus BRAM and URAM usage are increased. As Table 4 shows, extending sparsity to feature transformation over inter-layer pipeline has further reduced the kernel time by 31%, effectively achieving a 2.27× speedup over the baseline, while decreasing the DSP usage by 4.09×. The results clearly suggest that, since this is a memory-bounded application, throwing more resources to the architecture is not helpful. Instead, the memory access latency should be reduced and the computation units shall be used more efficiently. Since a large number of zero elements and required DSPs are excluded, the latency-area metric Kernel×DSP is greatly improved by 3.88× over the baseline.

### 5.4 End-to-end Acceleration of SimGNN

*5.4.1 Flexibility of Mapping to Different FPGAs:* We implement the whole pipeline of SimGNN on 2 HBM FPGAs and KU15P that uses DDR memory. Fig. 10 compares the resource breakdown of the modules at the top hierarchy of our design when mapped to U280.

**Table 4: Impact of GCN architecture optimizations on U280. The meaning of design parameters is summarized in Table 2. Baseline shows a single set of design parameters because it uses the same hardware for all layers.**

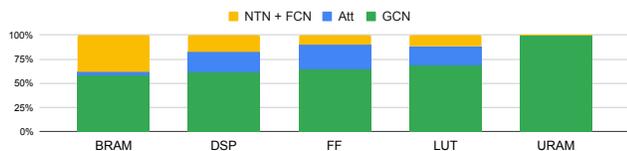
Architecture	Design Parameters (L1 / L2 / L3)				LUT / FF / DSP / BRAM / URAM (%)	Freq. (MHz)	Kernel (ms)	Kernel × DSP
	$SIMD_{FT}$	$SIMD_{Agg}$	DF	P				
Baseline	16	32	8	—	9.8 / 7.7 / 7.4 / 6.8 / 0	265	0.599 (1×)	4.46 (1×)
+Inter-Layer Pipeline	32/16/16	32/32/16	8/8/8	—	14 / 12 / 18 / 3.6 / 2.5	271	0.383 (1.56×)	6.74 (0.66×)
+Extended Sparsity	32/32/16	32/32/16	2/1/1	8/2/2	4.8 / 6.0 / 4.4 / 4.8 / 3.1	300	0.264 (2.27×)	1.15 (3.88×)

As the figure exhibits, we allocate most of the resources to the GCN stage as it is the computation-intensive part of the network. Table 5 shows the resource usage and performance for the three FPGA platforms. We can see that the kernel runs faster on HBM FPGAs compared to KU15P. This is mainly due to the fact that HBM FPGAs can achieve a better frequency as they have more resources and the Vitis tool has more freedom in PnR to optimize the timing. The fact that the multiplication and addition units have different latencies on these boards (5 and 8 cycles for KU15P; 4 and 7 cycles for U280, respectively) further increases the difference in the runtimes. In fact, the cycle count of the same kernel when it uses different types and number of banks for global memory is almost the same. This suggests that after our optimizations the bottleneck is no longer at the memory level. Besides, we can see that the end-to-end time is much larger than the kernel time (runtimes measured on the host using CPU clock). This is mainly because of the overhead of launching the kernel and PCIe transfers since all the pre-processing steps take less than 5% of the end-to-end time. This suggests the potential for batching queries to further improve the end-to-end throughput (Section 5.4.3).

**Table 5: Performance and resource utilization of a SPA-GCN design on different target FPGAs.**

FPGA	LUT / FF / DSP / BRAM / URAM (%)	Freq. (MHz)	Kernel (ms)	E2E (ms)	E2E (query/s)
KU15P	34 / 29 / 35 / 30 / 23	201	0.786	1.135	881
U50	17 / 16 / 12 / 16 / 4.7	279	0.423	0.538	1858
U280	11 / 10 / 7.7 / 10 / 3.1	290	0.327	0.509	1965

E2E: End-to-End



**Figure 10: Resource breakdown of the whole pipeline of SimGNN on U280.**

**Table 6: Performance comparison of running SimGNN on different hardware. The CPU and GPU runtimes are based on state-of-the-art PyG [20] implementation. Low computation complexity of small graphs and the coarse-grained execution of GPUs impact the efficiency of GPU significantly.**

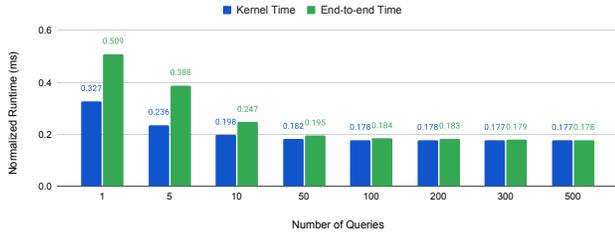
Platform	Max BW (GB/s)	Kernel (ms)	E2E (ms)	Speedup (Over CPU)	Speedup (Over GPU)
KU15P	19.2	0.786	1.135	8.2	12.1
U50	316	0.423	0.538	17.2	25.5
U280	460	<b>0.327</b>	<b>0.509</b>	<b>18.2</b>	<b>26.9</b>
PyG-CPU	76.8	5.85	9.27	1	1.5
PyG-GPU (V100)	900	9.68	13.7	0.68	1

BW: Bandwidth, E2E: End-to-End

**5.4.2 SPA-GCN vs CPU and GPU:** We test the performance of the whole pipeline of SimGNN on the CPU and GPU described in Section 5.2. In this section, we are assuming that the inputs are already stored in the host memory, and we want to offload the graph comparison queries to either of the target platforms. The goal is to compare the performance of these platforms for processing a graph matching query. Table 6 summarizes the results. The results are averaged over 10,000 queries on each of the platforms. The queries are started sequentially, and the end-to-end time of all the platforms is the time interval between two consecutive queries are started. This contains the runtime for any pre-processing steps as well. For FPGA and GPU, it also includes the host-kernel communication via the PCIe link, writing data to FPGA/GPU’s global memory, kernel computation, reading the results from that, and the overheads for using the APIs (OpenCL for FPGA and PyTorch for CPU/GPU). We use the end-to-end time for comparison since these overheads are inevitable and should be accounted for. The kernel time on CPU/GPU is measured with the PyTorch profiler.

The results demonstrate that our FPGA solution can outperform both CPU and GPU significantly. As discussed in Section 1, this is partly because of the dynamic load balance and the irregular memory access of the graph structure. Furthermore, since we target small graphs, it results in extreme under-utilization of GPU. In fact, the profiling results indicate that the GPU utilization does not go higher than 6% and, for the most part, the PyG-GPU only uses 1 streaming multiprocessor (SM) since the matrices are small. Because of this and the fact that GPU runs at a lower frequency (1.3GHz) compared to CPU (2.2 GHz), the GPU version of this application is even slower than the CPU. The *nvprof* profiling results show that PyG-GPU runs 225 kernels for accelerating this application that on average have 4.6 *KFLOPs*. With this low computation intensity, the overhead of running the kernel (such as *cudaLaunchKernel*) is larger than the actual kernel runtime that greatly impacts the GPU performance. Designing the GPU kernel manually can alleviate some of these shortcomings, but the underlying problem still exists due to the coarse-grained execution model of GPU. In contrast, our FPGA solution suffers from the kernel initialization overheads only once since we develop a deep pipeline across all stages of the computation by fusing them in one kernel. This pipelining has several other benefits as explained in Section 3.3.

**5.4.3 Batching Queries on FPGA:** The difference in the end-to-end time and the kernel time on FPGA is mainly due to both the DMA transfers and overhead of OpenCL APIs. In fact, our experiments using the Profile Summary provided by Vitis Runtime library [65] show that OpenCL APIs can take 10-100  $\mu$ s which is comparable to the kernel execution for one query. Therefore, SPA-GCN implements batching graph queries to amortize these overheads. We



**Figure 11: Effect of batching queries. Runtime numbers are the average over all queries.**

can send the data for processing multiple queries to FPGA, process each query serially, and send the results back together. Fig. 11 demonstrates the results of this experiment on U280 board. The results suggest that by transferring the data for ~300 queries, we can amortize the setup time and achieve a 2.8 $\times$  speedup.

Batching queries makes it possible to further improve the throughput by adding another level of parallelization. When targeting the HBM FPGAs, we use 4 of the HBM PCs for processing one query of graph matching. As each HBM FPGA has 32 PCs that can be accessed independently, we can further improve the performance by replicating the logic of SPA-GCN for one query to process up to 8 different queries in parallel as long as we have enough resources available. Table 3 illustrates that the available resources allow us to instantiate 6 SPA-GCN pipelines with U280 before hitting the 80% resource usage upper-bound. This does not change the latency of each graph query, but it would further increase the throughput by 6 $\times$  resulting in a throughput of 33522 query/s on U280.

## 6 Related Work

**GCN Accelerators:** Because of the popularity of GCN, there is a growing interest in developing an accelerator for it [21, 35, 69, 74, 76]. As summarized in Table 1, HyGCN [69], GraphACT [74], and Prasanna et al. [76] develop a fixed hardware for all the layers of GCN and process them sequentially. This is an undesirable feature particularly when we target small graphs. In fact, our baseline architecture, in which we reuse the same architecture for all the GCN layers, exploit only the sparsity of the Aggregation step, treat the Feature Transformation step as a regular matrix multiplication, and employ a 2D computation unit for it, has the same design principles as these works. However, as the experimental results in Table 4 show, not only should we execute the GCN layers in a pipelined fashion, but also we should exploit the sparsity of the node embeddings to enhance both area and performance. In addition, GraphACT and Prasanna et al. [76] mainly rely on *redundancy reduction* in the Aggregation step to decrease the number of operations by pre-computing the repeated aggregations. This is possible only when the adjacency matrix is *binary*. However, not all GCNs can benefit from this feature since they typically work with the *normalized* adjacency matrix meaning that they need *weighted* additions in this step. AWB-GCN [21] proposes an architecture that supports inter-layer pipelining and considerations for sparsity of the node embeddings for accelerating GCN. However, partitioning the computation by the nodes in their approach complicates the design of the task distributor since the node embeddings are sparse and special consideration is needed to prevent PEs from doing

unnecessary operations on the zero elements. On the other hand, feature-level parallelization deals better with workload imbalance as we discussed in Section 3. In addition, AWB-GCN is developed for large graphs and is based on the inner-product matrix multiplication, whereas, as we discuss in Section 3.2, the outer-product multiplication is preferred to reduce the RAW dependency especially when we are targeting small graphs and do not have enough (non-zero) nodes to fill the dependency window.

Moreover, when targeting small graphs, we need more optimizations. For example, it is important to reduce the number of accesses to the global memory as much as possible. This is not the focus of the previous works. Besides, in addition to fitting the design for the whole network into an FPGA, we have the option to process parallel queries of the network. Table 1 summarizes the differences of our design compared to the aforementioned works on GCN acceleration.

**SpMM and SCNN Accelerators:** Apart from the works focusing on GCN, there has been a lot of research on sparse matrix multiplication (MM) either for pruned CNNs or normal MM [17, 24, 32, 42, 52, 53, 60, 78, 83]. They all rely on the fact that the sparse matrix is known offline and they can pre-process it. For example, EIE [24] propose a sparse matrix vector multiplier for the fully-connected layers. It reorganizes the sparse matrix in compressed sparse column (CSC) format and pre-loads that into on-chip memory. As another example, Kung et al. [32] pre-process the data by merging multiple sparse columns of the weight matrix into one and pruning all the weights except for the most-significant ones, resulting in some accuracy loss. These approaches are not feasible for GCN in which the sparse matrix (i.e. the node embeddings) is generated *while running* the algorithm; whereas, we proposed a technique to prune the zeros *on-the-fly*.

## 7 Conclusion

In this paper, we analyzed and examined the optimization opportunities when GCN is applied to small graphs. We presented an efficient architecture, SPA-GCN, and developed an accelerator for SimGNN based on that as an end-to-end application. SimGNN is a neural-network-based graph matching algorithm for calculating an approximation of the GED between two graphs. The computation disparity existing in the network calls for a customized accelerator. Besides, since the GPU has coarse-grained execution, we cannot have improvement beyond the optimizations applied for each phase since different phases are executed separately. However, on the FPGA side, we can exploit a deep pipeline across the phases by enabling a dataflow architecture. Not only does it help us reduce the global memory transactions, we can also eliminate the overhead of running different kernels. Furthermore, we showed that since this is a memory-bounded application, instantiating many computation units (as in GPU) is not beneficial. The experimental results demonstrate that SPA-GCN can outperform CPU and GPU by 18.2 $\times$  and 26.9 $\times$ , respectively, because of the use of a very deep pipeline with different levels of parallelization.

## Acknowledgments

We would like to thank Marci Baun for editing the paper. A significant part of this work was conducted while the first author was

interning at Samsung Semiconductor Inc. It is also supported by the CAPA award jointly funded by NSF (CCF-1723773) and Intel (3688881), the RTML award funded by NSF (CCF-1937599), and CDSC industrial partners<sup>1</sup>.

## References

- [1] 2020. PyTorch Scatter. (2020). [https://github.com/rusty1s/pytorch\\_scatter](https://github.com/rusty1s/pytorch_scatter)
- [2] Johnathan Alsop, Marc S Orr, Bradford M Beckmann, and David A Wood. 2016. Lazy release consistency for GPUs. In *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 1–14.
- [3] Yunsheng Bai, Hao Ding, Song Bian, Ting Chen, Yizhou Sun, and Wei Wang. 2019. Simgnn: A neural network approach to fast graph similarity computation. In *WSDM*. 384–392.
- [4] Yunsheng Bai, Hao Ding, Ken Gu, Yizhou Sun, and Wei Wang. 2020. Learning-based efficient graph similarity computation via multi-scale convolutional set matching. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 3219–3226.
- [5] Abanti Basak, Shuangchen Li, Xing Hu, Sang Min Oh, Xinfeng Xie, Li Zhao, XiaoWei Jiang, and Yuan Xie. 2019. Analysis and optimization of the memory hierarchy for graph processing workloads. In *HPCA*.
- [6] Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, Talapady N Bhat, Helge Weissig, Ilya N Shindyalov, and Philip E Bourne. 2000. The protein data bank. *Nucleic acids research* 28, 1 (2000), 235–242.
- [7] Evan E Bolton, Yanli Wang, Paul A Thiessen, and Stephen H Bryant. 2008. PubChem: integrated platform of small molecules and biological activities. In *Annual reports in computational chemistry*, Vol. 4. Elsevier, 217–241.
- [8] Aydin Buluc and John R Gilbert. 2008. On the representation and multiplication of hypersparse matrices. In *IPDPS*. IEEE, 1–11.
- [9] Horst Bunke and Kim Shearer. 1998. A graph distance metric based on the maximal common subgraph. *Pattern Recognition Letters* 19, 3-4 (1998), 255–259.
- [10] Guangyong Chen, Pengfei Chen, Chang-Yu Hsieh, Chee-Kong Lee, Benben Liao, Renjie Liao, Weiweng Liu, Jiezhong Qiu, Qiming Sun, Jie Tang, et al. 2019. Alchemy: A quantum chemistry dataset for benchmarking ai models. *arXiv preprint arXiv:1906.09427* (2019).
- [11] Xiaoyang Chen, Hongwei Huo, Jun Huan, and Jeffrey Scott Vitter. 2019. An efficient algorithm for graph edit distance computation. *Knowledge-Based Systems* 163 (2019), 762–775.
- [12] Yu-Hsin Chen, Joel Emer, and Vivienne Sze. 2016. Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks. *ACM SIGARCH Computer Architecture News* 44, 3 (2016), 367–379.
- [13] Yuze Chi, Guohao Dai, Yu Wang, Guangyu Sun, Guoliang Li, and Huazhong Yang. 2016. Nxgraph: An efficient graph processing system on a single machine. In *ICDE*. IEEE, 409–420.
- [14] Vidushi Dadu, Jian Weng, Sihao Liu, and Tony Nowatzki. 2019. Towards General Purpose Acceleration by Exploiting Common Data-Dependence Forms. In *MICRO*. <https://doi.org/10.1145/3352460.3358276>
- [15] Guohao Dai, Yuze Chi, Yu Wang, and Huazhong Yang. 2016. Fpgp: Graph processing framework on fpga a case study of breadth-first search. In *FPGA*. 105–110.
- [16] Guohao Dai, Tianhao Huang, Yuze Chi, Ningyi Xu, Yu Wang, and Huazhong Yang. 2017. Foregraph: Exploring large-scale graph processing on multi-fpga architecture. In *FPGA*. 217–226.
- [17] Caiwen Ding, Siyu Liao, YanZhi Wang, Zhe Li, Ning Liu, Youwei Zhuo, Chao Wang, Xuehai Qian, Yu Bai, Geng Yuan, et al. 2017. Circnn: accelerating and compressing deep neural networks using block-circulant weight matrices. In *MICRO*. 395–408.
- [18] Philippe Dosch and Ernest Valveny. 2005. Report on the second symbol recognition contest. In *International Workshop on Graphics Recognition*. Springer, 381–397.
- [19] David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P Adams. 2015. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in neural information processing systems*. 2224–2232.
- [20] Matthias Fey and Jan Eric Lenssen. 2019. Fast graph representation learning with PyTorch Geometric. *arXiv preprint arXiv:1903.02428* (2019).
- [21] Tong Geng, Ang Li, Runbin Shi, Chunshu Wu, Tianqi Wang, Yanfei Li, Pouya Haghi, Antonino Tumeo, Shuai Che, Steve Reinhardt, et al. 2020. AWB-GCN: A graph convolutional network accelerator with runtime workload rebalancing. In *MICRO*. IEEE, 922–936.
- [22] Winston Haaswijk, Edo Collins, Benoit Seguin, Mathias Soeken, Frédéric Kaplan, Sabine Süssstrunk, and Giovanni De Micheli. 2018. Deep learning for logic optimization algorithms. In *ISCAS*. IEEE, 1–4.
- [23] T. J. Ham, L. Wu, N. Sundaram, N. Satish, and M. Martonosi. 2016. Graphicionado: A high-performance and energy-efficient accelerator for graph analytics. In *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. 1–13. <https://doi.org/10.1109/MICRO.2016.7783759>
- [24] Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A Horowitz, and William J Dally. 2016. EIE: Efficient inference engine on compressed deep neural network. *ACM SIGARCH Computer Architecture News* 44, 3 (2016), 243–254.
- [25] Intel. 2020. Intel MKL. (2020). <https://software.intel.com/en-us/mkl>
- [26] Sho Ishida, Tomo Miyazaki, Yoshihiro Sugaya, and Shinichiro Omachi. 2021. Graph Neural Networks with Multiple Feature Extraction Paths for Chemical Property Estimation. *Molecules* 26, 11 (2021), 3125.
- [27] JEDEC. 2020. High Bandwidth Memory (HBM) FPGA. (2020). <https://www.jedec.org/standards-documents/docs/jesd235a>
- [28] Patrick Judd, Jorge Albericio, Tayler Hetherington, Tor M Aamodt, and Andreas Moshovos. 2016. Stripes: Bit-serial deep neural network computing. In *MICRO*. IEEE, 1–12.
- [29] Viggo Kann. 1992. On the approximability of the maximum common subgraph problem. In *Annual Symposium on Theoretical Aspects of Computer Science*. Springer, 375–388.
- [30] Jongik Kim, Dong-Hoon Choi, and Chen Li. 2019. Inves: Incremental Partitioning-Based Verification for Graph Similarity Search.. In *EDBT*. 229–240.
- [31] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [32] HT Kung, Bradley McDanel, and Sai Qian Zhang. 2019. Packing sparse convolutional neural networks for efficient systolic array implementations: Column combining under joint optimization. In *ASPLOS*. 821–834.
- [33] Yujia Li, Chenjie Gu, Thomas Dullien, Oriol Vinyals, and Pushmeet Kohli. 2019. Graph matching networks for learning the similarity of graph structured objects. In *International conference on machine learning*. PMLR, 3835–3845.
- [34] Zhaoshi Li, Leibo Liu, Yangdong Deng, Shouyi Yin, Yao Wang, and Shaojun Wei. 2017. Aggressive Pipelining of Irregular Applications on Reconfigurable Hardware. In *ISCA*. 575–586.
- [35] Shengwen Liang, Ying Wang, Cheng Liu, Lei He, LI Huawei, Dawen Xu, and XiaoWei Li. 2020. EnGN: A High-Throughput and Energy-Efficient Accelerator for Large Graph Neural Networks. *IEEE Trans. Comput.* (2020).
- [36] Hehuan Ma, Yatao Bian, Yu Rong, Wenbing Huang, Tingyang Xu, Weiyang Xie, Geyan Ye, and Junzhou Huang. 2020. Multi-view graph neural networks for molecular property prediction. *arXiv preprint arXiv:2005.13607* (2020).
- [37] Huiyu Mo, Leibo Liu, Wenjing Hu, Wenping Zhu, Qiang Li, Ang Li, Shouyi Yin, Jian Chen, XiaoWei Jiang, and Shaojun Wei. 2020. Tfe: Energy-efficient transferred filter-based engine to compress and accelerate convolutional neural networks. In *MICRO*. IEEE, 751–765.
- [38] Anurag Mukkara, Nathan Beckmann, Maleen Abeydeera, Xiaosong Ma, and Daniel Sanchez. 2018. Exploiting Locality in Graph Analytics through Hardware-Accelerated Traversal Scheduling. In *MICRO*. 1–14.
- [39] Hiroki Nakahara, Zhiqiang Que, and Wayne Luk. 2020. High-Throughput Convolutional Neural Network on an FPGA by Customized JPEG Compression. In *FCCM*. IEEE, 1–9.
- [40] NCI/NIH. 2004. AIDS Antiviral Screen Data. (2004). <https://wiki.nci.nih.gov/display/NCIDTPdata/AIDS+Antiviral+Screen+Data>
- [41] NVIDIA. 2020. NVIDIA cuBLAS. (2020). <https://developer.nvidia.com/cublas>
- [42] Angshuman Parashar, Minsoo Rhu, Anurag Mukkara, Antonio Pugliesi, Rangharajan Venkatesan, Bruce Khailany, Joel Emer, Stephen W Keckler, and William J Dally. 2017. Scnn: An accelerator for compressed-sparse convolutional neural networks. *ACM SIGARCH Computer Architecture News* 45, 2 (2017), 27–40.
- [43] Zongyue Qin, Yunsheng Bai, and Yizhou Sun. 2020. GHashing: Semantic Graph Hashing for Approximate Similarity Search in Graph Databases. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2062–2072.
- [44] Kaspar Riesen and Horst Bunke. 2008. IAM graph database repository for graph based pattern recognition and machine learning. In *Joint IAPR International Workshops on Statistical Techniques in Pattern Recognition (SPR) and Structural and Syntactic Pattern Recognition (SSPR)*. Springer, 287–297.
- [45] Benedek Rozemberczki. 2018. SimGNN. (2018). <https://github.com/benedekrozmberczki/SimGNN>
- [46] Alberto Sanfeliu and King-Sun Fu. 1983. A distance measure between attributed relational graphs for pattern recognition. *IEEE Transactions on Systems, Man, and Cybernetics* 3 (1983), 353–362.
- [47] Yakun Sophia Shao, Jason Clemons, Rangharajan Venkatesan, Brian Zimmer, Matthew Fojtik, Nan Jiang, Ben Keller, Alicia Klinefelter, Nathaniel Pinckney, Priyanka Raina, et al. 2019. Simba: Scaling deep-learning inference with multi-chip-module-based architecture. In *MICRO*. 14–27.
- [48] Hardik Sharma, Jongse Park, Naveen Suda, Liangzhen Lai, Benson Chau, Vikas Chandra, and Hadi Esmaeilzadeh. 2018. Bit fusion: Bit-level dynamically composable architecture for accelerating deep neural network. In *ISCA*. IEEE, 764–775.
- [49] Yongming Shen, Michael Ferdman, and Peter Milder. 2017. Escher: A CNN accelerator with flexible buffering to minimize off-chip transfer. In *FCCM*. IEEE, 93–100.

<sup>1</sup><https://cdsc.ucla.edu/partners/>

- [50] Atefeh Sohrabizadeh, Jie Wang, and Jason Cong. 2020. End-to-End Optimization of Deep Learning Applications. In *FPGA*. 133–139.
- [51] Linghao Song, Xuehai Qian, Hai Li, and Yiran Chen. 2017. Pipelayer: A pipelined ream-based accelerator for deep learning. In *HPCA*. IEEE, 541–552.
- [52] Nitish Srivastava, Hanchen Jin, Jie Liu, David Albonesi, and Zhiru Zhang. 2020. Matraptor: A sparse-sparse matrix multiplication accelerator based on row-wise product. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, 766–780.
- [53] Nitish Srivastava, Hanchen Jin, Shaden Smith, Hongbo Rong, David Albonesi, and Zhiru Zhang. 2020. Tensaurus: A Versatile Accelerator for Mixed Sparse-Dense Tensor Computations. In *HPCA*.
- [54] Teague Sterling and John J Irwin. 2015. ZINC 15–ligand discovery for everyone. *Journal of chemical information and modeling* 55, 11 (2015), 2324–2337.
- [55] Robert A Van De Geijn and Jerrell Watts. 1997. SUMMA: Scalable universal matrix multiplication algorithm. *Concurrency: Practice and Experience* 9, 4 (1997), 255–274.
- [56] Xilinx. 2020. Alveo U280 Data Center Accelerator Card Data Sheet. (2020). [https://www.xilinx.com/support/documentation/data\\_sheets/ds963-u280.pdf](https://www.xilinx.com/support/documentation/data_sheets/ds963-u280.pdf)
- [57] Xiaoli Wang, Xiaofeng Ding, Anthony KH Tung, Shanshan Ying, and Hai Jin. 2012. An efficient graph indexing method. In *ICDE*. IEEE, 210–221.
- [58] Yangzihao Wang, Andrew Davidson, Yuechao Pan, Yuduo Wu, Andy Riffel, and John D Owens. 2016. Gunrock: A high-performance graph processing library on the GPU. In *Proceedings of the 21st ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*. 1–12.
- [59] Yu Wang, James C Hoe, and Eriko Nurvitadhi. 2019. Processor assisted workload scheduling for FPGA accelerated graph processing on a shared-memory platform. In *FCCM*. IEEE, 136–144.
- [60] Yang Wang, Chen Zhang, Zhiqiang Xie, Cong Guo, Yunxin Liu, and Jingwen Leng. 2021. Dual-side Sparse Tensor Core. In *ISCA*.
- [61] Craig I Watson and Charles L Wilson. 1992. NIST special database 4. *Fingerprint Database, National Institute of Standards and Technology* 17, 77 (1992), 5.
- [62] Xuechao Wei, Yun Liang, Xiuhong Li, Cody Hao Yu, Peng Zhang, and Jason Cong. 2018. TGPA: tile-grained pipeline architecture for low latency CNN inference. In *ICCAD*. 1–8.
- [63] Xuechao Wei, Cody Hao Yu, Peng Zhang, Youxiang Chen, Yuxin Wang, Han Hu, Yun Liang, and Jason Cong. 2017. Automated systolic array architecture synthesis for high throughput CNN inference on FPGAs. In *DAC*. 1–6.
- [64] Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. 2018. MoleculeNet: a benchmark for molecular machine learning. *Chemical science* 9, 2 (2018), 513–530.
- [65] Xilinx. 2019. Vitis Manual - UG1393. (2019). [https://www.xilinx.com/support/documentation/sw\\_manuals/xilinx2019\\_2/ug1393-vitis-application-acceleration.pdf](https://www.xilinx.com/support/documentation/sw_manuals/xilinx2019_2/ug1393-vitis-application-acceleration.pdf)
- [66] Xilinx. 2020. Alveo U50 Data Center Accelerator Card Data Sheet. (2020). [https://www.xilinx.com/support/documentation/data\\_sheets/ds965-u50.pdf](https://www.xilinx.com/support/documentation/data_sheets/ds965-u50.pdf)
- [67] Xilinx. 2020. Vitis Unified Software Platform. (2020). <https://www.xilinx.com/products/design-tools/vitis/vitis-platform.html>
- [68] Xilinx. 2020. Vivado Design Suite User Guide - High-Level Synthesis (UG902). (2020). <https://www.xilinx.com/>
- [69] Mingyu Yan, Lei Deng, Xing Hu, Ling Liang, Yujing Feng, Xiaochun Ye, Zhimin Zhang, Dongrui Fan, and Yuan Xie. 2020. Hygen: A gen accelerator with hybrid architecture. In *HPCA*. IEEE, 15–29.
- [70] Pinar Yanardag and SVN Vishwanathan. 2015. Deep graph kernels. In *SIGKDD*. 1365–1374.
- [71] Yifan Yang, Zhaoshi Li, Yangdong Deng, Zhiwei Liu, Shouyi Yin, Shaojun Wei, and Leibo Liu. 2020. GraphABCD : Scaling Out Graph Analytics with Asynchronous Block Coordinate Descent. In *ISCA*.
- [72] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018. Graph convolutional neural networks for web-scale recommender systems. In *SIGKDD*. 974–983.
- [73] Hanqing Zeng, Ren Chen, Chi Zhang, and Viktor Prasanna. 2018. A framework for generating high throughput CNN implementations on FPGAs. In *FPGA*. 117–126.
- [74] Hanqing Zeng and Viktor Prasanna. 2020. Graphact: Accelerating gen training on cpu-fpga heterogeneous platforms. In *FPGA*. 255–265.
- [75] Zhiping Zeng, Anthony KH Tung, Jianyong Wang, Jianhua Feng, and Lizhu Zhou. 2009. Comparing stars: On approximating graph edit distance. *Proceedings of the VLDB Endowment* 2, 1 (2009), 25–36.
- [76] Bingyi Zhang, Hanqing Zeng, and Viktor Prasanna. 2020. Accelerating large scale GCN inference on FPGA. In *FCCM*. IEEE, 241–241.
- [77] Chen Zhang, Guangyu Sun, Zhenman Fang, Peipei Zhou, Peichen Pan, and Jason Cong. 2018. Caffeine: Toward uniformed representation and acceleration for deep convolutional neural networks. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 38, 11 (2018), 2072–2085.
- [78] Shijin Zhang, Zidong Du, Lei Zhang, Huiying Lan, Shaoli Liu, Ling Li, Qi Guo, Tianshi Chen, and Yunji Chen. 2016. Cambricon-X: An accelerator for sparse neural networks. In *MICRO*. IEEE, 1–12.
- [79] Xiaofan Zhang, Junsong Wang, Chao Zhu, Yonghua Lin, Jinjun Xiong, Wen-mei Hwu, and Deming Chen. 2018. DNNBuilder: an automated tool for building high-performance DNN hardware accelerators for FPGAs. In *ICCAD*. IEEE, 1–8.
- [80] Ling Zhao, Yujiao Song, Chao Zhang, Yu Liu, Pu Wang, Tao Lin, Min Deng, and Haifeng Li. 2019. T-gcn: A temporal graph convolutional network for traffic prediction. *IEEE Transactions on Intelligent Transportation Systems* (2019).
- [81] Shijie Zhou, Charalampos Chelmiss, and Viktor K Prasanna. 2015. Accelerating large-scale single-source shortest path on FPGA. In *IPDPSW*. IEEE, 129–136.
- [82] Shijie Zhou, Charalampos Chelmiss, and Viktor K Prasanna. 2016. High-throughput and energy-efficient graph processing on FPGA. In *FCCM*. IEEE, 103–110.
- [83] Maohua Zhu, Tao Zhang, Zhenyu Gu, and Yuan Xie. 2019. Sparse Tensor Core: Algorithm and Hardware Co-Design for Vector-wise Sparse Neural Networks on Modern GPUs. In *MICRO*. 359–371.