

TC-GNN: Bridging Sparse GNN Computation and Dense Tensor Cores on GPUs

Yuke Wang, Boyuan Feng, Zheng Wang, Guyue Huang, and Yufei Ding, University of California, Santa Barbara

https://www.usenix.org/conference/atc23/presentation/wang-yuke

This paper is included in the Proceedings of the 2023 USENIX Annual Technical Conference.

July 10-12, 2023 • Boston, MA, USA

978-1-939133-35-9

Open access to the Proceedings of the 2023 USENIX Annual Technical Conference is sponsored by









TC-GNN: Bridging Sparse GNN Computation and Dense Tensor Cores on GPUs

Yuke Wang, Boyuan Feng, Zheng Wang, Guyue Huang, and Yufei Ding University of California, Santa Barbara

Abstract

Recently, graph neural networks (GNNs), as the backbone of graph-based machine learning, demonstrate great success in various domains (e.g., e-commerce). However, the performance of GNNs is usually unsatisfactory due to the highly sparse and irregular graph-based operations. To this end, we propose TC-GNN, the first GNN acceleration framework based on GPU Tensor Core Units (TCUs). The core idea is to reconcile the "Sparse" GNN computation with the highperformance "Dense" TCUs. Specifically, we conduct an indepth analysis of the sparse operations in mainstream GNN computing frameworks. We introduce a novel sparse graph translation technique to facilitate TCU processing of the sparse GNN workload. We implement an effective CUDA core and TCU collaboration design to fully utilize GPU resources. We integrate TC-GNN with the PyTorch framework for high programmability. Rigorous experiments show an average of 1.70× speedup over the state-of-the-art DGL framework across various models and datasets.

1 Introduction

Over the recent years, with the increasing popularity of graphbased learning, graph neural networks (GNNs) [27, 51, 59] become dominant in the computing of essential tasks across a wide range of domains, like e-commerce, financial services, and etc. Compared with standard methods for graph analytics, such as random walk [18, 22, 50] and graph laplacians [8, 32, 33], GNNs highlight themselves with significantly higher accuracy [27, 54, 59] and better generality [19]. From the computation perspective, GNNs feature an interleaved execution phase of both graph operations (scatter-andgather [17]) at the *Aggregation* phase and Neural Network (NN) operations (matrix multiplication) at the *Update* phase. Our experimental studies further show that the aggregation phase which involves highly sparse computation on irregular input graphs generally takes more than 80% of the running time for both GNN training and inference. Existing GNN

frameworks, e.g., Deep Graph Library [55] and PyTorch Geometric [13], are mostly built upon the popular NN frameworks that are originally optimized for dense operations, such as general matrix-matrix multiplication (GEMM). To support sparse computations in GNNs, their common strategy is to incorporate sparse primitives (such as cuSPARSE [38]) for their backend implementations. However, cuSPARSE leverages the sparse linear algebra (LA) algorithm which involves lots of high-cost indirect memory accesses on non-zero elements of a sparse matrix. Therefore, cuSPARSE cannot enjoy the same level of optimizations (e.g., data reuse) as its dense counterpart, such as cuBLAS [40]. Moreover, cuSPARSE is designed to only utilize CUDA cores. Therefore, It cannot benefit from advancements in GPU hardware features, like Tensor Core Units (TCUs) on the recent NVIDIA Ampere and Hopper GPUs. Such a design is also the trend of many other AI-tailored accelerators/units (e.g., Google TPU [24] and Matrix Core [2] on AMD GPUs) and can significantly boost the performance of dense LA algorithms (e.g., GEMM and Convolution) in most conventional deep-learning applications (e.g., CV [20] and NLP [10]).

This work focuses on exploring the potential of TCUs for accelerating such GNN-based graph learning and our design/optimization principles will also benefit other similar AI hardware [2, 24] for sparse deep-learning workloads. We remark that making TCUs effective for general GNN computing is a non-trivial task. Our initial study shows that naively applying the TCU to sparse GNN computation would even result in inferior performance compared with the existing sparse implementations on CUDA cores. There are several challenges. *First*, directly resolving the sparse GNN computing problem with the pure dense GEMM solution is impractical due to the extremely large memory cost ($O(N^2)$), where N is the number of nodes). Besides, traversing the matrix tiles already known to be filled with all-zero elements would cause excessive unnecessary computations and memory access. Second, simply employing TCUs to process non-zero matrix tiles of the sparse graph adjacency matrix would still waste most of the TCU computation and memory access efforts. This is because TCU

input matrix tiles are defined with fixed dimension settings $(e.g., height(16) \times width(8))$, whereas the non-zero elements of a sparse graph adjacency matrix are distributed irregularly. Thus, it requires intensive zero-value padding to satisfy such a rigid input constraint. *Third*, although the recent CUDA release update enables TCUs to exploit the benefit of certain types of sparsity [37], it only supports blocked SpMM, where non-zero elements must first fit into well-shaped blocks and the number of blocks must be the same across different rows. Such an input restriction makes it hard to handle highly irregular sparse graphs in real-world GNN applications.

To this end, we introduce, **TC-GNN**¹, the first TCU-based GNN acceleration design on GPUs. Our key insight is to let the sparse input graph fit the dense computation of TCUs. At the input level, instead of exhaustively traversing all sparse matrix tiles and determining whether to process each tile, we develop a new sparse graph translation (SGT) technique that can effectively identify those non-zero tiles and condense nonzero elements from these tiles into fewer number of "dense" tiles. Our major observation is that neighbor sharing is very common among nodes in real-world graphs. Therefore, applying SGT can effectively merge the unnecessary data loading of the shared neighbors among different nodes to avoid highcost memory access. SGT is generic to any kind of sparse pattern of input graphs and can always yield the correct results as the original sparse algorithm. At the GPU kernel level, for efficiently processing GNN sparse workloads, TC-GNN exploits the benefits of CUDA core and TCU collaboration. The major design idea is that the CUDA core, which is more powerful at fine-grained thread-level execution, would be a good candidate for managing memory-intensive data access. While TCU, which is more powerful in handling simple arithmetic operations (e.g., multiplication and addition), would be well-suited for compute-intensive GEMM on dense tiles generated from SGT. At the framework level, we integrate TC-GNN with the popular PyTorch [49] framework. Thereby, users only need to interact with their familiar PyTorch programming environment by using TC-GNN APIs. This can significantly reduce extra learning efforts, and improve user productivity and code portability.

To sum up, we summarize our contributions as follows:

- We conduct a detailed analysis (§3) of existing solutions (*e.g.*, SpMM on CUDA cores) and identify the potentials of TCUs for accelerating sparse GNN workloads.
- We introduce a sparse graph translation technique (§4.1).
 It can make the sparse and irregular GNN input graphs easily fit the dense computing of TCUs for acceleration.
- We build a TCU-tailored algorithm (§4.2) and GPU kernel design (§4.3) for CUDA core and TCU collaboration on GPUs to handle different sparse GNN computation.

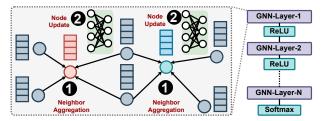


Figure 1: GNN General Computation Flow.

 Extensive experiments show TC-GNN achieves 1.70× speedup on average over the state-of-the-art GNN computing framework, Deep Graph Library, across various mainstream GNN models and dataset settings.

2 Background

2.1 Graph Neural Networks

Graph neural networks (GNNs) are an effective tool for graph-based machine learning. The detailed computing flow of GNNs is illustrated in Figure 1. GNNs basically compute the node feature vector (embedding) for node v at layer k+1 based on the embedding information at layer k ($k \ge 0$), as shown in Equation 1,

$$a_{v}^{(k+1)} = Aggregate^{(k+1)} (h_{u}^{(k)} | u \in \mathbf{N}(v) \cup h_{v}^{(k)})$$

$$h_{v}^{(k+1)} = Update^{(k+1)} (a_{v}^{(k+1)})$$
(1)

where $h_{\nu}^{(k)}$ is the embedding vector for node ν at layer k; $a_{\nu}^{(k+1)}$ is the aggregation results through collecting neighbors' information (e.g., node embeddings); N(v) is the neighbor set of node v. The aggregation method and the order of aggregation and update could vary across different GNNs. Some methods [19, 27] just rely on the neighboring nodes while others [54] also leverage the edge properties that are computed by applying vector dot-product between source and destination node embeddings. The update function is generally composed of standard NN operations, such as a fully connected layer or a multi-layer perceptron (MLP) in the form of $w \cdot a_v^{(k+1)} + b$, where w and b are the weight and bias parameters, respectively. The common choices for node embedding dimensions are 16, 64, and 128, and the embedding dimension may change across different layers. After several iterations of aggregation and update (i.e., several GNN layers), we will get the output feature embedding of each node, which can be used for various downstream graph learning tasks, such as node classification [11, 16, 25] and link prediction [6, 28, 53].

The sparse computing in the aggregation phase is generally formalized as the sparse-matrix dense-matrix multiplication (SpMM), as illustrated in Figure 2a, and is handled by many sparse libraries (*e.g.*, cuSPARSE [38]) in many state-of-the-art GNN frameworks [55,57]. These designs only count on GPU CUDA cores for computing, which waste the modern GPUs

¹https://github.com/YukeWang96/TC-GNN_ATC23.git

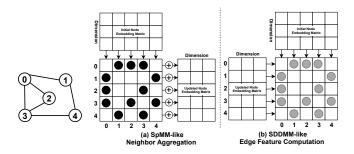


Figure 2: (a) SpMM-like and (b) SDDMM-like Operation in GNNs. Note that " \rightarrow " indicates loading data; " \oplus " indicates neighbor embedding accumulation.

with diverse computing units, such as the Tensor Core Unit (TCU). Specifically, we formalized the neighbor aggregation as SpMM-like operations (Equation 2)

$$\mathbf{\hat{X}} = (\mathbf{F}_{N \times N} \odot \mathbf{A}_{N \times N}) \cdot X_{N \times D}) \tag{2}$$

where A is the graph adjacency matrix stored in CSR format. X is a node feature embedding matrix stored in dense format. N is the number of nodes in the graph, and D is the size of node feature embedding dimension; \odot is the elementwise multiplication and \cdot is the standard matrix-matrix multiplication; F is the edge feature matrix in CSR format and can be computed by Sampled Dense-Dense Matrix Multiplication (SDDMM)-like operations (Equation 3 and Figure 2b).

$$\mathbf{F} = (\mathbf{X}_{N \times D} \cdot \mathbf{X}_{N \times D}^{T}) \odot \mathbf{A}_{N \times N} \tag{3}$$

Note that the computation of F is optional in GNNs, which is generally adopted by the Attention-based Graph Neural Network in PyTorch [51] for identifying more complicated graph structural information. Other GNNs, such as the Graph Convolutional Network [27] and Graph Isomorphism Network [59], only use the adjacency matrix for neighbor aggregation.

2.2 GPU Tensor Core

In the most recent GPU architectures (since Volta [43]), NVIDIA announced a new type of computing unit, Tensor Core Unit (TCU), for accelerating dense deep-learning operations (e.g., Dense GEMM). A GPU Streaming-Multiprocessor (w/TCU) is illustrated in Figure 3. Note that FP64, FP32, INT, and SFU are for double-precision, single-precision, integer, and special function units, respectively. Different from scalar computation on CUDA cores, TCU provides tile-based matrixmatrix computation primitives on register fragments, which can deliver more than 10× throughput improvement. In particular, TCU supports the compute primitive of $\mathbf{D} = \mathbf{A} \cdot \mathbf{B} + \mathbf{C}$, where **A** and **B** are required to be a certain type of precision (e.g., half, TF-32), while C and D are stored in FP32. Depending on the data precision and GPU architecture version, the matrix size (MMA shape) of $\mathbf{A}(M \times K)$, $\mathbf{B}(K \times N)$, and $\mathbb{C}(M \times N)$ should follow some principles [41]. For example, TF-32 TCU computing requires M = N = 16 and K = 8.

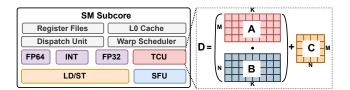


Figure 3: A Subcore of GPU SM with TCUs.

Listing 1: WMMA APIs for TCUs in CUDA C.

```
wmma::fragment<matrix_a, M, N, K, tf32, row_major> a_frag;
// Load tiles (global/shared mem. -> register fragments).
wmma::load_matrix_sync(a_frag, A, M);
// Execute GEMM on loaded tiles on register fragments.
wmma::mma_sync(c_frag, a_frag, b_frag, c_frag);
// Move results (register fragments -> global/shared mem).
wmma::store_matrix_sync(C, c_frag, N, mem_row_major);
```

In the recent CUDA release (>=11.0) on Ampere (*sm*>=80), TF-32 serves as a good alternative to float/double on TCU-based GPU computing for modern deep-learning applications, according to NVIDIA's in-depth studies [45].

Different from the CUDA cores that operate at the thread level (*e.g.*, allowing the "if" branch among threads), TCU supports only the operation at the warp level (*e.g.*, forbidding the "if" branch among threads within a warp). Before calling TCUs, all registers in a warp need to collaboratively store matrix tiles into a new memory hierarchy *Fragment* [48], which allows data sharing across registers. This intra-warp sharing provides opportunities for fragment-based memory optimizations. TCU can be utilized in several ways. The simplest way is to call cuBLAS [40] by using the cublasSgemmEX API. The second way is to call the Warp Matrix Multiply-Accumulate (WMMA) (nvcuda::wmma) API [47] in CUDA C to operate TCUs directly with four major operations (Listing 1).

Since the appearance of the TCU, research efforts have been devoted to accelerating deep-learning (DL) workloads with TCUs. Ang and Simon [31] leverage 1-bit GEMM capability on Turing TCUs for accelerating binary Neural Network inference. Boyuan et al. [12] introduce GEMM-based scientific computing on TCUs with extended precision and high performance. Yuke et al. [56] treat batched quantized GNNs (partitioning large graphs into small graphs as batches) as batched dense GEMM computation and accelerate it on TCUs for inference. These prior efforts use TCUs in the dense DL applications that TCU is initially designed for, while TC-GNN jumps out of the scope defined by TCU designers and accelerates the sparse full-graph GNNs using TCUs.

3 Motivation

In this section, we will discuss the major technical thrust for us to leverage TCUs for accelerating sparse GNN computation. We use the optimization of SpMM as the major example in this discussion, and the acceleration of SDDMM would also benefit from similar optimization principles.

Table 1: Profiling of GCN Sparse Operations.

Dataset	Aggr. (%)	Update (%)	Cache(%)	Occ.(%)
Cora	88.56	11.44	37.22	15.06
Citeseer	86.52	13.47	38.18	15.19
Pubmed	94.39	5.55	37.22	16.24

SpMM on CUDA cores 3.1

As the major component of sparse linear algebra operation, SpMM has been incorporated in many off-the-shelf libraries [1, 3, 5, 21, 38]. The close-sourced cuSPARSE [38] library developed by NVIDIA is the most popular solution and it can deliver state-of-the-art performance for most GPUbased SpMM computation. cuSPARSE has also been widely adopted by many GNN frameworks, such as Deep Graph Library (DGL) [55], as the backend for sparse operations. To understand its characters, we profile DGL on one layer of a GCN [27] model (neighbor aggregation + node update) on NVIDIA RTX3090. We report two key kernel matrices for only neighbor aggregation kernel, including L1/texture cache hit rate (Cache) and the achieved Streaming-Multiprocessor (SM) occupancy (Occ.). We select three representative GNN datasets: Cora with 3,327 nodes, 9,464 edges, and 3,703 node embedding dimensions; Citeseer with 2,708 nodes, 10,858 edges, and 1,433 dimensions; Pubmed with 19,717 nodes, 88,676 edges, and 500 dimensions.

From Table 1, we have several observations: *First*, the aggregation phase usually dominates the overall execution of the GNN execution. From these three commonly used GNN datasets, we can see that the aggregation phase usually takes more than 80% of the overall execution time, which demonstrates the key performance bottleneck of the GNNs is to improve the performance of the sparse neighbor aggregation. **Second**, sparse operations in GNNs show very low memory performance. The column Cache of Table 1 shows GNN sparse operations could not well benefit from the GPU cache system, thus, showing a low cache-hit ratio (around 37%) and frequent global memory access. *Third*, sparse operations of GNNs show very inefficient computation. As described in the column Occupancy of Table 1, the sparse operation of GNNs could hardly keep the GPU busy because 1) its low computation intensity (the number of non-zero elements in the sparse matrix is generally small); 2) its highly irregular memory access for fetching rows of the dense matrix during the computation, resulting in memory-bound computation; 3) it currently can only leverage CUDA cores for computation, which naturally has limited throughput performance. On the other side, this study also points out several potential directions for improving the SpMM performance on GPUs, such as improving the computation intensity (e.g., assigning more workload to each thread/warp/block), boosting memory access efficiency (e.g., crafting specialized memory layout for coalesced memory access), and breaking the computation performance ceiling (e.g., using TCUs).

Table 2: Medium-size Graphs in GNNs.

Dataset	# Nodes	# Edges	Memory	Eff.Comp
OVCR-8H	1,890,931	3,946,402	14302.48 GB	0.36%
Yeast	1,714,644	3,636,546	11760.02 GB	0.32%
DD	334,925	1,686,092	448.70 GB	0.03%

Dense GEMM on CUDA Cores/TCUs

While the dense GEMM is mainly utilized for dense NN computation (e.g., linear transformation and convolution), it can also be leveraged for GNN aggregation under some circumstances. For example, when an input graph has a very limited number of nodes, we can directly use the dense adjacency matrix of the graph and accelerate the intrinsically sparse neighbor aggregation computation on CUDA cores/TCUs by calling cuBLAS [40]. However, such an assumption may not hold even for medium-size graphs in real-world GNNs.

As shown in Table 2, for these selected datasets, the memory consumption of their dense graph adjacent matrix (as a 2D float array) would easily exceed the device memory constraint of today's GPU (less than 100GB). Even if we assume the dense adjacent matrix can fit into the GPU memory, the extremely low effective computation (the last column of Table 2) would also be a major obstacle for us to achieve high performance. We measure the effective computation as $\frac{nnz}{N\times N}$, where nnz is the number of the non-zero elements (indicating edges) in the graph adjacent matrix and N is the number of nodes in the graph. The number of *nnz* is tiny in comparison with the $N \times N$. Therefore, computation and memory access on zero elements are wasted.

3.3 Hybrid Sparse-Dense Solution

Another type of work [29, 37] takes the path of mixing the sparse control (tile-based iteration) with Dense GEMM computation. They first apply a convolution-like (2D sliding window) operation on the adjacent matrix and traverse all possible dense tiles that contain non-zero elements. Then, for all identified non-zero tiles, they invoke GEMM on CUDA cores/TCUs for computation. However, this strategy has two shortcomings. First, the sparse control itself would cause a high overhead. Based on our empirical study, the non-zero elements are highly scattered on the adjacent matrix of a sparse graph. Therefore, traversing all blocks in a super large adjacent matrix would be time-consuming. **Second**, the identified sparse tiles would still waste lots of computation. The irregular edge connections of the real-world graphs could hardly fit into these fixed-shape tile frames. Therefore, most of the dense tiles would still have very few non-zero elements.

Inspired by the above studies, we make several design choices in order to achieve high-performance sparse GNN operations. First, we choose the hybrid sparse-dense solution as the starting point. This can give us more flexibility for optimizations at the sparse control (e.g., traversing fewer tiles)

Table 3: Comparison among Sparse GEMM, Dense GEMM, Hybrid Sparse-Dense, and TC-GNN. Note that MC: Memory Consumption, EM: Effective Memory Access, CI: Computation Intensity, EC: Effective Computation.

Solution	MC	EM	CI	EC
Sparse GEMM (§3.1)	Low	Low	Low	High
Dense GEMM (§3.2)	High	High	High	Low
Hybrid Sparse-Dense (§3.3)	High	Low	Low	High
TC-GNN (This work)	Low	High	High	High

and dense computation (*e.g.*, increasing the effective computation/memory access when processing each tile). <u>Second</u>, we employ shared memory as the key space for GPU kernel-level data management. It can help us to re-organize the irregular GNN input data in a more "regularized" way such that both the memory access efficiency and computing performance can be well improved. <u>Third</u>, we choose TCUs as our major computing unit since they can bring significantly higher computing throughput performance in comparison with CUDA cores. This also indicates the great potential of using TCUs for harvesting more performance gains.

Finally, we crystallize all of our ideas and insights into TC-GNN that effectively coordinates the execution of GNN sparse operations on dense TCU. We show a brief qualitative comparison among TC-GNN and the above three solutions in Table 3. Note that *Memory Consumption* is the size of memory used by the sparse/dense graph adjacency matrix; The *Effective Memory Access* is the ratio between the size of the accessed data that is actually involved in the later computation and the total data being accessed; The *Computation Intensity* is the ratio of computing operations versus the data being accessed; The *Effective Computation* is the operations for generating the final result versus the total operations.

4 TC-GNN Design

In this section, we will first give an overview of TC-GNN through its high-level programming interface and then detail the TCU-aware GNN algorithm design. As detailed in Listing 2, TC-GNN consists of several key components to facilitate the programming of GNN models on GPU TCUs. TC-GNN introduces a set of pre-built popular GNN layers (e.g., TCGNN.GCNConv) that can be easily connected with some other existing neural network layers (e.g., ReLU and softmax), to help users define their own GNN model quickly. For those non-conventional GNN layers, users can directly use our low-level APIs (e.g., TCGNN.spmm and TCGNN.sddmm) to express the GNN computation easily. TC-GNN introduces an input Loader to load the GNN input graph as a rawGraph and capture the key input information for system-level optimizations. TC-GNN incorporates a Preprocessor to build tiles from rawGraph and generate TCU-aware tiledGraph (§4.1), and optimize runtime configuration (e.g., warps per block) for

Listing 2: Example of a 2-layer GCN in TC-GNN.

```
import TCGNN, torch
  # include other packages ...
  class GCN(torch.nn.Module):
      def .
            _init__(self, inDim, hiDim, outDim):
          self.layer1 = TCGNN.GCNConv(inDim, hiDim)
          self.layer2 = TCGNN.GCNConv(hiDim, outDim)
          self.softmax = torch.nn.Softmax()
      def forward(self, tiledGraph, param):
          tiled_adj, X = tiledGraph.adj, tiledGraph.X
          X = self.layer1(X, tiledAdj, param)
11
          X = self.ReLU(X)
          X = self.layer2(X, tiledAdj, param)
          X = self.softmax(X)
          return X
  # Define a two-layer GCN model in TC-GNN.
  model = GCN(inDim=100, hiDim=16, outDim=10)
  # Load graph and extract input information.
  rawGraph, info = TCGNN.Loader(graphFilePath)
  # Generate TCU tile and runtime configuration.
  tiledGraph, config = TCGNN.Preprocessor(rawGraph, info)
  # Run model through forward computation.
  predict_y = model(tiledGraph, config)
  # Compute loss and accuracy
  # Gradient backpropagation for training.
```

our TCU-tailored GPU kernel (§4.2 and §4.3) based on input. Finally, we train the initialized GNN model defined in TC-GNN as the regular GNN models defined in other frameworks through forward and backward computation.

4.1 TCU-aware Sparse Graph Translation

As the major component of TC-GNN, we introduce a novel *Sparse Graph Translation* (SGT) technique to facilitate the TCU acceleration of GNNs. Our core idea is that *the pattern of the graph sparsity can be well-tuned for TCU computation through effective graph structural manipulation meanwhile guaranteeing output correctness*. Our key observation is that neighbor sharing is common in real-world graphs and has been exploited for various tasks like link prediction [63]. Our evaluated datasets (Section 5) have 18% to 47% (averaged 29%) neighbor similarity. Specifically, we condense (remap) the highly-scattered neighbor ids into highly-condensed new neighbor ids that can facilitate the dense TCU computation paradigm. Also, such condensing should not compromise any original information (*e.g.*, edge connections) and can generate the exact output as the conventional design.

As exemplified in Figure 4a and Figure 4b, we take the regular graph in CSR format as the input and condense the columns of each row window (in the red-colored rectangular box) to build TCU blocks (TC_block) (a.k.a., the input operand shape of a single MMA instruction), in the orange-colored rectangular box. nodePointer is the row pointer array edgeList is the edges of each node stored continuously. In this paper, we demonstrate the use of standard MMA shape for TF-32 of TCU on Ampere GPU architecture, and other MMA shapes [41] can also be used under different precision (e.g., half and int8) and GPU architecture (e.g., Turing).

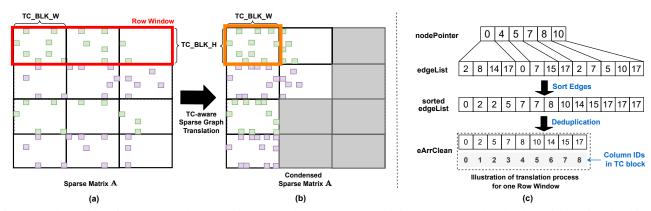


Figure 4: Illustration of Sparse Graph Translation. Note that the grey area indicates the TCU blocks that will be directly skipped.

```
Algorithm 1: TCU-aware Sparse Graph Translation.
  input: Graph adjacent matrix A (nodePointer, edgeList).
  output: Result of winPartition and edgeToCol.
  /* Compute the total number of row windows.
1 numRowWin = ceil(numNodes/winSize):
2 for winId in numRowWin do
      /* EdgeIndex range of the current rowWindow.
      winStart = nodePointer[winId * winSize];
3
      winEnd = nodePointer[(winId + 1) * winSize];
4
      /* Sort the edges of the current rowWindow.
      eArray = Sort(winStart, winEnd, edgeList);
      /* Deduplicate edges of the current rowWindow.
      eArrClean = Deduplication(eArray);
      /* #TC blocks in the current rowWindow.
      winPartition[winId] =
7
       ceil(eArrClean.size/TC_BLK_w);
      /* Edges-to-columnID mapping in TC Blocks.
      for eIndex in [winStart, winEnd] do
          eid = edgeList[eIndex];
          edgeToCol[eIndex] = eArrClean[eid];
10
      end
11
12 end
```

SGT takes several steps for processing each row window, as detailed in Algorithm 1 and visualized in Figure 4c. win-Partition is an array for maintaining the number of TC blocks in each row window. edgeToCol is an array for maintaining the mapping between the edges and their corresponding position in the graph after SGT. Note that edgeToCol has the same length as edgeList but with column-id from eArrClean. colToRow maps column-id of adjacency matrices to the row-id of embedding matrices. We choose the size of the row window (winSize=TC BLK H) and column width (TC_BLK_W) according to TCU MMA specification (e.g., TC_BLK_H=16, TC_BLK_W=8 in TF-32). After condensing the graph within each row window, the time complexity of sliding the TC_block can be reduced from $O(\frac{N}{TC_BLK_W})$ to only $O(\frac{nnz_{unique}}{TC_BLK_W})$, where *N* is the total number of nodes in the graph and $\overline{nnz_{unique}}$ is the size of the unique neighbor within the current row window, which equals eArrClean.size

```
Algorithm 2: TC-GNN Neighbor Aggregation.
  input: Condensed graph structure (nodePointer, edgeList,
          edgeToCol, winPartition) and node embedding matrix (X).
  output : Updated node embedding matrix (\hat{\mathbf{X}}).
  /* Traverse through all row windows.
1 for winId in numRowWindows do
       /* #TC blocks of the row window.
       numTCblocks = winPartition[winId];
       /* Edge range of TC blocks of the row window.
       edgeRan = GetEdgeRange(nodePointer, winId);
3
       for TCblkId in numTCblocks do
4
           /* The edgeList chunk in current TC block. */
           edgeChunk = GetChunk(edgeList, edgeRan, TCblkId);
5
           /* Neighbor node Ids in current TC block.
           colToNId = GetNeighbors(edgeChunk, edgeToCol);
           /* Initiate a dense tile (ATile).
           ATile = InitSparse(edgeChunk, winId);
           /* Initiate a dense tile (XTile)
           XTile, colId = FetchDense(colToNId, X);
           /* Compute XnewTile via TCU GEMM.
           XnewTile = TCcompute(ATile, XTile);
           /* Store XnewTile of \hat{X}.
           \hat{X} = StoreDense(XNewTile, winId, colId);
10
       end
11
12 end
```

in Algorithm 1. The density (computation intensity) of each identified TCU block can be largely improved. Considering the case in Figure 4, after the sparse graph translation, we can achieve $2 \times$ higher density on individual TCU blocks (Figure 4b) compared with the original one (Figure 4a).

Compared to existing sparse matrix formats (e.g., Blocked-Ellpack [37]) which use the regular matrix tiles to cover the irregularly scattered non-zero elements, SGT reduces the irregularity of non-zero-elements layout to fit them into fewer number TCU blocks, thus, reducing the unnecessary computation and memory overhead. SGT is applicable for both the SpMM and SDDMM in GNN sparse operations and can be easily parallelized because the processing of individual row windows is independent. In most cases, SGT only needs to execute once and its result can be reused across many epochs/rounds of GNN training/inference.

Additionally, SGT can be generally used with other ac-

celerators (e.g., AMD-GPUs with matrixCore and TPUs) that offer similar dense MM primitives. CPUs have no direct alternative to TensorCore-like MM primitives. However, with AVX-vectorized instructions, CPUs can benefit from SGT by setting *BLK_H*=1 and *BLK_W*=(#elements-per-AVX-instruction). TC-GNN currently targets GNN training. SGT is conducted once before training. SGT cost can be offset by training iterations (averaged 2% for 200 iterations as DGL).

4.2 TCU-tailored GNN Computation

Besides the effective way to condense the sparse tiles, the next major challenge is *how to tailor the computation schedule of GNN algorithms* so that we can capitalize on the performance of condensed sparse graphs and the powerful TCUs. We focus on two major types of computation in GNNs.

Neighbor Aggregation The major part of GNN sparse computing is neighbor aggregation, which can generally be formalized as SpMM operations by many state-of-the-art frameworks [55]. And they employ the cuSPARSE [38] on CUDA cores as a black-box technique for supporting sparse GNN computation. In contrast, our TC-GNN design targets at TCU for the major neighbor aggregation computation which demands a specialized algorithmic design. TC-GNN focuses on maximizing the net performance gains by gracefully batching the originally highly irregular SpMM as dense GEMM computation and solving it on TCU effectively. As illustrated in Algorithm 2, the node aggregation processes all TC blocks from each row window. nodePointer and edgeList are directly from graph CSR, while edgeToCol and winPartition are generated from SGT discussed in the previous section. Note that **InitSparse** is to initialize a sparse tile in dense format according to the translated graph structure of the current TC block. Meanwhile, FetchDense returns a dense node embedding matrix tile XTile for TCU computation, and the corresponding column range *colld* (embedding dimension range) of matrix **X**. This is to handle the case that the width of one *XTile* could not cover the full-width (all dimensions) of **X**. Therefore, the colld will be used to put the current TCU computation output to the correct location in the updated embedding matrix $\hat{\mathbf{X}}$.

Edge Feature Computing Previous research [51,54] has demonstrated the great importance of incorporating the edge feature for a better GNN model algorithmic performance (e.g., accuracy, and F1-score). The underlying building block to generate edge features is the Sampled Dense-Dense Matrix Multiplication (SDDMM)-like operation. In TC-GNN, we support SDDMM with the collaboration of the above sparse graph translation and TCU-tailored algorithm design, as described in Algorithm 3. The overall algorithm structure and inputs are similar to the above neighbor aggregation. The major difference is the output. In the case of neighbor aggregation, our output is the updated dense node embedding matrix ($\hat{\mathbf{X}}$), where edge feature computing will generate a sparse output with the same shape as the graph edge lists.

Algorithm 3: TC-GNN Edge Feature Computation.

```
input :Condensed graph data (nodePointer, edgeList, edgeToCol,
          winPartition) and node embedding matrix (X).
  output: Edge Feature List (edgeValList).
  /* Traverse through all row windows.
1 for winId in numRowWin do
       /* #TC blocks in the row window.
       numTCblocks = winPartition[winId]:
       /* Edge range of TC blocks of the row window.
       edgeRan = \mathbf{GetEdgeRange}(nodePointer, winId);
3
       for TCblkId in numTCblocks do
           /* EdgeList chunk in current TC block.
           edgeChunk = GetChunk(edgeList, edgeRan, TCblkId);
           /* Neighbor node Ids in current TC block.
           colToNId = GetNeighbors(edgeChunk, edgeToCol);
           /* Fetch a dense tile (XTile_A).
           XTile_A = \mathbf{FetchDenseRow}(winId, TCblkId, X);
           /* Fetch a dense tile (XTile_R).
           XTile_B = \mathbf{FetchDenseCol}(colToNId, edgeToCol, X);
           /* Compute edgeValTile via TCU GEMM.
           edgeValTile = TCcompute(XTile_A, XTile_B);
           /* Store edgeValTile to edgeValList.
           StoreSparse(edgeValList, edgeValTile,
10
11
                        edgeList, edgeToCol);
12
       end
13 end
```

Note that fetching the $XTile_A$ only needs to consecutively access the node embedding matrix \mathbf{A} by rows while fetching the $XTile_B$ requires first computing the TCU block column-id to node-id (colToNId) to fetch the corresponding neighbor node embeddings from the same node embedding matrix \mathbf{X} .

Despite the dataflow similarity with dense-GEMM computation (e.g., CUTLASS [39]), TC-GNN has to overcome the limited parallelism (imbalance workload) and sparse/irregular access with novel algorithmic and kernel designs. While these challenges are absent in dense-GEMM computation with naturally high parallelism and data-access locality.

4.3 TCU-centric Workload Mapping

In collaborating with our TCU-tailored algorithm design, an effective mapping of our algorithmic design to low-level GPU primitives is indispensable for high-performance delivery. We discuss two key techniques: *GPU-aware Workload Decomposition* and *TCU-optimized dataflow design*.

4.3.1 GPU-aware Workload Decomposition

Different from previous work [13, 55] focusing on CUDA cores only, TC-GNN highlights itself with CUDA core and TCU collaboration through effective two-level workload mapping. The idea is based on the fact that CUDA cores work in SIMT fashion and are operated by individual threads, while TCU designated for GEMM computation requires collaboration from a warp of threads (32 threads). Our key design principle is to *mix these two types of computing units as a single GPU kernel*, which can efficiently coordinate the kernel

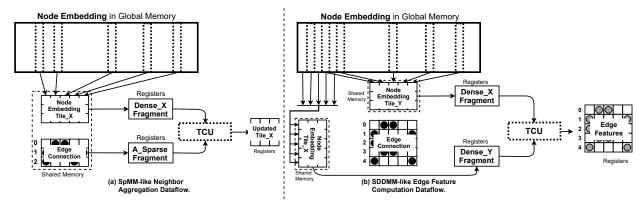


Figure 5: TCU-optimized Dataflow Design for (a) Neighbor Aggregation and (b) Edge Feature Computing in GNNs.

execution at different levels of execution granularity.

In TC-GNN, we operate CUDA cores by thread blocks and manage TCU by thread warps. Specifically, threads running CUDA cores from the same thread block will load data (e.g., edges) from the global memory to shared memory. Note that in our design we assign each row window (discussed in §4.1) to one thread block. The number of threads in each block should be divisible by the number of threads in each warp (32) for better performance. Once threads running on CUDA cores (CUDA-core threads) finish the data loading, threads from each warp (TCU threads) will operate TCU for GEMM computation (including loading the data from the shared memory to thread-local registers (fragments), applying GEMM computation on data in registers, accumulating results on registers, and storing the final results back to global memory). Note that there would be a large overlap of the CUDA-core threads and TCU threads, both of which are threads from the same blocks but running at a different time frames. In general, we use more CUDA-core threads than TCU threads considering that global memory access demands more parallelization.

There are two major benefits of such two-level workload decomposition. First, threads from the same block can work together to improve the memory access parallelization to better utilize memory bandwidth. Second, warps from the same block can reuse the loaded data, including the information (*e.g.*, column index mapping) of the translated graph and the tiles from the dense node embedding matrix. Therefore, we can avoid redundant high-cost global memory operations.

4.3.2 TCU-optimized Dataflow Design

As the major technique to improve the GPU performance, shared memory is customized for our TCU-based sparse kernel design for re-organizing data layout for dense TCU computation and reducing the redundant global memory traffic. Our design takes the TCU specialty into careful consideration from two aspects, 1) the input matrix tile size of the TCU, which is $M(16)\times N(16)\times K(8)$ in the case of TF-32, and 2) the tile fragment layout for fast computation. The common

practice of the loaded tile A and B are stored in row-major and column-major for better performance. Next, we will detail our TCU-optimized dataflow design for both neighbor aggregation and edge feature computation.

Neighbor Aggregation In Figure 5a, shared memory is mainly used for caching several most frequently used information, including the tile of sparse matrix A (sparse_A), the column-id of the sparse matrix A to row-id of node embedding matrix **X** (sparse_AToX_index), and the dense tile of **X** (dense_X). When handling each TCU block, we assign all threads from the same block of threads for loading the sparse tile while allowing several warps to concurrently load the dense row tile from the matrix X. The reasons for enforcing such caching are two-fold. First, it can bridge the gap between the sparse graph data and the dense GEMM computing that requires continuous data layout. For example, the adjacent matrix A is input as CSR format that cannot be directed feed to TCU GEMM computation, therefore, we use a shared memory sparse_A to initialize its equivalent dense tile. Similarly, we cache rows of X according to the columns of A to the row of X mapping after our sparse graph translation, where originally scattered columns of A (the rows of X) are condensed. Second, it can enable data reuse on sparse_AToX_index and sparse_A. This is because in general, the BLK_H (16) cannot cover all dimensions of a node embedding (e.g., 64), multiple warps will be initiated of the same block to operate TCU in parallel to work on non-overlapped dense tiles while sharing the same sparse adjacency matrix tile.

Edge Feature Computation Similar to the shared memory design in neighbor aggregation, for edge feature computing, as visualized in Figure 5b, the shared memory is utilized for sparse_A, sparse_ATOX_index, and dense_X. We assign all threads from the same block of threads for loading the sparse tile while allowing several warps to concurrently load the dense row tile from the matrix **X**. Compared with dataflow design in neighbor aggregation, edge feature computing demonstrates several differences.

First, the sizes of sparse_A are different. In the neighbor aggregation computation, the sparse matrix **A** is used as

Table 4: Datasets for evaluation.

Type	Dataset	Abbr.	#Vertex	#Edge	Dim.	#Class
I	Citeseer	CR	3,327	9,464	3703	6
	Cora	CO	2,708	10,858	1433	7
	Pubmed	PB	19,717	88,676	500	3
	PPI	PI	56,944	818,716	50	121
	PROTEINS_full	PR	43,471	162,088	29	2
	OVCAR-8H	OV	1,890,931	3,946,402	66	2
II	Yeast	YT	1,714,644	3,636,546	74	2
	DD	DD	334,925	1,686,092	89	2
	YeastH	YH	3,139,988	6,487,230	75	2
Ш	amazon0505	AZ	410,236	4,878,875	96	22
	artist	AT	50,515	1,638,396	100	12
	com-amazon	CA	334,863	1,851,744	96	22
	soc-BlogCatalog	SC	88,784	2,093,195	128	39
	amazon0601	AO	403,394	3,387,388	96	22

one operand in the SpMM-like computation, therefore, the minimal processing granularity is 16 × 8, while in edge feature computing by following SDDMM-like operation, the sparse matrix A serves as the output matrix, thus, maintaining the minimum processing granularity is 16×16 . To reuse the same translated sparse graph as SpMM, we need to recalculate the total number of TC blocks. Second, iterations along the embedding dimension would be different. Compared with neighbor aggregation, edge feature computing requires the result accumulation along the embedding dimension. The result will only be output until all iterations have finished. In neighbor aggregation, the node embedding vector is divided among several warps, each of which will output their aggregation result to non-overlapped embedding dimension ranges in parallel. *Third*, the output format has changed. Compared with SpMM-like neighbor aggregation which directly output computing results as an updated dense matrix $\hat{\mathbf{X}}$, SDDMM-like edge feature computing requires a sparse format (the same shape as edgeList) output for compatibility with neighbor aggregation and memory space. Therefore, one more step of dense-to-sparse translation is employed.

5 Evaluation

Benchmarks: We choose two representative GNN models widely used by previous work [13, 34, 55] on node classification tasks. Specifically, 1) Graph Convolutional Network (GCN) [27] is one of the most popular GNN model architectures. It is also the key backbone for many other GNNs (e.g., GraphSAGE [19] and differentiable pooling (Diffpool) [61]). Therefore, improving the performance of GCN will also benefit a broad range of GNNs. For GCN evaluation, we use the setting: 2 layers with 16 hidden dimensions per layer, which is also the setting from the original paper [27]. 2) Attentionbased Graph Neural Network (AGNN) [51]. AGNN differs from GCN in its aggregation function, which computes edge features (via embedding vector dot-product between source and destination vertices) before the node aggregation. AGNN is also the reference architecture for many other recent GNNs for better model algorithmic performance. For AGNN, we

use: 4 layers with 32 hidden dimensions per layer.

Baselines: 1) Deep Graph Library (**DGL**) [55] is the stateof-the-art GNN framework on GPUs, which is built with the high-performance CUDA-core-based cuSPARSE [38] library as the backend and uses PyTorch [49] as its front-end programming interface. DGL significantly outperforms other existing GNN frameworks [13] over various datasets on many mainstream GNN model architectures. Therefore, we make an in-depth comparison with DGL. 2) PyTorch Geometric (PyG) [13] is another GNN framework. PyG leverages torchscatter [14] library (highly-engineered CUDA-core kernel) as the backend support, which highlights its performance on batched small graph settings; 3) Blocked-SpMM [37] (**bSpMM**) accelerates SpMM on TCU. It is included in the recent update on the cuSPARSE library. bSpMM requires the sparse matrix with Blocked-Ellpack format for computation. Its computation on non-zero blocks can be seen as the hybrid sparse-dense solution (§3.3). Note that the bSpMM has not been incorporated into any existing GNN frameworks. We also compare TC-GNN with tSparse [62] and Triton [52] for non-vendor-developed highly optimized kernels on TCUs.

Datasets, Platforms, and Metrics: We cover three types of datasets (Table 4), which have been used in previous GNNrelated work [13, 34, 55]. Specifically, Type I graphs are the typical datasets used by previous GNN algorithm papers [19, 27, 59]. They are usually small in the number of nodes and edges, but rich in node embedding information with high dimensionality. Type II graphs [26] are the popular benchmark datasets for graph kernels and are selected as the built-in datasets for PyG [13]. Each dataset consists of a set of small graphs, which only have intra-graph edge connections without inter-graph edge connections. Type III graphs [27, 30] are large in terms of the number of nodes and edges. These graphs demonstrate high irregularity in its structures, which are challenging for most of the existing GNN frameworks. The core design of TC-GNN consists of around 2.5K lines of code. TC-GNN backend is implemented with C++ and CUDA C, and its front end is implemented in Python. Our major evaluation platform is a server with an 8-core 16thread Intel Xeon Silver 4110 CPU and an NVIDIA RTX3090 GPU. To measure the performance speedup, we calculate the average latency of 200 end-to-end runs.

5.1 Compared with DGL

Figure 6a shows that TC-GNN achieves $1.70\times$ speedup on average compared to DGL over three types of datasets across GCN and AGNN models on end-to-end training. Our kernel profiling via Nsight Compute shows that TC-GNN achieves high SM occupancy (averaged 85.28%), which is on average 21.05% higher compared to DGL across all datasets.

Type I Graphs: The performance improvements against DGL are significantly higher for GCN (on average $2.23\times$) compared to AGNN (on average $1.93\times$). The major reason

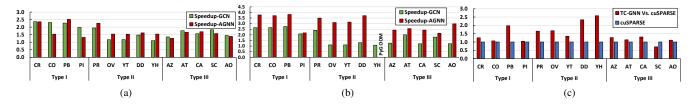


Figure 6: Speedup over (a) DGL and (b) PyG on GCN and AGNN; (c) Speedup over cuSPARSE bSpMM on TCUs.

is their different GNN computation patterns. GCN only consists of a neighbor aggregation phase (SpMM-like operation) and a node update phase (GEMM operation). Whereas in the AGNN, the aggregation phase would also require an additional edge attention value (feature) computation based on SDDMM-like operations. Compared with SpMM-like operations, edge attention computation in SDDMM is more sensitive to the irregular sparse graph structure because of much more intensive computations and memory access. Thus, the performance improvement is relatively lower.

Type II Graphs: TC-GNN achieves averaged 1.38× speedup on GCN and 1.70× speedup on AGNN for the Type II graphs. Speedup on Type II graphs is relatively lower compared with Type I, since Type II datasets consist of a set of small graphs with very dense intra-graph connections but no inter-graph edges. This leads to a lower benefit from the sparse graph translation that would show more effectiveness on highly irregular and sparse graphs. Such a clustered graph structure would also benefit cuSPARSE due to more efficient memory access, *i.e.*, less irregular data fetching from the sparse matrix. In addition, for AGNN, TC-GNN can still demonstrate evident performance benefits over the DGL (CUDA core only) that can mainly contribute to TCU-based SDDMM-like designs that can fully exploit the power of GPU through an effective TCU and CUDA core collaboration.

Type III Graphs: The speedup is also evident (on average 1.59× for GCN and average 1.51× for AGNN) on graphs with a large number of nodes and edges and irregular graph structures. The reason is the high overhead global memory access can be well reduced through our spare graph translation. Besides, our dimension-split strategy further facilitates efficient workload sharing among warps by improving the data spatial/temporal locality. On the dataset AT and SC, which have a higher average degree within Type III datasets, we notice a better speedup performance for both GCN and AGNN. This is because 1) more neighbors per node can lead to a higher density of non-zero elements within each tile/fragment. Thus, it can fully exploit the computation benefits of each TCU GEMM operation; 2) it can also facilitate more efficient memory access. For example, in AGNN, fetching one dense embedding x from the dense matrix X can be reused more times by applying a dot-product between x and many columns of the dense matrix X^T (neighbors embeddings).

Table 5: Compare TC-GNN with tSparse and Triton.

Dataset	tSparse (ms)	Triton (ms)	TC-GNN (ms)
AZ	18.60	31.64	4.09
AT	9.15	12.86	3.06
CA	13.84	15.50	3.26
SC	9.74	14.38	3.59
AO	11.93	21.78	3.41

Additionally, our performance breakdown analysis shows that for graphs with highly scattered and irregular edge distribution, such as Type I and III graphs, SGT would contribute more (averaged 64%) to the overall performance improvements since it helps significantly reduce the unnecessary workload. For graphs with highly dense and more regular edge connections, such as Type II datasets, SGT contributes relatively minor (averaged 23%) to the overall performance since it could not squeeze out more redundant computations from already condensed edge tiles.

5.2 Compared with other baselines

Compared with PyG Figure 6b shows TC-GNN can outperform PyG with an average of 1.76× speedup on GCN and an average of 2.82× speedup on AGNN. For GCN, TC-GNN achieves significant speedup on datasets with high-dimensional node embedding, such as *Yeast (YT)*, through effective TCU acceleration through a TCU-aware sparse graph translation while reducing the synchronization overhead by employing our highly parallelized TCU-tailored algorithm design. PyG, however, achieves inferior performance because its underlying GPU kernel can only leverage CUDA cores, thus, intrinsically bounded by CUDA core performance.

Compared with cuSPARSE bSpMM Figure 6c shows that TC-GNN outperforms bSpMM with on average 1.76× speedup on neighbor aggregation and improves effective computation by 75.8% on average. Our SGT technique can maximize the non-zero density of each non-zero tile and significantly reduce the number of non-zero tiles to be processed. However, bSpMM in cuSPARSE has to comply with the strict input sparse pattern (indicated in official API documentation [42]). For example, all rows in the arrays must have the same number of non-zero blocks. Thus, more redundant computations (on padding non-structural non-zero blocks) in bSpMM lead to inferior performance. We also notice that for

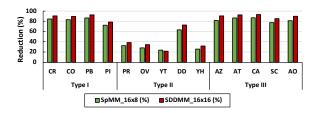


Figure 7: SGT Effectiveness on SpMM and SDDMM.

SC datasets with a high average node degree and clustered node distribution, bSpMM would benefit more due to its usage of a larger block size of 32×32 (fewer TCU invocations) compared to 16×8 in TC-GNN (more TCU invocations).

Compared with tSparse and Triton From Table 5, TC-GNN can outperform tSparse with on average 3.60× speedup on SpMM. The major reason behind this is that TC-GNN can well reduce the graph structural-level irregularity through our novel SGT strategy to benefit the dense TCU-based computation. In contrast, tSparse only considers partitioning the input sparse matrix into dense/sparse tiles based on their non-zero elements but ignores the potential of compressing non-zero elements into fewer tiles to reduce the workload. TC-GNN also outperforms Triton with on average 5.42× speedup on SpMM. Triton's block-sparse GEMM for TCU acceleration is designed for dense neural networks (focusing on feature maps' sparsity), which is quite different from GNNs (focusing on the graph adjacency matrix's sparsity) with significantly larger sparse matrix size and more irregular pattern.

5.3 Additional Studies

SGT Effectiveness & Overhead We conduct a quantitive analysis of SGT in terms of the total number of TCU blocks between graphs w/o SGT and the graphs w/ SGT applied. Note that in the SpMM-based aggregation, the size of TCU blocks is 16×8 since it serves as one of the operands in TCU GEMM. While in SDDMM-based edge feature computation, the size of TCU blocks is 16×16 since it serves as the resulting matrix of TCU GEMM. Figure 7 shows that across all types of datasets, our SGT technique can significantly reduce the number of traversed TCU blocks (on average 67.47%). The major reason is that SGT can largely improve the density of non-zero elements within each TCU block. In contrast, the graphs w/o SGT would demonstrate a large number of highly sparse TCU blocks. What is also worth noticing is that on Type II graphs, such a reduction benefit is lower. The reason is that Type II graphs consist of a set of small subgraphs that only maintain the intra-subgraph connections, which already maintain dense columns. We evaluate the overhead of SGT (Figure 8), we find that its overhead is consistently low (on average 4.43%) compared with the overall training time (200 epoches as DGL [55]).

Sparsity Analysis We compare with bSpMM on synthetic matrix data with different sparsity (zero-element ratio). Note



Figure 8: The overhead analysis of SGT.

Table 6: Sparsity Analysis. Numbers for bSpMM/TC-GNN are in GFLOPs. "DB/W": dense blocks per row window.

DB/W	Sparsity (%)	bSpMM	TC-GNN
1	99.61	773.86	12,686.02
2	99.22	1,597.83	11,010.75
4	98.44	3,348.75	18,164.08
8	96.88	6,528.10	25,883.10
16	93.75	12,955.40	23,865.99
32	87.50	26,061.70	16,629.28

that we change the sparsity by varying the number of dense non-zero blocks (16×16) within each row window, the input adjacent matrix size is fixed to 4096×4096 while the dense embedding matrix dimension is fixed to 16. Table 6 shows that when sparsity increases from 93.75% to 99.61%, TC-GNN design demonstrates more throughput performance strength (averaged $6.9\times$) and this is also the common sparsity range (more than 95%) for most input graphs of GNNs. When the sparsity drops to around 87.50% the sparse would demonstrate more advantage due to more dense blocks for computation.

Warps per Block: Figure 9 shows that with the increase of the number of warps, the overall performance for training per epoch would first decrease due to the better parallelism for loading the graph data. However, the number of warps per block would decrease the overall performance under certain circumstances (e.g., 32). All three settings suffer from evident performance degradation. Because the global memory access contention will become severe, leading to lower execution performance. Different datasets would have different "optimal" choices of the warp-per-block parameter. For example, on the CA dataset, 2 warps per block can deliver the best performance, while AZ requires 8 warps per block. Based on our profiling and empirical study, the selection of this parameter should consider the average #edges per row window (avg.edges), which can be easily get during the preprocessing. Our **preprocessor** will generate $warpPerBlock = \lfloor \frac{avg.edge}{32} \rfloor$ to approach the "optimal" performance. For instance, the average edges per row window are 88 for CA, it reaches the best performance at 2 warps per block.

Throughput Analysis: For sparse matrix computations in GNNs, we measure the throughput performance of SpMM in TC-GNN when the dimension of node embedding increases for a roofline analysis. Because sparse matrix computation is largely limited by its memory access performance, which

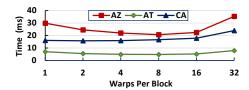


Figure 9: Performance Impact of Warps per Block.

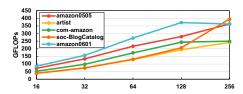


Figure 10: Analysis of TC-GNN kernel throughput when increasing the node embedding dimension from 16 to 256.

is quite different from the dense GEMM computation that is largely bounded by the computing performance. Figure 10 shows that the throughput of TC-GNN can scale proportionally with the growing number of node embedding dimensions. This also indicates that TC-GNN can effectively handle the graphs with high-dimensional node embeddings and well utilize GPU resources.

6 Related Work and Discussion

Other GPUs TC-GNN can easily generalize to other GPUs (e.g., A6000, H100, and RTX4090) with TCUs via recompilation (python setup.py install). TC-GNN also supports different TCU configurations (e.g., precision) by modifying $(BLK\ H,\ BLK\ W\ in\ TCGNN_conv/config.h)$ and four parameters (M, N, K, dataType) in wmma::fragment, then recompile. For future GPUs with more TCUs, our TC-GNN can also be adapted to accommodate such changes and maintain its performance advantage. There are two future GPU designs that we anticipate. The first direction is to place more TCUs per SM while keeping the total number of SMs unchanged. There will be more active warps per thread block (This is mainly because TCUs are operated by warps) and each warp will process fewer neighbors. The cost of decomposition and mapping can be offset by parallelism among more warps. The second direction is to place more SMs on GPUs while keeping TCUs per GPU unchanged. In this scenario, there will be more thread blocks and each thread block will process neighbors from fewer nodes. The cost can be offset by parallelism among more thread blocks.

Other GNN Frameworks Besides DGL and PyG, other single-GPU GNN frameworks like GNNAdvisor [57], GE-SpMM [21], and fuseGNN [7], tailor their own GNN layers manually with low-level GPU kernel optimizations. Unfortunately, these designs limit their kernel optimizations to CUDA cores, thus, missing the golden opportunities to exploit the full potential of widely deployed AI-tailored GPUs with TCUs.

Graph Partitioning/Reordering ROC [23] introduces a learning-based graph partitioning to reduce the data movement between CPU and GPU when processing large graphs. Rabbit Order [4] and Reverse Cuthill Mckee Algorithm [9] are focusing on *row reordering/clustering* to improve node/rowwise computation locality. Our sparse-graph translation (SGT) technique is orthogonal and complementary to these graph partitioning and reordering techniques since our SGT focuses on *column (neighbor) re-indexing* to improve neighbor-wise locality for TCU computation.

Distributed GNN Computation There are two major ways of scaling-up GNN computing: 1) *Distributed sampled graphs* [13, 35, 55, 60] (where graph nodes and their embeddings are on the same GPU): TC-GNN can be incorporated directly since all sampled graphs along with their node embeddings are presented at the same GPU. 2) *Distributed full-graph* [15, 23, 34, 58] (where graph nodes and their embeddings may be on different GPUs): TC-GNN needs to be modified slightly by incorporating inter-GPU communication techniques (e.g., Unified Virtual Memory [46] and NVSH-MEM [44]) to support the remote neighbor embedding access. We leave such exploration for our future work.

7 Conclusion

In this paper, we introduce TC-GNN, the first GNN acceleration framework on TCU of GPUs. We design a novel sparse graph translation technique to gracefully fit the sparse GNN workload on dense TCUs. Our TCU-tailored GPU kernel design maximizes the TCU performance gains for GNN computing through effective CUDA core and TCU collaboration and a set of memory/data flow optimizations. Our seamless integration with the PyTorch framework further facilitates end-to-end GNN computing with high programmability. Extensive experiments demonstrate the performance advantage of TC-GNN over the state-of-the-art frameworks. across diverse GNN models and datasets.

Furthermore, our TC-GNN design could also inspire potential TCU-like hardware features that can support (i) the dynamic shape of TCU input tiles and (ii) the dynamic structural sparsity of input tiles to yield higher performance benefits at the runtime. These proposed hardware features will help reduce more unnecessary computation in a more fine-grained and precise manner.

8 Acknowledgment

We would like to thank our shepherd, Asim Kadav, and the anonymous USENIX ATC reviewers. This work was supported in part by NSF-2124039 and CloudBank [36]. We also appreciate the generous help and support from NVIDIA Graduate Fellowship 2022-2023 for Yuke Wang and Amazon Faculty Award 2021-2022 for Yufei Ding.

References

- [1] Intel Math Kernel Library. Reference Manual. Intel Corporation. Santa Clara, USA.
- [2] AMD. All-new matrix core technology for hpc and ai. https://amd.com/en/technologies/cdna.
- [3] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen. LAPACK Users' Guide. Society for Industrial and Applied Mathematics, 1999.
- [4] Junya Arai, Hiroaki Shiokawa, Takeshi Yamamuro, Makoto Onizuka, and Sotetsu Iwamura. Rabbit order: Just-in-time parallel reordering for fast graph analysis. In 2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS), 2016.
- [5] Wieb Bosma, John Cannon, and Catherine Playoust. The Magma algebra system. I. The user language. J. Symbolic Comput., 1997. Computational algebra and number theory (London, 1993).
- [6] Hsinchun Chen, Xin Li, and Zan Huang. Link prediction approach to collaborative filtering. In *Proceedings* of the 5th ACM/IEEE-CS Joint Conference on Digital Libraries (JCDL), 2005.
- [7] Zhaodong Chen, Mingyu Yan, Maohua Zhu, Lei Deng, Guoqi Li, Shuangchen Li, and Yuan Xie. fusegnn: accelerating graph convolutional neural network training on gpgpu. In 2020 IEEE/ACM International Conference On Computer Aided Design (ICCAD), 2020.
- [8] De Cheng, Yihong Gong, Xiaojun Chang, Weiwei Shi, Alexander Hauptmann, and Nanning Zheng. Deep feature learning via structured graph laplacian embedding for person re-identification. Pattern Recognition, 2018.
- [9] Elizabeth Cuthill and James McKee. Reducing the bandwidth of sparse symmetric matrices. In *Proceedings of* the 1969 24th national conference, 1969.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- [11] Alberto Garcia Duran and Mathias Niepert. Learning graph representations with embedding propagation. In Advances in neural information processing systems (NeurIPS), 2017.
- [12] Boyuan Feng, Yuke Wang, Guoyang Chen, Weifeng Zhang, Yuan Xie, and Yufei Ding. Egemm-tc: Accelerating scientific computing tensor cores with extended

- precision. ACM SIGPLAN Symposium on Principles & Practice of Parallel Programming (PPoPP), 2021.
- [13] Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In ICLR Workshop on Representation Learning on Graphs and Manifolds (ICLR), 2019.
- [14] Matthias Fey and Jan E. Lenssen. Pytorch extension library of optimized scatter operations, 2019.
- [15] Swapnil Gandhi, Anand Padmanabha Iyer, Henry Xu, Theodoros Rekatsinas, Shivaram Venkataraman, Yuan Xie, Yufei Ding, Keval Vora, Ravi Netravali, Miryung Kim, et al. P3: Distributed deep graph learning at scale. In 15th USENIX Symposium on Operating Systems Design and Implementation (OSDI), 2021.
- [16] Jaume Gibert, Ernest Valveny, and Horst Bunke. Graph embedding in vector spaces by node attribute statistics. Pattern Recognition, 2012.
- [17] Joseph E Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. Powergraph: Distributed graph-parallel computation on natural graphs. In 10th USENIX Symposium on Operating Systems Design and Implementation (OSDI), 2012.
- [18] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In Proceedings of the 22nd ACM international conference on Knowledge discovery and data mining (SIGKDD), 2016.
- [19] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In Advances in neural information processing systems (NeurIPS), 2017.
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR), 2016.
- [21] Guyue Huang, Guohao Dai, Yu Wang, and Huazhong Yang. Ge-spmm: General-purpose sparse matrix-matrix multiplication on gpus for graph neural networks. In International Conference for High Performance Computing, Networking, Storage and Analysis (SC), 2020.
- [22] Zexi Huang, Arlei Silva, and Ambuj Singh. A broader picture of random-walk based graph embedding. In Proceedings of the 27th ACM International Conference on Knowledge Discovery & Data Mining (SIGKDD),
- [23] Zhihao Jia, Sina Lin, Mingyu Gao, Matei Zaharia, and Alex Aiken. Improving the accuracy, scalability, and

- performance of graph neural networks with roc. Proceedings of Machine Learning and Systems (MLSys), 2020.
- [24] Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. Indatacenter performance analysis of a tensor processing unit. In Proceedings of the 44th annual international symposium on computer architecture (ISCA), 2017.
- [25] Riesen Kaspar and Bunke Horst. Graph classification and clustering based on vector space embedding. World Scientific, 2010.
- [26] Kristian Kersting, Nils M. Kriege, Christopher Morris, Petra Mutzel, and Marion Neumann. Benchmark data sets for graph kernels, 2016.
- [27] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. International Conference on Learning Representations (ICLR), 2017.
- [28] Jérôme Kunegis and Andreas Lommatzsch. Learning spectral graph transformations for link prediction. In Proceedings of the 26th Annual International Conference on Machine Learning (ICML), 2009.
- [29] Süreyya Emre Kurt, Aravind Sukumaran-Rajam, Fabrice Rastello, and P. Sadayyapan. Efficient tiled sparse matrix multiplication through matrix signatures. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC), 2020.
- [30] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. http://snap. stanford.edu/data, June 2014.
- [31] Ang Li and Simon Su. Accelerating binarized neural networks via bit-tensor-cores in turing gpus. IEEE Transactions on Parallel and Distributed Systems (TPDS), 2020.
- [32] Dijun Luo, Chris Ding, Heng Huang, and Tao Li. Nonnegative laplacian embedding. In 2009 Ninth IEEE International Conference on Data Mining (ICDM), 2009.
- [33] Dijun Luo, Feiping Nie, Heng Huang, and Chris H Ding. Cauchy graph embedding. In Proceedings of the 28th International Conference on Machine Learning, 2011.
- [34] Lingxiao Ma, Zhi Yang, Youshan Miao, Jilong Xue, Ming Wu, Lidong Zhou, and Yafei Dai. Neugraph: parallel deep neural network computation on large graphs. In 2019 USENIX Annual Technical Conference (USENIX ATC), 2019.

- [35] Seung Won Min, Kun Wu, Sitao Huang, Mert Hidayetoğlu, Jinjun Xiong, Eiman Ebrahimi, Deming Chen, and Wen-mei Hwu. Pytorch-direct: Enabling gpu centric data access for very large graph neural network training with irregular accesses. arXiv preprint arXiv:2101.07956, 2021.
- [36] Michael Norman, Vince Kellen, Shava Smallen, Brian DeMeulle, Shawn Strande, Ed Lazowska, Naomi Alterman, Rob Fatland, Sarah Stone, Amanda Tan, et al. Cloudbank: Managed services to simplify cloud access for computer science research and education. In Practice and Experience in Advanced Research Computing. 2021.
- [37] Nvidia. Accelerating matrix multiplication with block sparse format and nvidia tensor https://developer.nvidia.com/blog/ cores. accelerating-matrix-multiplication-with\ -block-sparse-format-and-nvidia-tensor-cores/.
- [38] Nvidia. Cuda sparse matrix library (cusparse). https: //developer.nvidia.com/cusparse.
- [39] NVIDIA. Cuda template library for dense linear algebra at all levels and scales (cutlass).
- [40] Nvidia. Dense linear algebra on gpus. https:// developer.nvidia.com/cublas.
- [41] NVIDIA. Improved tensor core operations. https:// docs.nvidia.com/cuda/ampere-tuning-guide/ index.html#tensor-operations.
- [42] Nvidia. Nvidia blocked-sparse api. https: //docs.nvidia.com/cuda/cusparse/index. html#cusparse-generic-function-spmm).
- [43] Nvidia. Nvidia volta. https://en.wikipedia.org/ wiki/Volta_(microarchitecture).
- [44] Nvidia. Nvshmem communication library. https:// developer.nvidia.com/nvshmem.
- [45] NVIDIA. Tensorfloat-32 in the a100 gpu accelerates ai training, hpc up to 20x.
- Unified memory [46] NVIDIA. for cuda. https://developer.nvidia.com/blog/ unified-memory-cuda-beginners/.
- [47] Nvidia. Warp matrix multiply-accumulate (wmma). https://docs.nvidia.com/cuda/ cuda-c-programming-guide/index.html#wmma.
- [48] NVIDIA. Programming tensor cores in cuda 9. https://devblogs.nvidia.com/ programming-tensor-cores-cuda-9/, 2017.

- [49] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, highperformance deep learning library. In Advances in Neural Information Processing Systems (NeurIPS). 2019.
- [50] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In Proceedings of the 20th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2014.
- [51] Kiran K Thekumparampil, Chong Wang, Sewoong Oh, and Li-Jia Li. Attention-based graph neural network for semi-supervised learning. 2018.
- [52] Philippe Tillet, H. T. Kung, and David Cox. Triton: An intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages (MAPL)*, 2019.
- [53] Tomasz Tylenda, Ralitsa Angelova, and Srikanta Bedathur. Towards time-aware link prediction in evolving social networks. In *Proceedings of the 3rd workshop on social network mining and analysis*, 2009.
- [54] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [55] Minjie Wang, Lingfan Yu, Da Zheng, Quan Gan, Yu Gai, Zihao Ye, Mufei Li, Jinjing Zhou, Qi Huang, Chao Ma, Ziyue Huang, Qipeng Guo, Hao Zhang, Haibin Lin, Junbo Zhao, Jinyang Li, Alexander J Smola, and Zheng Zhang. Deep graph library: Towards efficient and scalable deep learning on graphs. ICLR Workshop on Representation Learning on Graphs and Manifolds, 2019.
- [56] Yuke Wang, Boyuan Feng, and Yufei Ding. Qgtc: accelerating quantized graph neural networks via gpu tensor core. In *Proceedings of the 27th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP)*, 2022.
- [57] Yuke Wang, Boyuan Feng, Gushu Li, Shuangchen Li, Lei Deng, Yuan Xie, and Yufei Ding. Gnnadvisor: An efficient runtime system for gnn acceleration on gpus. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2021.

- [58] Yuke Wang, Boyuan Feng, Zheng Wang, Tong Geng, Ang Li, Kevin Barker, and Yufei Ding. Mgg: Accelerating graph neural networks with fine-grained intra-kernel communication-computation pipelining on multi-gpu platforms. In USENIX Symposium on Operating Systems Design and Implementation (OSDI), 2023.
- [59] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In International Conference on Learning Representations (ICLR), 2019.
- [60] Jianbang Yang, Dahai Tang, Xiaoniu Song, Lei Wang, Qiang Yin, Rong Chen, Wenyuan Yu, and Jingren Zhou. Gnnlab: a factored system for sample-based gnn training over gpus. In *Proceedings of the Seventeenth European Conference on Computer Systems (EuroSys)*, 2022.
- [61] Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L. Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. In *Proceedings of the 32nd International* Conference on Neural Information Processing Systems (NeurIPS), 2018.
- [62] Orestis Zachariadis, Nitin Satpute, Juan Gómez-Luna, and Joaquín Olivares. Accelerating sparse matrix—matrix multiplication with gpu tensor cores. *Computers & Electrical Engineering*, 2020.
- [63] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in neural information processing systems* (*NeurIPS*), 31, 2018.

Artifact Appendix

TC-GNN is the first TCU-based GNN acceleration design on GPUs. At the input level, TC-GNN is equipped with a new sparse graph translation (SGT) technique that can effectively identify those non-zero tiles and condense non-zero elements from these tiles into fewer number of "dense" tiles. At the GPU kernel level, TC-GNN exploits the benefits of CUDA core and TCU collaboration. The major design idea is that the CUDA core, which is more powerful at fine-grained thread-level execution, would be a good candidate for managing memory-intensive data access. While TCU, which is more powerful in handling simple arithmetic operations (e.g., multiplication and addition), would be well-suited for computeintensive GEMM on dense tiles generated from SGT. At the framework level, TC-GNN is integrated with the popular Py-Torch framework to reduce extra learning efforts and improve user productivity and code portability.

- Code repository: **Github**² and **Zenodo**³.
- Hardware, OS & Compiler:
 - Intel Xeon Sliver 4110 CPU (8-core 16-threads) with 64GB host memory, NVIDIA RTX3090 GPU with 24 GB device memory.
 - Operating systems and versions: Ubuntu 16.04+.
 - Compilers and versions: NVCC-11.1+, GCC-7.5.0+ Libraries and versions: CUDA-11.1+, Pytorch-1.8.0, DGL-v0.6.0, PyG-1.6.3 Input datasets and versions: SNAP network datasets.

Step-1: Environment Setup

- 1.1a. [Method-1] Install via Docker (Recommended).

```
cd docker/
./launch.sh
```

- 1.1b. [Method-2] Install via Conda.

```
curl -0 https://repo.anaconda.com/archive/Anaconda3
        -2021.05-Linux-x86_64.sh
  bash Anaconda3-2019.03-Linux-x86_64.sh
  source ~/.bashrc
  conda create -n env_name python=3.6
  conda install pytorch torchvision torchaudio cudatoolkit
       =11.1 -c pytorch -c conda-forge
  conda install -c dglteam dgl-cuda11.0
  pip install torch requests tqdm
  pip install torch-scatter -f https://pytorch-geometric.
       com/whl/torch-1.8.0+cull1.html
9 pip install torch-sparse -f https://pytorch-geometric.
       com/whl/torch-1.8.0+cu111.html
10 pip install torch-cluster -f https://pytorch-geometric.
       com/whl/torch-1.8.0+cull1.html
pip install torch-spline-conv -f https://pytorch-
       geometric.com/whl/torch-1.8.0+cu111.html
  pip install torch-geometric
```

- 1.2. Install TC-GNN.

```
cd TCGNN_conv/
./0_build_tcgnn.sh
```

- 1.3. Download Datasets.

```
wget https://storage.googleapis.com/graph_dataset/tcgnn-
     ae-graphs.tar.gz
tar -zxvf tcgnn-ae-graphs.tar.gz
rm -rf tcgnn-ae-graphs.tar.gz
```

Step-2. Run Major Experiments.

- 2.1. TC-GNN model End-to-End.

```
1 ./0_run_tcgnn_model.sh
```

Results: 1_bench_gcn.csv and 1_bench_agnn.csv.

- 2.2. DGL baseline (Fig-6a).

```
cd dgl_baseline/
./0_run_dgl.sh
```

Results: Fig_6a_dgl_gcn.csv and Fig_6a_dgl_agnn.csv.

- 2.3. TC-GNN single kernel.

```
./0_run_tcgnn_single_kernel.sh
```

Results: 1_bench_gcn.csv and 1_bench_agnn.csv.

- 2.4. cuSPARSE-bSpMM Baseline (Fig-6c).

```
cd TCGNN-bSpmm/cusparse
./0_run_bSpMM.sh
```

Results: Fig_6c_cuSPARSE_bSpMM.csv.

- 2.5. Dense Tile Reduction (Fig-7).

```
python 3_cnt_TC_blk_SDDMM.py
python 3_cnt_TC_blk_SpMM.py
```

Results: 3_cnt_TC_blk_SDDMM.csv and 3_cnt_TC_blk_SDDMM.csv.

- 2.6. tSparse Baseline (Table-5, column-2).

```
cd TCGNN-tsparse/
./0_run_tSparse.sh
```

Result: Table_5_tSparse.csv.

- 2.7. Triton Baseline (Table-5, column-3).

```
cd TCGNN-trition/python/bench
./0_run_triton.sh
```

Result: 1_run_triton.csv.

²https://github.com/YukeWang96/TC-GNN_ATC23.git

³https://doi.org/10.5281/zenodo.7893174