

# Accelerating Sparse Approximate Matrix Multiplication on GPUs

Xiaoyan Liu<sup>1</sup>, Yi Liu<sup>1</sup>, Ming Dun<sup>2</sup>, Bohong Yin<sup>1</sup>, Hailong Yang<sup>1</sup>, Zhongzhi Luan<sup>1</sup>, and Depei Qian<sup>1</sup>

School of Computer Science and Engineering<sup>1</sup>

School of Cyber Science and Technology<sup>2</sup>

Beihang University<sup>1,2</sup>, Beijing, China, 100191

{liuxiaoyan,yi.liu,dunming0301,yinbohong,hailong.yang,zhongzhi.luan,depeiqli}@buaa.edu.cn

## ABSTRACT

Although the matrix multiplication plays a vital role in computational linear algebra, there are few efficient solutions for matrix multiplication of the near-sparse matrices. The Sparse Approximate Matrix Multiply (SpAMM) is one of the algorithms to fill the performance gap neglected by traditional optimizations for dense/sparse matrix multiplication. However, existing SpAMM algorithms fail to exploit the performance potential of GPUs for acceleration. In this paper, we present *cuSpAMM*, the first parallel SpAMM algorithm optimized for multiple GPUs. Several performance optimizations have been proposed, including algorithm re-design to adapt to the thread parallelism, blocking strategies for memory access optimization, and the acceleration with the tensor core. In addition, we scale *cuSpAMM* to run on multiple GPUs with an effective load balance scheme. We evaluate *cuSpAMM* on both synthesized and real-world datasets on multiple GPUs. The experiment results show that *cuSpAMM* achieves significant performance speedup compared to vendor optimized *cuBLAS* and *cuSPARSE* libraries.

## KEYWORDS

sparse approximate matrix multiplication, performance optimization, multiple GPUs

## 1 INTRODUCTION

Generally, the existing GEMM algorithms can be classified into dense and sparse algorithms according to the ratio of non-zero elements of the input matrices. Given a matrix  $A \in \mathbb{R}^{N \times N}$ , the number of non-zero elements is  $O(N^2)$  and  $O(N)$  for dense and sparse algorithms, respectively. However, in real applications, there are a large number of matrices in the middle ground between dense and sparse matrices, a.k.a. near-sparse matrices, whose non-zero elements are between  $O(N^2)$  and  $O(N)$ . Near-sparse matrices are widely used in the field of scientific computing, such as computational chemistry [41], quantum physics [26], electronic structure calculation [44].

The near-sparse matrices also exist in emerging domains such as deep neural networks. Especially in convolutional neural networks (CNNs), the feature and weight matrices participated in the calculation are near-sparse [13, 23] due to weight pruning [31] and activation functions [33]. For example, the activation function of Rectified Linear Unit (ReLU) can lead to more than 50% sparsity of the feature matrices on average [13]. In CNNs, the convolution operations between feature and weight matrices are transformed to GEMM using the *im2col* algorithm [2]. In such case, the matrices involved in the GEMM calculation are also near-sparse.

There is also a special class of matrices that are inherently near-sparse, the matrices with decay [5] (a.k.a. decay matrices), whose

elements (values) decrease rapidly from diagonal to sides. The elements can be ignored if they are small enough and the corresponding matrices become near sparse. Due to the unique properties, there are many researches focusing on decay matrix itself, such as the decay rate [20], the left inverse [21, 48], high-dimensional statistics [4], and numerical analysis [52]. In addition, the decay matrices often appear in widely used matrix operations such as matrix inverse [7, 20], matrix exponential [32], Jacobi matrices [47], and etc [6]. Moreover, decay matrices are commonly adopted in application domains such as quantum chemistry [5, 11] and quantum information theory [18, 19, 22, 46].

However, existing research works [12, 37, 51] for dense and sparse GEMM are hardly efficient when applied to near-sparse matrices. On the one hand, the researches for dense GEMM focus on reducing the computation complexity. For example, Strassen's algorithm [37] and Williams' algorithm [51] achieve  $O(N^{2.8})$  and  $O(N^{2.3727})$ , respectively. Whereas, the complexity reduction is hardly useful for eliminating redundant computation of near-sparse matrices on the zero elements. On the other hand, the researches for sparse GEMM propose various storage formats such as CSR [12] to store the sparse matrices compactly. However, the sparse formats can hardly benefit the near-sparse GEMM due to its non-sparse nature. Therefore, both the dense GEMM and the sparse GEMM have limited performance potential for near-sparse GEMM.

Fortunately, the approximation provides a good opportunity to boost the performance of near-sparse GEMM. For example, skipping the calculation of small enough elements of near-sparse matrices is a profitable way for performance acceleration. Based on such idea, Sparse Approximate Matrix Multiply (SpAMM) [14] has been proposed for accelerating the decay matrix multiplication. For matrices with exponential decay, existing research [3] has demonstrated the absolute error of SpAMM can be controlled reliably.

In the meanwhile, with wide adoption in a large number of fields, GPUs have been proven with excellent speedup for matrix operations [45]. Especially with the advent of tensor core units provided by NVIDIA GPUs, mixed-precision techniques have been exploited to further accelerate matrix operations [42]. Although there are few research works optimizing SpAMM computation on CPUs [3, 9, 10], to the best of our knowledge, there is no GPU implementation available for accelerating SpAMM computation, especially exploiting the architectural features such as tensor core and scaling to multiple GPUs. This motivates our work in this paper to re-design the SpAMM algorithm for better adaption to GPU architecture and propose corresponding optimizations to achieve superior performance compared to the state-of-the-art GEMM libraries. Specifically, this paper makes the following contributions:

- We propose *cuSpAMM*, a re-designed SpAMM algorithm tailored for GPU. Specifically, we adapt the calculation steps

and the data access patterns of the SpAMM algorithm to the memory hierarchy and thread organization of GPU with increased parallelism and reduced memory accesses.

- We propose several optimization schemes such as blocking strategies for the calculation kernels and utilization of tensor core for accelerating the calculation. In addition, we present a scaling method to extend *cuSpAMM* to multiple GPUs.
- We compare *cuSpAMM* with highly optimized vendor libraries such as *cuBLAS* [15] and *cuSPARSE* [17] on GPUs. In addition, we evaluate on two real datasets from electronic structure calculation (*ergo*) and convolutional neural network (VGG13), to demonstrate the performance of *cuSpAMM* on real-world applications.

The paper is organized as follows. In Section 2, we introduce the background of SpAMM algorithm and GPU optimizations. In Section 3, we present our re-designed SpAMM algorithm *cuSpAMM*, and corresponding optimizations for performance acceleration on multiple GPUs. Section 4 compares our *cuSpAMM* with the-state-of-art GEMM libraries and evaluates on two real-world datasets using *cuSpAMM*. Section 5 discusses the related works and Section 6 concludes this paper.

## 2 BACKGROUND

### 2.1 Decay matrix and SpAMM algorithm

A matrix is defined as the decay matrix when its elements decrease following the decay rate from the diagonal to the sides. The decay rate can be exponential or algebraical, formulated as  $|A[i][j]| < c\lambda^{|i-j|}$  and  $|A[i][j]| < c/(|i-j|^\lambda + 1)$  respectively, where  $A[i][j]$  is the index of the element in matrix  $A$ . The  $|i-j|$  is the separation and can be replaced by other index-based distance function of the matrix or physical distance such as  $|\vec{r}_i - \vec{r}_j|$  in non-synthetic cases [3]. By mathematical definition, the decay matrix is quite dense due to few zero elements. However, under certain conditions (e.g., elements less than the threshold), a number of elements in the decay matrix can be treated as zeros, which renders the matrix as near-sparse.

SpAMM is an approximate matrix multiplication method that can be used on decay matrices. The problem solved by SpAMM can be described as  $C = \alpha AB + \beta C$ , where  $\alpha$  and  $\beta$  are parameters, and  $A$ ,  $B$ , and  $C$  are the matrices with exponential decay or fast algebraical decay [14]. For convenience, the rest of the paper takes  $\alpha = 1$ ,  $\beta = 0$ , and the square matrices  $N \times N$ . Besides,  $\tau$  is a parameter for controlling the extent of approximation. The algorithm divides the input into quad-tree recursively, depicted in Equation 1. Then, the algorithm performs multiplication of sub-matrices recursively. The density of sub-matrices is measured by the Frobenius norm (F-norm), depicted in Equation 2. The Algorithm 1 shows the pseudo-code of SpAMM. The algorithm performs multiplication only if the product of norms from two sub-matrices is no smaller than the parameter  $\tau$  (line 8 and line 13).

$$A^t = \begin{pmatrix} A_{0,0}^{t+1} & A_{0,1}^{t+1} \\ A_{1,0}^{t+1} & A_{1,1}^{t+1} \end{pmatrix} \quad (1)$$

$$\|A^t\|_F = \sqrt{\sum_{i=0}^{N-1} \sum_{j=0}^{N-1} (A_{i,j}^t)^2} \quad (2)$$

---

### Algorithm 1 SpAMM algorithm

---

```

1: Input: Matrices  $A, B$ , parameter  $\tau$ 
2: Output: Multiplication result  $C$ 
3: if lowest level then
4:   return  $C = AB$ 
5: for  $i=0$  to  $1$  do
6:   for  $j=0$  to  $1$  do
7:     if  $\|A_{i,0}\|_F \|B_{0,j}\|_F \geq \tau$  then
8:        $T_0 = \text{SpAMM}(A_{i,0}, B_{0,j}, \tau)$ 
9:     else
10:       $T_0 = 0$ 
11:     if  $\|A_{i,1}\|_F \|B_{1,j}\|_F \geq \tau$  then
12:        $T_1 = \text{SpAMM}(A_{i,1}, B_{1,j}, \tau)$ 
13:     else
14:       $T_1 = 0$ 
15:      $C_{i,j} = T_0 + T_1$ 
16: return  $C$ 

```

---

### 2.2 GPU architecture and optimization

**2.2.1 GPU architecture.** The CUDA [16] programming paradigm provides a classic definition of GPU architecture, with the thread and memory organized hierarchically.

**Thread hierarchy** - The threads are organized at four levels with coarsened granularity including *thread*, *warp*, *block* and *grid*. One *warp* usually contains 32 *threads*, and the *threads* within a *warp* are implicitly synchronized in execution. The *warp* is the most basic unit for calculation execution and hardware scheduling. The *block* consists of multiple *threads*, and the *threads* within the same *block* can be synchronized. The *grid* consists of multiple *blocks*, and the *blocks* within the *grid* are executed in the SIMD fashion.

**Memory hierarchy** - The memory hierarchy can be divided into three levels, including *register level*, *shared memory level* and *global memory level*. Each *thread* has private registers, which is the fastest but also the most limited storage on GPU. The second fastest memory is the shared memory for each *block*. The threads within the same *block* can be synchronized through the shared memory. The *global memory level* consists of global memory and texture memory that hosts the data transferred from the CPU.

**2.2.2 GPU optimization.** We briefly summarize the commonly used optimization strategies for high-performance matrix multiplication on GPU.

**Architecture targeted optimizations** - The blocking strategies [8] partition the matrices and performs calculations across GPU memory hierarchy. Memory prefetching strategies [35] utilize guiding statements for writing memory, explicitly creating buffers and calling primitives, which overlaps the data movement with computation. Register optimization strategies [45] achieve better performance by reducing the active registers and minimizing access to high-latency memory such as global memory. Other optimization approaches that avoid bank conflict and tune hyper-parameters [50] (e.g., block size) are also useful for accelerating GEMM on GPU.

**Tensor core adoption** - The tensor core [38] introduced from Nvidia Pascal GPU has already been explored in many fields for further performance optimization such as linear algebra [28] and weather simulation [29] recently. In general, tensor core is a computation unit for Matrix-Multiply-Accumulate (MMA), formulated as  $D_{m \times k} = A_{m \times n} \times B_{n \times k} + C_{m \times k}$ , where the maximum number of matrix elements is 256. The matrix  $A, B$  must be in FP16 precision, while matrix  $C, D$  can be in FP16 or FP32 precision. The programming of tensor core is based on a special data structure

named *fragment*, which stores the computation data for the tensor core. The *threads* in each *warp* operate on the *fragments* to perform MMA calculation on tensor core.

### 3 METHODOLOGY AND IMPLEMENTATION

In this section, we will first give an overview of our re-designed SpAMM algorithm tailored for GPU, *cuSpAMM*. Then, we introduce the design of two important kernels in *cuSpAMM*, which adopts several optimization strategies as well as leverages tensor core to optimize the performance. In addition, we scale our implementation to multiple GPUs for processing larger matrices. Finally, we propose load balance and accuracy searching optimizations that further improve the performance of *cuSpAMM*.

For the convenience of illustration, we use the following notations. The input of the algorithm are matrices  $A, B \in \mathbb{R}^{N \times N}$  and  $\tau$ , where  $A, B$  are decay matrices, and  $\tau$  is the approximation threshold. The output matrix is  $C$ . For optimization, we divide the input matrix into sub-matrices with size of  $LoNum \times LoNum$ . We use  $BDIM = N/LoNum$  to denote the number of sub-matrices per row/column, where  $N$  is divisible by  $LoNum$ . The coordinates of the sub-matrix are represented by a square bracket. For example,  $A[i, j]$  represents the sub-matrix with the starting index of  $A[i \times LoNum][j \times LoNum]$ . To avoid incomplete division, the matrices are padded with zeros to satisfy the above assumption.

#### 3.1 Overview of *cuSpAMM*

Figure 1 shows the design overview of *cuSpAMM*. To eliminate the GPU-unfriendly recursion in original SpAMM as well as exploit higher parallelism, we re-design the algorithm composed of two kernels, *Get-norm kernel* and *Multiplication kernel*. The first kernel is responsible for calculating the F-norm of input matrices, and the second kernel decides whether to multiply the matrices depending on the F-norm results from the first kernel. The array used to record the F-norm values is *normmap*, where  $A\_normmap[i][j] = \|A[i, j]\|_F$  and  $B\_normmap[i][j] = \|B[i, j]\|_F$  for matrix  $A, B$  respectively. The re-designed *cuSpAMM* algorithm is equivalent to the original SpAMM algorithm, because they both perform calculation on the sub-matrices that satisfy the F-norm threshold ( $\tau$ ). In addition, we propose several optimizations for the above two kernels, and utilize tensor core for further performance acceleration. Specifically, we apply the blocking optimization to *cuSpAMM* across the following memory hierarchies. At *device level*, we partition matrices  $A, B, C, A\_normmap$  and  $B\_normmap$  in GPU global memory. At *block, warp, and thread level*, we partition the intermediate results in corresponding memory hierarchy. The details of blocking optimization are presented in the following sections.

#### 3.2 *Get-norm kernel*

The *get-norm* kernel is responsible for calculating the F-norm (based on Equation 2) results for all sub-matrices. Each block of *get-norm* kernel calculates the F-norm of one sub-matrix. Considering the computation characteristics of F-norm, we adopt the reduction algorithm for better parallelization, as shown in Figure 2(a). Firstly, each thread takes an element from the input matrix, calculates its square value, and stores the results into shared memory. Then, the thread block performs the reduction on the shared memory.

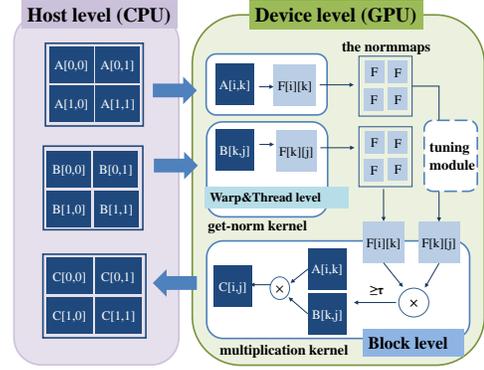


Figure 1: The overview of *cuSpAMM* algorithm.

To optimize reduction, we adopt sequential addressing instead of stride addressing to avoid the bank conflict on shared memory.

To further accelerate the performance with input matrices in FP16 precision, we use tensor core as MMA unit for reduction. Equation 3 and 4 show the reduction of  $m^2$  elements, where  $[1]_{m \times m}$  and  $[0]_{m \times m}$  represents a square matrix composed of 1 and 0 respectively,  $x_{11}, x_{12}, \dots, x_{mm}$  are the data waiting for summation. After two MMA operations, the reduction results are stored in matrix  $D'$ . This optimization can accelerate the reduction calculation compared to the traditional reduction on GPU [40]. Finally, the *thread 0* writes the result back to *normmap*.

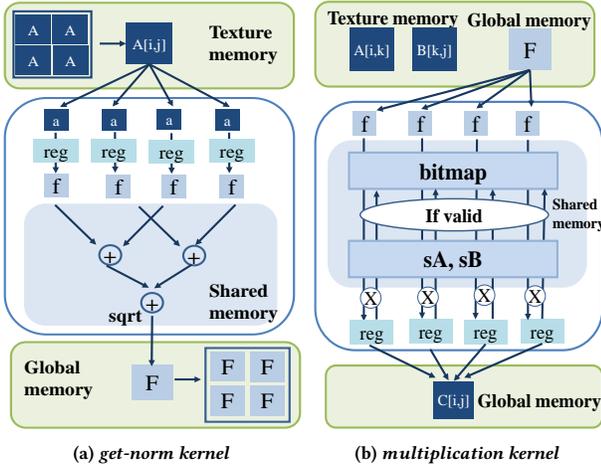
$$D = [1]_{m \times m} \times \begin{bmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mm} \end{bmatrix} + [0]_{m \times m} \quad (3)$$

$$D' = \begin{bmatrix} \sum_{i=1}^m x_{1i} & \cdots & \sum_{i=1}^m x_{1i} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^m x_{mi} & \cdots & \sum_{i=1}^m x_{mi} \end{bmatrix} \times [1]_{m \times m} + [0]_{m \times m} \quad (4)$$

Meanwhile, we apply additional optimizations to further boost the performance. Firstly, we increase the amount of data to be processed by each thread for coalescing the global memory access. We also use the vector operations such as *float2* to reduce the number of memory load instructions. Moreover, we perform loop unrolling on both algorithm level and warp level to reduces redundant jump and synchronization operations.

#### 3.3 *Multiplication kernel*

The *multiplication* kernel is responsible for performing the actual matrix multiplication depending on the F-norm results from *get-norm* kernel. Figure 2(b) shows the blocking strategy and execution flow for *multiplication* kernel. Each block has  $LoNum \times LoNum$  threads and is responsible for calculating one sub-matrix of matrix  $C$ . Supposing the block is responsible for  $C[i, j]$ , and  $C[i, j] = \sum A[i, k] \times B[k, j] \times bitmap[k]$ , where  $k$  ranges from 0 to  $(BDIM-1)$  and  $bitmap[k]$  is 1 or 0, indicating whether  $A[i, k]$  and  $B[k, j]$  satisfies the F-norm threshold, respectively. The *bitmap* is stored in shared memory. After that, all threads begin to go through the *bitmap*. If  $bitmap[k]$  is 1, the threads load the  $A[i, k]$  and  $B[k, j]$  into shared memory and then perform the dot product. We adopt double



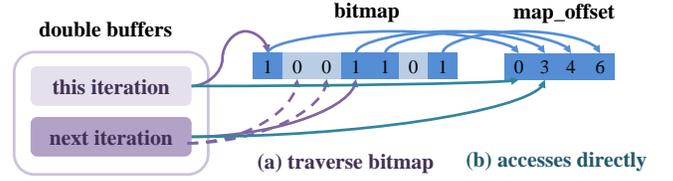
**Figure 2: The blocking optimizations and execution flow of *get-norm* kernel and *multiplication* kernel. To illustrate, assuming  $LoNum=2$ , bank number=2, and  $A$  is the decay matrix.**

buffering for hiding the memory access latency during the batched sub-matrix product.

However, as shown in Figure 3(a), it is inefficient to implement double buffering naively. This is because the thread needs to go through the *bitmap* to identify next valid sub-matrices for multiplication. The naive implementation introduces additional instructions (e.g., jump and comparison), which even leads to performance degradation. To address the above problem, we improve the double buffering technique as shown in Figure 3(b), which transforms the accessing of the valid flags in the *bitmap* from discontinuous to continuous for better locality. Although such an approach introduces additional calculations, it improves the efficiency of data prefetching with better locality, and thus accelerates the performance of *multiplication* kernel.

Algorithm 2 shows the optimized *multiplication* kernel. The threads identify the sub-matrices that requires actual multiplication and record corresponding indexes in *bitmap* in parallel (line 5~8). Specifically, threads calculate the F-norm condition using  $A\_normmap[i][k]$  and  $B\_normmap[k][j]$  and update  $bitmap[k]$  for each  $k$ . We use another array  $map\_offset$  to store the indexes of valid sub-matrices continuously (line 9~14). During each iteration (line 19~27), the first half of the block threads is responsible for matrix multiplication, and the second half is responsible for data prefetching. This strategy facilitates hiding the memory access latency by overlapping computation with data access on GPU. Besides, each thread calculates two elements of  $C$ , which can be stored in thread registers during dot product.

For input matrices in FP16 precision, we use the tensor core to further accelerate the matrix multiplication. Algorithm 3 shows the pseudo-code of *multiplication* kernel using tensor core for acceleration (with the same code in FP32 precision omitted). Each block has  $warpRow \times warpCol$  warps and each warp is responsible for the sub-matrix  $C[warpRow, warpCol]$ . The *fragment*  $a\_frag$  and



**Figure 3: To implement double buffering technique, (a) the naive way to traverse *bitmap* introduces discontinuous accesses and additional instructions, and (b) the optimized way uses additional *map\_offset* array to store the index of valid sub-matrices continuously.**

### Algorithm 2 *cuSpAMM*: multiplication kernel

```

1: Input: Matrix pointer *A, *B, *C, normmap pointer*A_normmap, *B_normmap, parameter  $\tau$ 
2: Shared memory: bitmap, map_offset sAW, sBW, sAR, sBR
3: for  $k=threadId$  to  $N/LoNum$  by  $blockDim.x$  do
4:    $norm\_mul = A\_normmap[i][k] \times B\_normmap[k][j]$ 
5:   if  $norm\_mul \geq \tau$  then
6:      $bitmap[i] = 1$ 
7:   else
8:      $bitmap[i] = 0$ 
9: for  $i=threadId$  to  $N/LoNum$  by  $blockDim.x$  do
10:  if  $bitmap[i] == 1$  then
11:     $t = 0$ 
12:    for  $j=0$  to  $i-1$  do
13:       $t = t + bitmap[j]$ 
14:     $map\_offset[i] = t$ 
15:   $\_syncthreads$ 
16:  reduce bitmap to get the amount of valid multiplication, save it in  $validNum$ 
17:  if  $validNum$  is not zero, fetch data in first block
18:   $\_syncthreads$ 
19:  for  $i=0$  to  $validNum-1$  do
20:     $b = map\_offset[i]$ 
21:     $\_syncthreads$ 
22:    Exchange pointer between read and write
23:    if the thread is in first half block then
24:      if  $i$  is less than  $validNum-1$ , let  $next = map\_offset[i+1]$  and perform prefetch
25:    else
26:      if the thread is in second half block then
27:        Calculate two values  $(c1, c2)$  of  $sAW \times sBW$ 
28:  write back  $c1, c2$  to matrix C

```

$b\_frag$  stores sub-matrices of  $A$  and  $B$ ,  $ab\_frag$  is the accumulator of intermediate results. The  $ab\_frag$  uses FP32 precision for obtaining better accuracy. The  $ab\_frag$  is initialized to 0. We also apply the double buffering optimization using *fragment*, which is similar to the implementation in FP32 precision.

### Algorithm 3 Multiplication kernel using tensor core

```

1: Input: Matrix pointer *A, *B, *C, normmap pointer*A_normmap, *B_normmap, parameter  $\tau$ 
2: shared memory: bitmap
3: fragment: a_frag, b_frag, ab_frag
4: .....
5: for  $i=0$  to  $validNum-1$  do
6:    $b = map\_offset[i]$ 
7:    $\_syncthreads$ 
8:   Exchange pointer between read and write
9:   if  $i$  is less than  $validNum-1$ , perform prefetch and let  $next = map\_offset[i+1]$ 
10:  load_matrix_sync(a_frag, A+warpRowOff+b×LoNum, N)
11:  load_matrix_sync(b_frag, B+warpColOff+b×LoNum×N, N)
12:  mma_sync(ab_frag, a_frag, b_frag, ab_frag)
13:  store_matrix_sync(C+warpRowOff×T+warpColOff, ab_frag)

```

### 3.4 Scaling to multiple GPUs

Modern servers are usually equipped with multiple GPUs (e.g., Nvidia DGX contains up to 16 GPUs). To leverage such performance potential, we extend the blocking optimizations to enable *cuSpAMM* scale to multiple GPUs on a single server. Note that our multiple GPU optimizations can be further integrated with distributed matrix multiplication optimizations such as CANNON [27] and SUMMA [49]. However, due to the time constraint, we focus on describing the multiple GPU optimizations on a single server, and leave the extension for distributed GPUs in future work.

Algorithm 4 presents the pseudo-code of scaling *cuSpAMM* to multiple GPUs. Supposing that there are  $M$  GPUs indexed from  $0$  to  $M-1$ . The calculation task is divided by row, and GPU  $i$  is responsible for the rows in the range of  $(i \times M/N, (i+1) \times M/N)$  of  $C$ . The data transfer is divided into  $P$  batches and implicitly managed by the use of UM [16] technique. We control the data transfer by ordered page faults. Firstly, several CUDA streams are created with each stream manipulating one GPU device. Then, the CPU transfers the whole matrix  $B$  to each GPU in batches, and each GPU obtains the *normmap* of  $B$  at the same time (line 4~6). After that, the CPU sends rows  $[i \times N/M, (i+1) \times N/M)$  of matrix  $A$  in batches to GPU  $i$ . When each GPU receives the corresponding rows of  $A$ , it invokes *get-norm* kernel and waits for the kernel to finish (line 9). After that, it invokes *multiplication* kernel for calculating the result (line 11). The batching approach is able to hide the data transfer latency as well as reduce the number of active blocks, which in turn mitigates the scheduling overhead.

---

#### Algorithm 4 Scaling to multiple GPUs

```

1: Input: matrix A, B, tau
2: Output: matrix C
3: create CUDA stream stream for devices
4: for i=0 to P do
5:   launch get-norm kernel for B
6:   synchronize at stream level
7: for i=0 to P do
8:   launch get-norm kernel for A
9:   synchronize at stream level
10:  launch multiplication kernel
11: synchronize at host level
12: output C

```

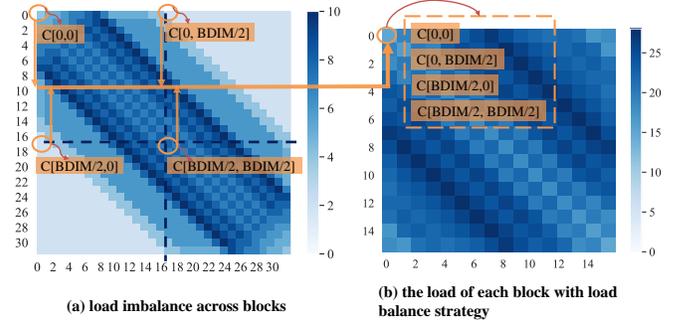
---

### 3.5 Additional optimizations

**3.5.1 Improving load balance.** The load imbalance could occur in *multiplication* kernel as shown in Figure 4(a). This is because each block calculates the *bitmap* dynamically to determine how many operations it needs to perform, which leads to block with less load staying idle and wasting resources. To measure the workload of each block, we propose the concept of valid multiplication  $v$ . For block responsible for calculating sub-matrix  $C[i, j]$ , its  $v$  equals to  $\sum_{i=0}^{BDIM} \text{bitmap}[i]$ . We organize the  $v$  values of all blocks into a matrix  $V$ , where  $V[i][j]$  is the  $v$  value of the sub-matrix  $C[i, j]$ . We observe that in matrix  $V$ , the closer to the diagonal, the greater the  $v$  is, which is determined by the property of decay matrix.

Based on the above observation, we propose the following load balance strategy. Each block of the *multiplication* kernel is responsible for the calculation of  $s$  (tunable parameter) sub-matrices with equal stride. For example, as shown in Figure 4(b), one block is responsible for sub-matrices  $C[0, 0]$ ,  $C[0, BDIM/2]$ ,  $C[BDIM/2, 0]$

and  $C[BDIM/2, BDIM/2]$  with  $s=2$ . The *multiplication* block can easily adopt the above strategy by adding a loop to change the index of its corresponding sub-matrices in order to achieve better load balance.



**Figure 4: The illustration of load balance strategy. The size of decay matrix is  $1024 \times 1024$ , the sub-matrix is  $32 \times 32$  and each multiplication block is responsible for  $16 \times 16$  sub-matrix.**

**3.5.2 Searching for customized accuracy.** For users using SpAMM algorithm to accelerate non-scientific applications such as deep neural networks (DNNs), adjusting  $\tau$  to control the extent of approximation is not intuitive. For example, the users of DNNs are more concerned about the accuracy of the entire network, other than the numerical accuracy of a single GEMM. In such case, we provide a tuning parameter *valid ratio*, formulated as  $\text{valid ratio} = \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} V[i, j] / BDIM^3$ , to control the actual multiplication of sub-matrices, which ensures that the sub-matrices with large and dense elements participate in calculation with higher priority. This tuning parameter can better adapt to the accuracy requirements of non-scientific applications.

Specifically, after the *normmaps* of  $A$  and  $B$  are obtained, we use a tuning kernel to calculate the average result (*ave*) of the norm products of all sub-matrices. The kernel then iterates to find the suitable value of norm  $\tau$  that satisfies the *valid ratio* given by the user. Binary search is applied during iterations with search space  $[0, k \times \text{ave}]$ , where the initial value of the norm is *ave*,  $k$  is the expansion coefficient, and the upper bound of the search space is dynamically extended. The initial value of  $k$  is one and will increase to  $k+1$  whenever the existing upper bound cannot satisfy the search demand. Besides, users can specify the number of iterations and tolerable error of *valid ratio* to balance the time cost and accuracy for searching. Since the searching algorithm is independent of the computation kernels, users can develop customized searching algorithms according to their application characteristics.

## 4 EVALUATION

### 4.1 Experiment setup

**Experiment platform** - The experiments are conducted on a CPU-GPU server, with two Intel Xeon E5-2680v4 processors and eight NVIDIA Volta V100 GPUs. Each GPU contains 32 GB memory. We use CUDA v10 [16] and nvcc compiler with -O3 option for our implementations.

**GEMM libraries** - We compare with *cuBLAS* that treats the decay matrix as dense matrix, and *cuSPARSE* that treats the decay matrix as sparse matrix through truncation. Note that, with truncation, the elements smaller than the threshold are treated as zero. For *cuBLAS*, we use *cublasSgemm* and *cublasHgemm* (with tensor core optimization) for matrix multiplication in FP32 and FP16 precision respectively. For *cuSPARSE*, we use *cusparseScsrgemm* for matrix multiplication in FP32 precision. However, *cuSPARSE* does not provide matrix multiplication in FP16 precision in CUDA v10.

**Evaluation criteria** - We use *cudaEvent* to record the execution time of the program in seconds, and the execution time ignores the overhead of input and output transfer (including format conversion) as well as warmup time. As for accuracy criteria, we use the F-norm of error matrix  $E_{n \times n}$  in Equation 5. For *cuSPARSE*, the decay matrix is truncated by setting the elements smaller than the threshold *TRUN* to zeros. The *nz ratio* represents the ratio of non-zero elements in the matrix. We use *valid ratio* to exhibit computation and memory patterns for *cuSpAMM*.

$$E_{n \times n} = A_{n \times n} B_{n \times n} - SpAMM(A_{n \times n}, B_{n \times n}, \tau) \quad (5)$$

**Synthesized matrix dataset** - For performance analysis, we synthesize the matrices with algebraical decay where  $a_{i,j} = b_{i,j} = 0.1/(|i - j|^{0.1} + 1)$ , and we control the *valid ratio* of the matrix indirectly by tuning the norm threshold  $\tau$ . Specifically, we use the tuning method in Section 3.5.2 to select the threshold and constrain the number of iterations to 20. The errors between actual and expected *valid ratio* are less than 1%.

**Table 1: The synthesized matrices with algebraical decay.**

<i>valid ratio</i> \ <i>N</i>	1,024	2,048	4,096	8,192	16,384	32,768
≈30%	1.434815	1.310666	1.195803	1.093354	0.997847	0.905539
≈25%	1.456555	1.330525	1.222981	1.113983	1.012852	0.919156
≈20%	1.489164	1.360312	1.250158	1.138739	1.03536	0.939582
≈15%	1.521774	1.40003	1.277335	1.171746	1.06537	0.966816
≈10%	1.586993	1.449676	1.322631	1.204753	1.110386	1.007668
≈5%	1.695691	1.548969	1.413222	1.28727	1.170407	1.062136

## 4.2 Comparison with GEMM libraries

In this section, we use synthesized matrices with algebraically decay listed in Table 1 for comparing with vendor optimized GEMM libraries including *cuBLAS* and *cuSPARSE*.

**4.2.1 Comparison with *cuBLAS*.** Table 2 presents the speedup of *cuSpAMM* on a single GPU compared to *cuBLAS*. The maximum speedup under each *valid ratio* is highlighted in red and blue for FP32 and FP16, respectively. When the *valid ratio* is 5%, the highest speedup is achieved with 13.4× (FP32) and 16.1× (FP16). Figure 5 shows the performance comparison when scaling our optimized *cuSpAMM* to multiple GPUs. It can be seen that *cuSpAMM* can accelerate matrix multiplication across all matrix sizes when the *valid ratio* is below a certain threshold. The reason is that when the *valid ratio* is below the threshold (25% for FP32 and 30% for FP16), our optimizations adopted in *cuSpAMM* can leverage the property of decay matrix multiplication for better parallelism compared to dense matrix multiplication adopted in *cuBLAS*. Moreover,

*cuSpAMM* achieves better performance speedup when scaling to multiple GPUs across all matrix sizes. For example, when *valid ratio* = 5%, *cuSpAMM* achieves the highest speedup of 51.4× with matrix size 4,096 in FP16 running on eight GPUs, compared to *cuBLAS*.

**Table 2: The speedup on matrices with algebraical decay on a single GPU. The first row under each *valid ratio* is the speedup for FP32, and the second line is for FP16.**

<i>valid ratio</i> \ <i>N</i>	1,024	2,048	4,096	8,192	16,384	32,768
≈30%	5.7	3.1	1.0	0.9	0.9	1.3
	4.3	5.2	2.3	1.1	1.3	1.6
≈25%	6.4	3.6	1.2	1.2	1.0	1.5
	4.6	5.8	2.9	1.6	1.5	1.8
≈20%	7.6	4.3	1.5	1.4	1.3	1.8
	5.2	6.9	3.7	2.0	1.9	2.2
≈15%	8.7	5.6	1.9	1.9	1.7	2.4
	5.8	8.7	4.7	2.7	2.5	2.9
≈10%	10.8	7.6	2.6	2.6	2.7	3.8
	6.5	11.4	6.8	3.5	3.8	4.5
≈5%	13.4	11.7	5.0	5.2	4.8	6.8
	7.6	16.1	11.9	7.0	6.5	7.6

**4.2.2 Comparison with *cuSPARSE*.** We choose the appropriate settings for  $\tau$  so that both implementations reach the same level of error. Since determining the appropriate settings ( $\tau$  and *TRUN*) is time consuming, we choose matrices matrix size of 1,024 and 8,192 for illustration (larger matrix causes out-of-memory error with *cuSPARSE* on a single GPU). From Table 3, it is clear that at the same error level, *cuSpAMM* is much faster than *cuSPARSE*, and the highest speedup reaches more than 601.0×. In addition, the speedup becomes larger as the *nz ratio* increases. Especially as the number of GPUs increases, the performance advantage of *cuSpAMM* becomes even larger (e.g., 3,985.1x speedup on eight GPUs). The phenomenon further demonstrates the incapability of *cuSPARSE* for handling near-sparse matrices with large non-zero ratio (e.g., more than 24%). Note that, the execution time of *cuSPARSE* does not include the time of format conversion. Thus, the performance speedup of our *cuSpAMM* will be higher than the results in Table 3 when compared in real-world application.

**Table 3: Performance comparison with *cuSPARSE*. We choose the matrix size of 1,024 (no.1) and 8,192 (no.2) with  $\tau$  and *TRUN* in *cuSpAMM* determined for the same level of error.**

<i>no.</i>	<i>nz ratio</i>	<i>valid ratio</i>	$\ E\ _F^*$	$\ E\ _F^\dagger$	speedup (1/2/4/8 GPUs)
1	52.13%	26.83%	1,020	996	232.3/379.5/586.8/875.8
	24.37%	6.70%	1,324	1,302	34.6/53.2/71.6/88.4
	10.91%	1.87%	1,400	1,387	11.0/14.6/22.9/21.4
2	59.59%	10.35%	38,173	37,090	589.9/1,171.4/2,127.4/3,985.1
	26.95%	0.73%	46,340	46,340	601.0/1,150.6/1,900.8/3,097.0
	2.12%	0.28%	46,340	46,340	71.2/130.9/220.8/336.9

\* the error of *cuSPARSE*

† the error of *cuSpAMM*

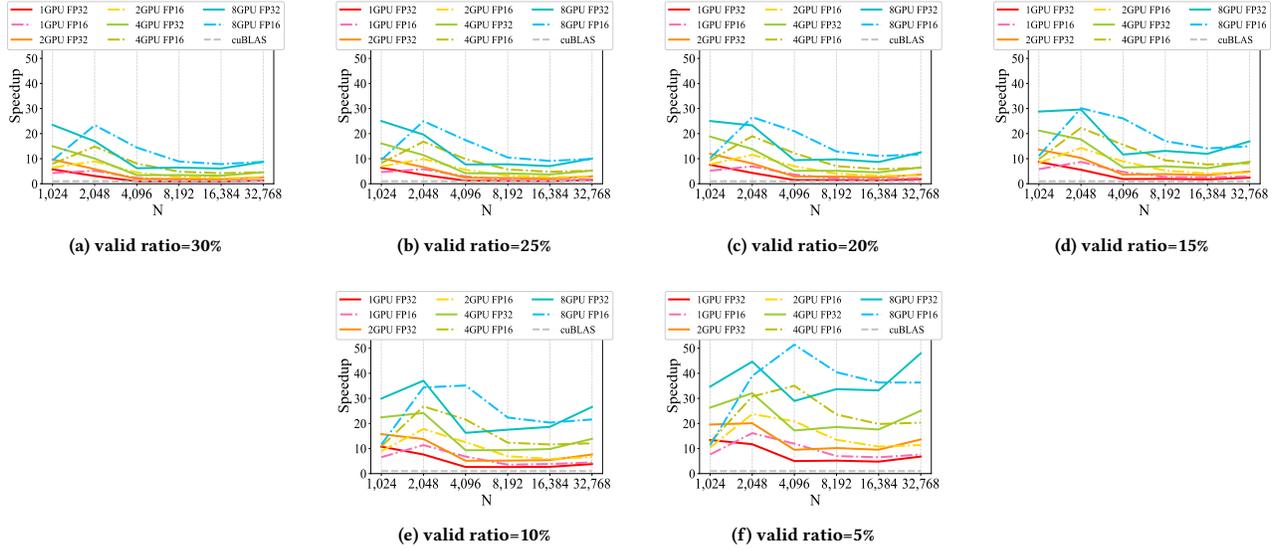


Figure 5: Performance comparison with *cuBLAS* on matrices with algebraical decay. We change the *valid ratio* from 30% to 5% with matrix size increasing from 1,024 to 32,768. In addition, we evaluate *cuSpAMM* scaling from one to eight GPUs.

### 4.3 Case study

We choose two applications widely used in scientific computing and deep neural network to further demonstrate the performance speedup. we only compare with *cuBLAS* in our case study since the experiment results in Section 4.2 indicate *cuBLAS* achieves better performance with near-sparse matrices compared to *cuSPARSE*.

**4.3.1 ergo application.** *ergo* [44] is an electronic structure computing program widely used in a range of scientific disciplines. We use *ergo* and the water cluster XYZ file [1] to derive the decay matrices directly. The program generates four decay matrices with exponential rate, and the size of each matrix is 13,656×13,656. We use *cuSpAMM* to calculate the power of these matrices, and we use parameter  $\tau$  to control the error ( $\|E_{n \times n}\|_F$ ) of the results.

Table 4 and Figure 6 present the F-norm of the matrices, the error of *cuSpAMM* with different  $\tau$ , and the performance speedup. *cuSpAMM* achieves increasing speedup when  $\tau$  becomes larger across all matrices. The performance speedup of *cuSpAMM* also scales when parallelizing on multiple GPUs. Especially for matrices with large F-norm ( $\|C\|_F > 1e^7$ ) and  $\tau=1e^{-2}$ , the average speedup ranges from  $3.0 \times$  to  $9.8 \times$  when scaling to multiple GPUs. In the meanwhile, the error introduced by *cuSpAMM* is much smaller than the data involved in the calculation ( $\|E\|_F/\|C\|_F < 8.9e^{-7}$ ). For the matrix with small F-norm such as the matrix no.1 and no.2, acceptable error ( $\|E\|_F/\|C\|_F < 1.6e^{-5}$  when  $\tau = 1e^{-4}$ ) can be achieved with average speedup of  $1.7 \times / 2.9 \times / 3.4 \times / 6.5 \times$  (1/2/4/8 GPUs). In the extreme case with no errors introduced (when  $\tau = 1e^{-10}$ ), *cuSpAMM* can still provide average speedup of  $1.3 \times / 1.5 \times / 2.3 \times / 4.0 \times$  across all matrices.

**4.3.2 VGG13 application.** We use VGG13 model on dataset MNIST. We use the *im2col* algorithm to convert the trained weights and input data into matrices. We use 80% of the dataset for training and

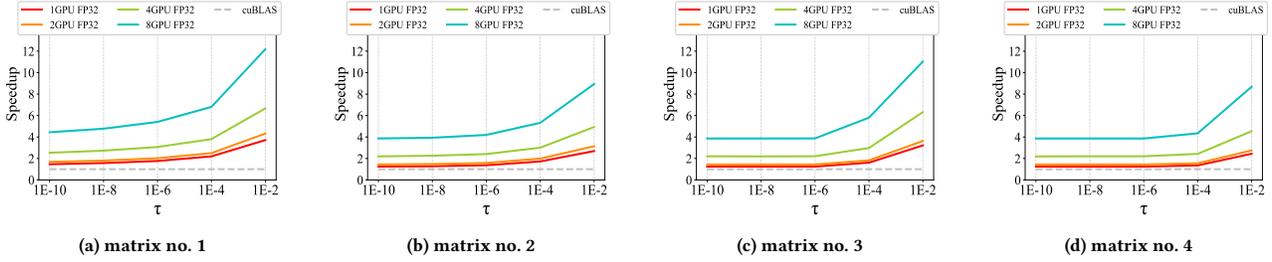
Table 4: The F-norm and error of the matrices from *ergo* application under different settings of  $\tau$  when using *cuSpAMM*.

matrix no.	$\ C\ _F$	$\tau$				
		$1e^{-10}$	$1e^{-8}$	$1e^{-6}$	$1e^{-4}$	$1e^{-2}$
1	755	0.0	$1e^{-06}$	$9e^{-05}$	0.01139	1.492293
2	10,406	0.0	$1e^{-06}$	$8.2e^{-05}$	0.013414	1.571806
3	3,169,858	0.0	0.0	$3.3e^{-05}$	0.021516	2.835374
4	17,171,990	0.0	0.0	$3e^{-06}$	0.013709	2.102697

the rest for validation. The size of the input figure is  $32 \times 32$ , and the number of channels is three. The *batch size* is set to 100 in both training and testing. The prediction accuracy of the original model is 96.6%.

Due to the time constraint, we only choose the two largest convolution layers *conv21* and *conv31* from VGG13 for detailed evaluation with *cuSpAMM*. Other convolution layers should expect similar tendency on performance speedup. After applying *im2col* operation, the scale of the matrix multiplication is  $128 \times 576 \times 25,600$  and  $256 \times 1,152 \times 6,400$  for these two layers, respectively. We apply *cuSpAMM* to accelerate the two layers and use the same dataset to test the performance and accuracy. Since the size of matrices is not large enough to occupy more than four GPUs, we only perform performance evaluation on four GPUs at the largest scale for VGG13.

Table 5 shows the evaluation results on both accuracy and performance when using *cuSpAMM*. In general, *cuSpAMM* is more effective for improving the performance of matrices multiplication in convolution neural network due to its insensitivity to matrix approximation. For both *conv21* and *conv31*, we can observe significant performance speedup under different settings of  $\tau$  with negligible accuracy loss (less than 1.1%). The performance speedup



**Figure 6: Performance comparison with *cuBLAS* on four matrices from *ergo* application. The four matrices exhibit exponential decay. In addition, we evaluate *cuSpAMM* scaling from one to eight GPUs.**

also scales well when increasing from one to four GPUs. Particularly for layer *conv31*, *cuSpAMM* achieves  $2.6\times$  to  $7.8\times$  performance speedup with the prediction accuracy unaffected. With 0.1% accuracy loss, *cuSpAMM* achieves  $2.7\times$  to  $9.0\times$  performance speedup. The highest performance speedup  $12.4\times$  is achieved when scaling to four GPUs with 1.1% accuracy loss.

**Table 5: The accuracy and speedup of *cuSpAMM* on VGG13 application. The *acc loss* measures the difference of prediction accuracy between the *cuSpAMM* optimized and original models, where negative results indicate accuracy loss.**

layer	valid ratio	acc loss	$\tau$	speedup (1/2/4 GPUs)
conv21	97.47%	0	0.1	2.8/5.0/8.4
	96.84%	0%	0.05	2.8/5.1/8.5
	85.00%	-0.1%	2.5	3.1/5.6/9.3
	82.90%	-0.1%	3.0	3.2/5.7/9.4
	63.41%	-0.9%	4.5	3.9/6.8/10.8
conv31	97.92%	0	2.5	2.6/4.8/7.6
	94.21%	0	3.5	2.6/4.8/7.8
	87.44%	-0.1%	4.5	2.7/5.0/8.1
	74.36%	-0.1%	5.5	3.1/5.6/9.0
	43.38%	-1.1%	7.5	4.8/8.1/12.4

## 5 RELATED WORK

### 5.1 Optimizing SpAMM

The SpAMM was first introduced by Challacombe *et al.* [14]. They proved that the time complexity of the algorithm is  $O(N \log N)$  at worst for matrices with a sufficient decay rate. Compared to the existing approximate methods, such as truncation and rank reduction, SpAMM requires much less floating-point operations. However, Challacombe *et al.* [14] only presented empirical experiments on error behavior, other than a detailed analysis on numerical error. In addition, they provided limited study on the decay matrices generated by the Heaviside matrix function from electronic structure domain [26]. Artemov *et al.* [3] studied the SpAMM algorithm with the matrices of exponential decay. They proved the absolute error behavior of SpAMM is  $\|E_{n \times n}\|_F = O(N^{(1/2)} \times O(\tau^{p/2}))$ , where  $p < 2$ .

Until recently, the performance optimization of SpAMM mainly focuses on CPUs. Bock *et al.* [9] reformed the SpAMM computation recursively and implemented the algorithm in parallel. In

addition, they used hashed (linkless) tree structure [24] to map the data and computation of sub-matrices. Moreover, the authors adopted hardware prefetching and optimized the performance using Intel Math Kernel Library (MKL) and AMD Core Math Library (ACML). Inspired from N-Body method [25], Bock *et al.* [10] and Artemov *et al.* [3] leveraged parallel programming models such as Charm++ [34] and Chunks/Tasks [43] to accelerate SpAMM, respectively. However, none of the existing works has optimized SpAMM on GPUs, which leaves the widely available GPU performance unexploited.

### 5.2 High performance GEMM

The high performance GEMM libraries on GPU have inspired our optimizations to SpAMM. *cuBLAS* and *cuSPARSE* are widely used high performance GEMM libraries heavily optimized by NVIDIA. Several optimization strategies have been proposed in *cuBLAS* and *cuSPARSE*. For example, the blocking strategies, instruction-level parallelism, parameter tuning, and tensor core acceleration have been successfully adopted for optimizing dense matrix multiplication in *cuBLAS*. For sparse matrix multiplication in *cuSPARSE*, various matrix storage formats such as COO, CSR and CSC are supported to optimize both sparse-sparse and sparse-dense matrix multiplication.

Moreover, Huang *et al.* [30] reviewed the blocking strategies of GEMM in *cuBLAS* and used tensor core to optimize the Strassen algorithm. Combined with optimizations such as software prefetching and parameter tuning, their implementation achieves  $1.11\times$  speedup on matrix multiplication compared to *cuBLAS*. Mukunoki *et al.* [39] evaluated the parallelized linear algebra kernels with multiple data precisions on GPUs. Ryoo *et al.* [45] summarized the general principles of matrix multiplication optimizations on GPU. Abdu *et al.* [40] used tensor core as matrix multiply-accumulate unit and proposed a chained reduction strategy. The above works have inspired the re-design and further optimization of SpAMM algorithm on GPU. Except the optimization of a single GEMM, batched GEMM [36] is a widely adopted technology that addresses small scale matrix multiplication and has already been integrated in vendor library such as *cuBLAS*. However, batched GEMM is not applicable to SpAMM due to its time-consuming reduction operations for accumulating the final result.

## 6 CONCLUSION

In this paper, we propose the first SpAMM algorithm *cuSpAMM*, tailored for acceleration on multiple GPUs. We re-design the SpAMM algorithm to parallelize the *get-norm* kernel and *multiplication* kernel. In addition, we apply several optimizations to improve the performance and accuracy of *cuSpAMM*, including blocking strategy, tensor core acceleration, double buffering, load balance and parameter searching. Moreover, we scale *cuSpAMM* to multiple GPUs for handling ever-increasing large near-sparse matrices. Our experiment results on both synthesized and real-world datasets show that *cuSpAMM* can achieve significant speedup compared to vendor optimized *cuBLAS* and *cuSPARSE* libraries.

## REFERENCES

- [1] 2020. ErgoSCF. (2020). <http://www.ergoscf.org/xyz/h2o.php>
- [2] Andrew Anderson, Aravind Vasudevan, Cormac Keane, and David Gregg. 2017. Low-memory gemm-based convolution algorithms for deep neural networks. *arXiv preprint arXiv:1709.03395* (2017).
- [3] Anton G. Artemov. 2019. Sparse approximate matrix multiplication in a fully recursive distributed task-based parallel framework. *CoRR* abs/1906.08148 (2019). [arXiv:1906.08148](http://arxiv.org/abs/1906.08148) <http://arxiv.org/abs/1906.08148>
- [4] Erlend Aune. 2012. Computation and modeling for high dimensional Gaussian distributions. Norges teknisk-naturvitenskapelige universitet, Fakultet for ...
- [5] Michele Benzi, Paola Boito, and Nader Razouk. 2013. Decay Properties of Spectral Projectors with Applications to Electronic Structure. *Siam Review* 55, 1 (2013), 3–64.
- [6] Michele Benzi and Gene H Golub. 1999. Bounds for the Entries of Matrix Functions with Applications to Preconditioning. *Bit Numerical Mathematics* 39, 3 (1999), 417–438.
- [7] Michele Benzi and Miroslav Tuma. 1999. Orderings for Factorized Sparse Approximate Inverse Preconditioners. *SIAM Journal on Scientific Computing* 21, 5 (1999), 1851–1868.
- [8] Christian H. Bischof and Philippe G. Lacroute. 1990. An adaptive blocking strategy for matrix factorizations. In *CONPAR 90 – VAPP IV*, Helmar Burkhart (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 210–221.
- [9] Nicolas Bock and Matt Challacombe. 2012. An Optimized Sparse Approximate Matrix Multiply for Matrices with Decay. *SIAM Journal on Scientific Computing* 35 (03 2012). <https://doi.org/10.1137/120870761>
- [10] Nicolas Bock, Matt Challacombe, and Laxmikant V. Kalé. 2016. Solvers for  $O(N)$  Electronic Structure in the Strong Scaling Limit. *SIAM J. Scientific Computing* 38, 1 (2016). <https://doi.org/10.1137/140974602>
- [11] D R Bowler and T Miyazaki. 2012.  $\mathcal{O}(N)$  methods in electronic structure calculations. *Reports on Progress in Physics* 75, 3 (feb 2012), 036503. <https://doi.org/10.1088/0034-4885/75/3/036503>
- [12] Aydin Buluç and John R. Gilbert. 2012. Parallel Sparse Matrix-Matrix Multiplication and Indexing: Implementation and Experiments. *SIAM J. Scientific Computing* 34, 4 (2012). <https://doi.org/10.1137/110848244>
- [13] Shijie Cao, Lingxiao Ma, Wencong Xiao, Chen Zhang, Yunxin Liu, Lintao Zhang, Lanshun Nie, and Zhi Yang. 2019. SeerNet: Predicting convolutional neural network feature-map sparsity through low-bit quantization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 11216–11225.
- [14] Matt Challacombe and Nicolas Bock. 2010. Fast Multiplication of Matrices with Decay. *CoRR* abs/1011.3534 (2010). [arXiv:1011.3534](http://arxiv.org/abs/1011.3534) <http://arxiv.org/abs/1011.3534>
- [15] NVIDIA Corporation. 2020. Nvidia cuBLAS. (2020). <https://developer.nvidia.com/cublas>
- [16] NVIDIA Corporation. 2020. NVIDIA CUDA C programming guide. (2020).
- [17] NVIDIA Corporation. 2020. Nvidia cuSPARSE. (2020). <https://docs.nvidia.com/cuda/cuspars>
- [18] M. Cramer and J. Eisert. 2006. Correlations, spectral gap, and entanglement in harmonic quantum systems on generic lattices. *New Journ. Phys* (2006).
- [19] M Cramer, Jens Eisert, Martin B Plenio, and J Dreissig. 2006. Entanglement-area law for general bosonic harmonic lattice systems. *Physical Review A* 73, 1 (2006), 012309.
- [20] Stephen Demko, William F Moss, and Philip W Smith. 1984. Decay rates for inverses of band matrices. *Math. Comp.* 43, 168 (1984), 491–499.
- [21] Victor Eijkhout and Ben Polman. 1988. Decay rates of inverses of banded  $M$ -matrices that are near to Toeplitz matrices. *Linear Algebra Appl.* 109 (1988), 247–277.
- [22] Jens Eisert, M Cramer, and Martin B Plenio. 2010. Area laws for the entanglement entropy - a review. *Reviews of Modern Physics* 82, 1 (2010), 277–306.
- [23] Trevor Gale, Matei Zaharia, Cliff Young, and Erich Elsen. 2020. Sparse GPU kernels for deep learning. *arXiv preprint arXiv:2006.10901* (2020).
- [24] Irene Gargantini. 1982. An effective way to represent quadrees. *Communications of The ACM* 25, 12 (1982), 905–910.
- [25] Alexander G. Gray and Andrew W. Moore. 2000. ‘N-Body’ Problems in Statistical Learning. In *NIPS*.
- [26] Stefan Grimme, Jens Antony, Stephan Ehrlich, and Helge Krieg. 2010. A consistent and accurate ab initio parametrization of density functional dispersion correction (DFT-D) for the 94 elements H-Pu. *Journal of Chemical Physics* 132, 15 (2010), 154104.
- [27] Anshul Gupta and Vipin Kumar. 1993. Scalability of parallel algorithms for matrix multiplication. In *1993 International Conference on Parallel Processing-ICPP’93*, Vol. 3. IEEE, 115–123.
- [28] Azzam Haidar, Stanimire Tomov, Jack Dongarra, and Nicholas J Higham. 2018. Harnessing GPU tensor cores for fast FP16 arithmetic to speed up mixed-precision iterative refinement solvers. In *SC18: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 603–613.
- [29] Sam Hatfield, Matthew Chantry, Peter Düben, and Tim Palmer. 2019. Accelerating high-resolution weather models with deep-learning hardware. In *Proceedings of the Platform for Advanced Scientific Computing Conference*. 1–11.
- [30] Jianyu Huang, Chenhan D. Yu, and Robert A. van de Geijn. 2018. Implementing Strassen’s Algorithm with CUTLASS on NVIDIA Volta GPUs. *ArXiv abs/1808.07984* (2018).
- [31] Yani Ioannou, Duncan Robertson, Roberto Cipolla, and Antonio Criminisi. 2017. Deep roots: Improving cnn efficiency with hierarchical filter groups. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1231–1240.
- [32] Arieh Iserles. 1999. How Large is the Exponential of a Banded Matrix? (1999).
- [33] Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. 2014. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866* (2014).
- [34] Laxmikant V. Kalé and Sanjeev Krishnan. 1993. CHARM++: a portable concurrent object oriented system based on C++. In *OOPSLA ’93*.
- [35] Nagesh B Lakshminarayana and Hyesoon Kim. 2014. Spare register aware prefetching for graph algorithms on GPUs. In *2014 IEEE 20th International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 614–625.
- [36] Xiuhong Li, Yun Liang, Shengen Yan, Liancheng Jia, and Yinghan Li. 2019. A coordinated tiling and batching framework for efficient GEMM on GPUs. In *Proceedings of the 24th Symposium on Principles and Practice of Parallel Programming*. 229–241.
- [37] Wei-Fen Lin, Steven K. Reinhardt, and Doug Burger. 2001. Reducing DRAM Latencies with an Integrated Memory Hierarchy Design. In *Proceedings of the Seventh International Symposium on High-Performance Computer Architecture (HPCA’01)*, Nuevo Leone, Mexico, January 20–24, 2001. IEEE Computer Society, 301–312. <https://doi.org/10.1109/HPCA.2001.903272>
- [38] Stefano Markidis, Steven Wei Der Chien, Erwin Laure, Ivy Bo Peng, and Jeffrey S Vetter. 2018. Nvidia tensor core programmability, performance & precision. In *2018 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. IEEE, 522–531.
- [39] Daichi Mukunoki and Takeshi Ogita. 2020. Performance and energy consumption of accurate and mixed-precision linear algebra kernels on GPUs. *J. Comput. Appl. Math.* 372 (2020), 112701.
- [40] Cristobal A Navarro, Roberto Carrasco, Ricardo J Barrientos, Javier A Riquelme, and Raimundo Vega. 2020. GPU Tensor Cores for fast Arithmetic Reductions. *arXiv: Distributed, Parallel, and Cluster Computing* (2020).
- [41] Mariana M. Odashima, Beatriz G. Prado, and E. Vernek. 2017. Pedagogical introduction to equilibrium Green’s functions: condensed-matter examples with numerical implementations. *Revista Brasileira de Ensino de Física* 39 (00 2017). [http://www.scielo.br/scielo.php?script=sci\\_arttext&pid=S1806-11172017000100402&nrm=iso](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S1806-11172017000100402&nrm=iso)
- [42] Roberto Olivaresamaya, Mark A Watson, Richard G Edgar, Leslie Vogt, Yihan Shao, and Alan Aspurguzik. 2010. Accelerating Correlated Quantum Chemistry Calculations Using Graphical Processing Units and a Mixed Precision Matrix Multiplication Library. *Journal of Chemical Theory and Computation* 6, 1 (2010), 135–144.
- [43] Emanuel H. Rubensson. 2017. Chunks and Tasks.
- [44] Elias Rudberg, Emanuel H Rubensson, Pawel Salek, and Anastasia Kruchinina. 2018. Ergo: An open-source program for linear-scaling electronic structure calculations. *SoftwareX* 7, 107–111.
- [45] Shane Ryoo, Christopher I Rodrigues, Sara S Baghsorkhi, Sam S Stone, David B Kirk, and Wenmei W Hwu. 2008. Optimization principles and application performance evaluation of a multithreaded GPU using CUDA. 73–82.
- [46] Norbert Schuch, J Ignacio Cirac, and Michael M Wolf. 2006. Quantum States on Harmonic Lattices. *Communications in Mathematical Physics* 267, 1 (2006), 65–92.
- [47] Barry Simon. 1982. Some Jacobi matrices with decaying potential and dense point spectrum. *Communications in Mathematical Physics* 87, 2 (1982), 253–258.
- [48] Romain Tessera. 2010. Left inverses of matrices with polynomial decay. *Journal of Functional Analysis* 259, 11 (2010), 2793–2813.

- [49] Robert A Van De Geijn and Jerrell Watts. 1997. SUMMA: Scalable universal matrix multiplication algorithm. *Concurrency: Practice and Experience* 9, 4 (1997), 255–274.
- [50] Vasily Volkov and James Demmel. 2008. Benchmarking GPUs to tune dense linear algebra. (2008), 31.
- [51] Virginia Vassilevska Williams. 2012. Multiplying Matrices Faster than Coppersmith-Winograd. In *Proceedings of the Forty-Fourth Annual ACM Symposium on Theory of Computing (STOC '12)*. Association for Computing Machinery, New York, NY, USA, 887–898. <https://doi.org/10.1145/2213977.2214056>
- [52] Qiang Ye. 2013. Error Bounds for the Lanczos Methods for Approximating Matrix Exponentials. *SIAM J. Numer. Anal.* 51, 1 (2013), 68–87.