

FINAL REPORT

SIKAI CHENG & PRISCILLA ZHANG

April 16, 2023

CONTENTS

1	Introduction	3
2	Literature Review	3
2.1	Different Matrix Format	3
2.2	Row Merging	4
2.3	Layered Graph Model	5
2.4	Sparse Matrix-Vector Multiplication	6
2.5	ESC Algorithm	6
3	Algorithm & Implementation	7
3.1	Setup	8
3.2	Expansion	9
3.3	Sorting	10
3.4	Contraction	11
4	Computational Result	13
5	Conclusion	14
6	Reference	16
7	Appendix	17
7.1	The Modified ESC Algorithm: Implementation, Test, and Sample Output	17
7.2	The <i>cusparseSpGEMM_compute</i> : Implementation, Test, and Sample Output	24
7.3	The <i>cusparse::multiply</i> : Implementation, Test, and Sample Output	28
7.4	The Modified ESC Algorithm: Demonstration	30

LIST OF FIGURES

Figure 1	Example: The Compressed Sparse Row (CSR) format.	4
Figure 2	Example: The ELLPACK (ELL) format.	4
Figure 3	Example: The Row Merging Operation.	5
Figure 4	Comparison: The performance of MKL, Cusp, Cusparse, and RMerge.	5
Figure 5	Example: A visualization of the Layered Graph Model.	6
Figure 6	Example: Represent a non-zero entry of C by the ensemble of paths.	6
Figure 7	Example: Expansion	7
Figure 8	Example: Sorting	7
Figure 9	Example: Compression	8
Figure 10	Performance Analysis: The ESC Algorithm	8
Figure 11	Example: Sparse matrix-matrix multiplication by using the modified ESC algorithm with all intermediate results.	8

Figure 12	Algorithm: Setup Phase	9
Figure 13	Code: Counting Kernel	10
Figure 14	Code: Setup Phase	10
Figure 15	Algorithm: Expansion Phase	11
Figure 16	Algorithm: Sorting Phase	11
Figure 17	Code: Bubble Sort Network	12
Figure 18	Algorithm: Contraction Phase	12
Figure 19	Code: Contraction Kernel	13
Figure 20	Code: Expansion Sorting and Contraction Kernel	13
Figure 21	Computational Results: The Modified ESC Algorithm	14
Figure 22	Sample Output	23
Figure 23	Sample Output	27
Figure 24	Sample Output	29

LIST OF TABLES

Table 1	Computational Results	14
---------	---------------------------------	----

1 INTRODUCTION

General Matrix-Matrix Multiplication (GEMM), the most fundamental and crucial linear operation, is widely used in all aspects of computational science. In most cases, GEMM is trivial and straightforward. It calculates the resulting matrix C ($\in \mathbb{Q}^{n \times k}$) based on two input matrices, A ($\in \mathbb{Q}^{n \times m}$) and B ($\in \mathbb{Q}^{m \times k}$), where $c_{ij} = \sum_l a_{il}b_{lj}$. However, in several particular scientific computing scenarios and combinatorial, such as multi-source breadth-first searching (Then et al., 2014), algebraic multigrid solvers (Bell et al., 2012), triangle counting (Davis, 2018), the shortest path finding (Chan, 2007), and colored intersecting (Kaplan et al., 2006), matrices A and B will contain ample zeros. These zeros and corresponding numerous computations of zero entries during the multiplication result in a high cost of both storage and operations. Therefore, SpGEMM (General Sparse Matrix-Matrix Multiplication) was formally proposed, quickly attracted the attention of many researchers, and became a standard computational block in these fields. In contrast to GEMM, SpGEMM requires exploiting the sparsity pattern in the matrix and obtaining a computational performance enhancement.

As opposed to traditional CPUs, which are tailored for optimizing the performance of sequential tasks in a single thread and minimizing the latency, modern GPUs contain massive parallelism, with dozens of multiprocessors each is capable of executing hundreds of hardware-scheduled threads and emphasizes total task throughput (Garland & Kirk, 2010). These properties equip GPUs with higher peak floating-point operations and memory bandwidth than CPUs, allowing them to impart tremendous performance in many high-performance computing applications, including SpGEMM definitely. However, realizing the full potential of GPUs for high-performance computing is not trivial. For multiple high-performance computing concerns, Dalton et al. (2015) underscore that the pervasive problem is that the existing methods, even those proven suitable for multicore CPUs, are not immediately applicable to the GPU. Meanwhile, they accentuate that effective use of GPUs requires substantial fine-grained parallelism at all stages of the computation. More specifically, for SpGEMM, critics have also argued that the superiority of current methods cannot escape from the dependence on the particular sparsity pattern of the matrix. Liu and Vinter (2014) claim that the methods either only work best for fairly regular sparse matrices (with most of the nonzero entries on the diagonal) or bring extra high memory overhead for matrices with some specific sparsity structures. Gremse et al. (2015) obtained a similar conclusion.

2 LITERATURE REVIEW

Therefore, the desire for a fast and general SpGEMM implementation has driven researchers to explore various matrix formats and computing algorithms. Some of the sections associated with this paper will be reviewed in the following.

2.1 Different Matrix Format

A natural idea, which is intuitively considered to have a direct effect on reducing the storage of sparse matrices, is to change the form of the matrices, in particular, to omit a large number of zeros from the matrix. Meanwhile, almost all previous explorations prove that it significantly impacts the algorithms' performance. Therefore, researchers have identified various operative formats for storing sparse matrices. One such format is the Coordinate list (COO), which involves storing the data in a tuple consisting of row index, column index, and entry value. The COO format is intuitive and does not depend on the number of zero values in a row. However, it generally performs worse in the SpGEMM problem compared to other formats (Li

et al., 2015). The Compressed Sparse Row (CSR) format is another commonly used format, Fig. 1, which consists of three arrays. The first array stores the entry values, the second array stores the corresponding column indices, and the third array records the index where a new row begins.

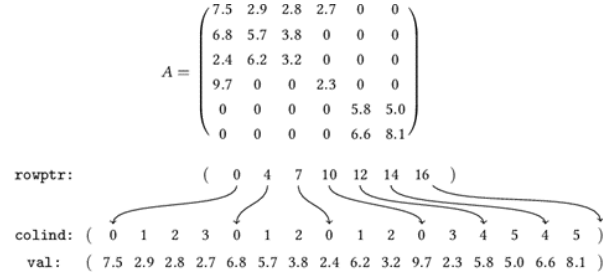


Figure 1: Example: The Compressed Sparse Row (CSR) format.

One significant observation is that the *rowptr* array is typically much smaller in size than the *colind* and *val* arrays, indicating that CSR format requires less storage space for matrix representation, particularly when dealing with large sparse matrices. Nonetheless, it should be noted that the CSR format's performance is highly dependent on the frequency of zero entries in a row. In addition to the commonly used COO and CSR formats, there exist other matrix formats, such as the ELLPACK (ELL) and hybrid (HYB) formats. Although these formats are not utilized in our study, they can be efficient options when the number of non-zero and zero elements are comparable.

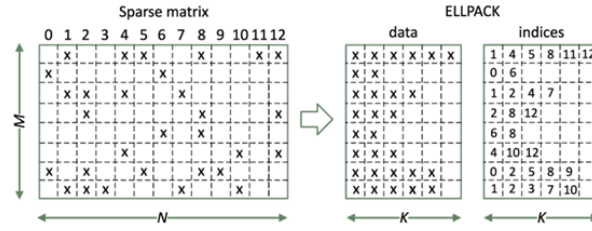


Figure 2: Example: The ELLPACK (ELL) format.

2.2 Row Merging

In addition to the ESC algorithm, alternative methodologies have been investigated in our literature review, among which is the iterative row merging algorithm using the CSR format (Gremse et al., 2015). In the context of matrix multiplication, namely $C = AB$, a possible approach involves translating the problem into $c = aB$, where c and a represent the corresponding rows of matrices C and A , respectively. Subsequently, c can be expressed as a linear combination of columns in B , with the weights determined by the entries in a . As a result, a row merging operation can be employed to compute each row in matrix C , as depicted in the Fig. 3. For the sake of simplicity, we assume the values of B and a to be 2 and 0.5, respectively. During each iteration of the row merging operation, the values with column indices are added with respect to the values of a , leading to the merging of two rows into a single row with a sorted column index. This process reduces the number of rows while potentially increasing the row size, depending on the frequency of identical indices present in the original rows.

In the case of computing $C = AA$, the authors assert that the row merging algorithm exhibits superior performance compared to other SpGEMM algorithms, such as MLK, Cusp, and Cuspars, in the Fig. 4, especially when working with large sparse matrices (Gremse et al., 2015). It is worth noting that the authors of the

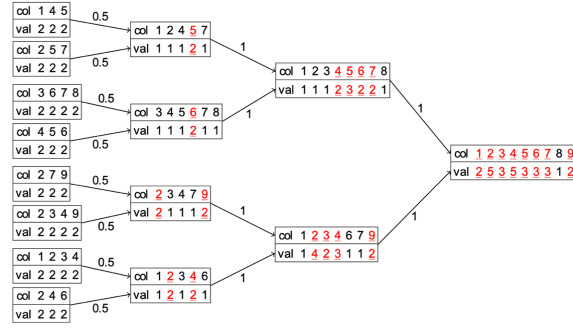


Figure 3: Example: The Row Merging Operation.

paper we are reproducing developed Cusp, a SpGEMM library. Subsequently, they made improvements to Cusp, which were published around the same time as the row merging algorithm. However, no additional comparison or analysis has been conducted to compare the performance of these two algorithms.

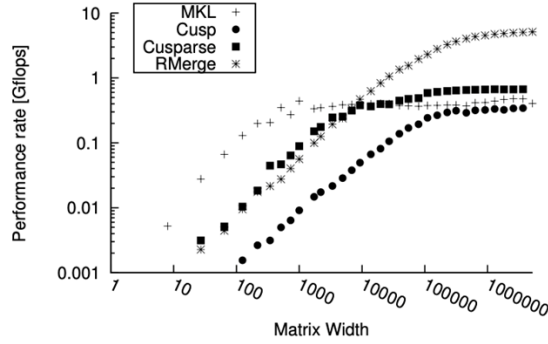


Figure 4: Comparison: The performance of MKL, Cusp, Cusp sparse, and RMerge.

2.3 Layered Graph Model

In the Expansion phase of the ESC algorithm, the authors utilize a geometric formulation of SpGEMM known as the Layered Graph Model (Cohen, 1996). This model employs a bipartite graph structure, where the vertices represent the rows and columns of the matrix. In the case of SpGEMM, where $C = AB$, the Layered Graph Model considers the rows of A as the base level vertices and the columns of A as the second level vertices. Notably, the number of columns in A is identical to the number of rows in B , which is crucial for leveraging this model in SpGEMM computations. The model then introduces a third layer of vertices representing the columns of B . In the Layered Graph Model, for each non-zero entry of matrix A with row index i , column index j , and entry value v , an edge is created between the i -th vertex of the base layer and the j -th vertex of the second layer with a weight of v . Similarly, for every nonzero entry of matrix B , a weighted edge is established between a second-layer vertex and a third-layer vertex. A visualization of such a model is illustrated in the Fig. 5 below as an example.

One of the primary benefits of the Layered Graph Model for SpGEMM problems is that it facilitates the formulation of each non-zero element in the resulting matrix C as a set of paths from a base-layer vertex to a top-layer vertex. More specifically, a non-zero entry of C with row index i , column index j , and entry value v , can be represented by the ensemble of paths from the i -th vertex of the base layer to the j -th vertex of the top layer, where the sum of the edge weights along all feasible paths is equivalent to v , Fig. 6.

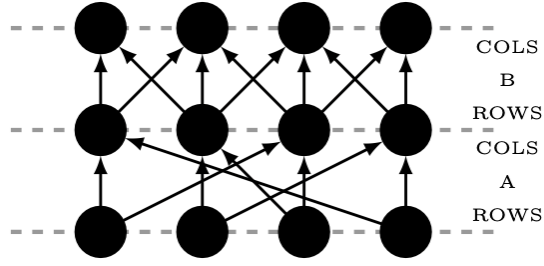


Figure 5: Example: A visualization of the Layered Graph Model.

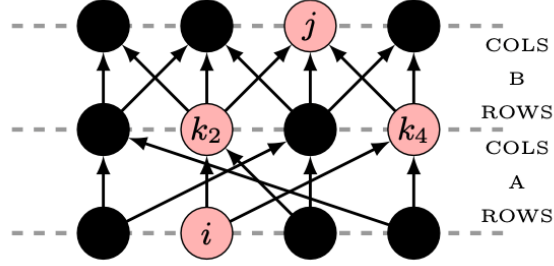


Figure 6: Example: Represent a non-zero entry of C by the ensemble of paths.

2.4 Sparse Matrix-Vector Multiplication

It is plausible that Sparse Matrix-Vector Multiplication (SpMV) is substantially similar to SpGEMM. Thus, the near-optimal algorithm for SpMV (Choi et al., 2010) gives the impression that a broad-spectrum and speedy algorithm for SpGEMM is close to being found. However, the study of Williams et al. (2009) points out that SpMV and SpGEMM are qualitatively different problems. Dalton et al. (2015) support the idea and explain it in detail. A common strategy for improving SpMV performance is to exploit a priori knowledge of the sparsity structure of the matrix in order to minimize expensive off-chip memory operations. This approach is considerably profitable when computing specific iterative algorithms that use the same sparse matrix multiple times or when the number of SpMV operations is sizeable. However, this approach is not entirely appropriate for SpGEMM since the interaction of sparsity from two matrices complicates the situation, and the SpGEMM is generally a fleeting operation, meaning that they are called at most once for a given set of matrices in most applications.

2.5 ESC Algorithm

The method proposed by Dalton et al., which our team wants to implement in the project, is strictly based on the Expansion, Sorting, and Compression (ESC) algorithms and consists of many finely designed bandwidths optimizing operations. The high-level structure of the ESC algorithm includes three steps:

- *Expansion* of AB into a intermediate coordinate format T .
- *Sorting* of T by row and column indices to form \hat{T} .
- *Compression* of \hat{T} by summing duplicate values for each matrix entry.

Dalton et al. (2015) analyzed the algorithm's performance, as shown in Fig. 10, and accordingly optimized each part, especially the sorting phase. In the Setup phase, the researchers recognized that the static assignment of computational units

$$\mathbf{A} = \begin{bmatrix} 5 & 10 & 0 \\ 15 & 0 & 20 \end{bmatrix} \quad \text{and} \quad \mathbf{B} = \begin{bmatrix} 25 & 0 & 30 \\ 0 & 35 & 40 \\ 45 & 0 & 50 \end{bmatrix} \quad \mathbf{T} = \begin{bmatrix} (0, 0, 125) \\ (0, 2, 150) \\ (0, 1, 350) \\ (0, 2, 400) \\ (1, 0, 375) \\ (1, 2, 450) \\ (1, 0, 900) \\ (1, 2, 1000) \end{bmatrix}$$

(a) Matrix Multiplication (b) Expansion

Figure 7: Example: Expansion

$$\hat{\mathbf{T}} = \begin{bmatrix} (0, 0, 125) \\ (0, 1, 350) \\ (0, 2, 150) \\ (0, 2, 400) \\ (1, 0, 375) \\ (1, 0, 900) \\ (1, 2, 450) \\ (1, 2, 1000) \end{bmatrix}$$

Figure 8: Example: Sorting

to process a fixed number of rows of matrix \mathbf{A} would cause an acute load imbalance, closely related to poor performance. So, they permute the rows of matrix \mathbf{A} to guarantee that the adjacent threads can handle rows with similar workloads. In the expansion phase, they adopt a formulation of SpGEMM as a layered graph model and employ the parallel bread-first-search to compute it. In the sorting and compression phase, they conclude that the parallel primitives used in the previous ESC algorithm, although general and avoiding excessive imbalance, force many data operations in global memory, inevitably leading to serious inefficiencies. Therefore, they focus on sorting and compressing within the GPU’s higher-bandwidth shared memory for increased efficiency.

3 ALGORITHM & IMPLEMENTATION

In this section, we will explain in detail the four critical parts of the algorithm, attach their implementation, and explain the code. Before starting the explanation, I want to declare the following three names:

- **The ESC algorithm:** The sparse matrix-matrix multiplication algorithm proposed in *EXPOSING FINE-GRAINED PARALLELISM IN ALGEBRAIC MULTI-GRID METHODS*. It is implemented by combining multiple fine-grained parallel primitives in the Thrust library, most of which are optimized at the global memory level.
- **The Modified ESC algorithm:** The modified version of the ESC algorithm proposed in *Optimizing Sparse Matrix-Matrix Multiplication for the GPU*. In addition, it makes full use of shared memory to improve the efficiency of sorting and contraction. This is the algorithm we want to reproduce.
- **The Referenced paper:** *Optimizing Sparse Matrix-Matrix Multiplication for the GPU*, the article we want to reproduce.

Two other points to note are that in the description of the referenced article, the authors repeatedly use functions from the open-source libraries Thrust and CUB. In our implementation, we avoid using existing functions except where explicitly

$$C = \begin{bmatrix} (0, 0, 125) \\ (0, 1, 350) \\ (0, 2, 550) \\ (1, 0, 1275) \\ (1, 2, 1450) \end{bmatrix}$$

Figure 9: Example: Compression

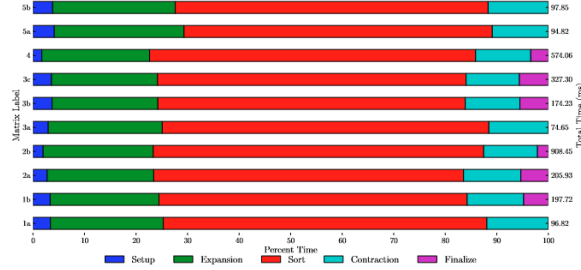


Figure 10: Performance Analysis: The ESC Algorithm

stated in the article. In addition, the referenced article considers that each part of expansion, sorting, and contraction may be applied separately to other algorithms, and therefore their implementations are independent of each other. In our implementation, I have integrated them into the **expansionSortingContractionKernel**, Fig. 20.

To plainly demonstrate the algorithm, we show the full intermediate process of multiplying two 7×7 sparse matrices as an example, Fig. 11.

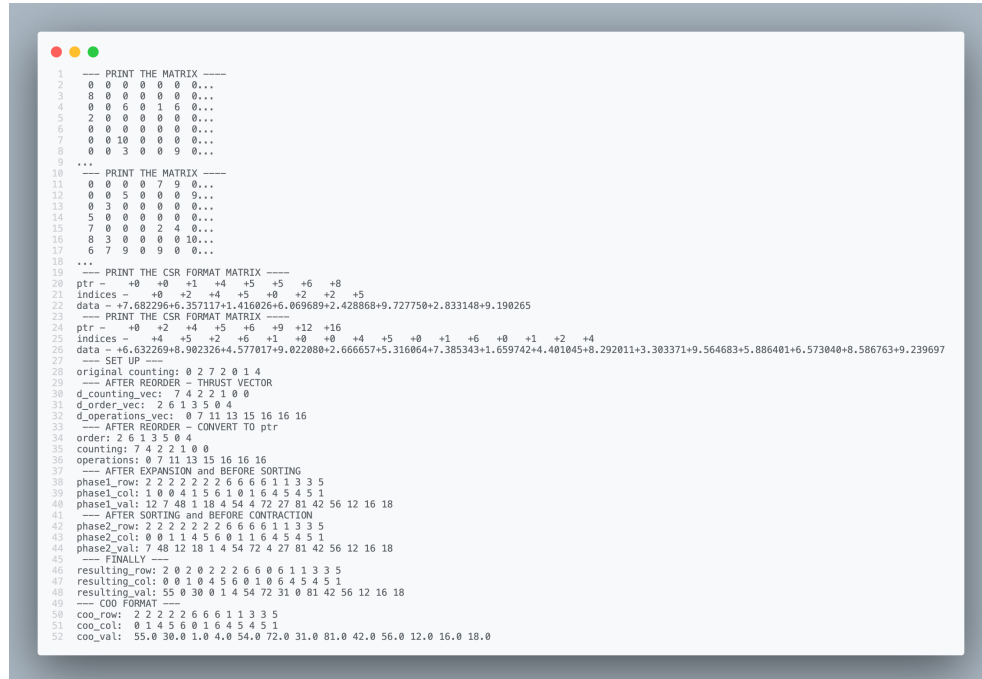


Figure 11: Example: Sparse matrix-matrix multiplication by using the modified ESC algorithm with all intermediate results.

3.1 Setup

The primary objective during the setup phase of the modified ESC algorithm is to determine the optimal allocation of threads and memory per row for the output

matrix C . The most naïve approach is to assign a uniform number of threads and memory per row. Nevertheless, this method is not suitable for sparse matrices since the number of operations needed in each row may vary significantly. Therefore, if a fixed number of threads and memory is allocated to all rows, many rows with high sparsity will waste a significant amount of resources, which is referred to as load imbalance.

Prior studies have examined various methods to mitigate load imbalances, but these methods have generally required considerable data movement between stages or some specific shape of nonzero entries in the sparse matrix. Therefore, in the referenced paper, the authors have proposed a novel approach based on reordering the A 's rows with the C 's rows based on the amount of computational work in the model. It is noteworthy that the minimum number of floating-point operations (FLOPs) associated with forming C_{row_i} is directly proportional to

$$\sum_{j \in \text{NNZ}(A_{\text{row}_i})} \text{nnz}(B_{\text{row}_j}) \quad (1)$$

By sorting the amount of computational work required per row in a non-decreasing order, a permutation matrix P can be derived for C , such that $PC = PAB$. This implies that the rows of matrix A can be processed in a permuted order. Utilizing this permutation technique makes it possible to group together rows that necessitate comparable computational work and allocate the appropriate amount of resources to each group. Nonetheless, the authors acknowledged that this approach has a disadvantage: the requirement of sorting the final output matrix back to its original order. However, the reassembly cost associated with this sorting process is typically low and can be deemed acceptable in most cases.

```

parameters:  $A, B$                                       $\{A \in \mathbb{R}^{m \times k} \text{ and } B \in \mathbb{R}^{k \times n}\}$ 
return:  $P$                                               {reordering vector}

 $F_i = 0$  for  $i = 1$  to  $m$ 
foreach row  $i$  in  $C$  do
    for  $j \in \text{NNZ}(A_{\text{row}_i})$ 
         $F_i \leftarrow F_i + \text{nnz}(B_{\text{row}_j})$            {gather  $B$  row lengths based on  $A$  column indices}
 $P \leftarrow \text{sort}(F)$                                 {set  $P$  to permutation of  $F$  in non-decreasing order}

```

Figure 12: Algorithm: Setup Phase

The codes used in the setup phase are shown in Fig. 13 & 14. First, I implemented a **countingKernel** to calculate the number of operations needed for each row of C , where one thread corresponds to one row. Then, **thrust::stable_sort_by_key** (line 15) is used to sort the number of operations obtained in the previous step. And **thrust::exclusive_scan** (line 17) is used to calculate the prefix sum of the sorted array, which will help to allocate a reasonable amount of memory for each row in the subsequent operations. The results of the setup phase are rendered in Fig. 11 line 27-36.

3.2 Expansion

The essential aim of the expansion phase is to perform computations on the rows of matrix B to produce an intermediate buffer, denoted as \hat{C} , which is an expanded and scaled version of B . A representation of such an expansion can be seen in the accompanying figure, Fig. 15, where B_{row_k} is scaled by $A_{i,k}$. This scaled expansion enables the computation of \hat{C} , which is similar to the COO format and consists of row indices \hat{I} , column indices \hat{J} , and values $\hat{V} = A_{i,k} * B_{k,j}$. In the literature review section, we introduced the layered-graph model, which was utilized by the authors during the expansion phase. An example of this model is presented in Fig. 5. As discussed previously, this model effectively demonstrates the formation of every non-zero entry of the final output matrix C as a group of paths from the initial

```

1  __global__
2  void countingKernel(int n, int* d_counting, int* A_ptr, int* A_ind, int* B_ptr) {
3      int idx = blockIdx.x * blockDim.x + threadIdx.x;
4      __syncthreads();
5      if (idx < n) {
6          int start = A_ptr[idx];
7          int end = A_ptr[idx + 1];
8          for (int i = start; i < end; i++) {
9              int j = A_ind[i];
10             d_counting[idx] = d_counting[idx] + B_ptr[j + 1] - B_ptr[j];
11         }
12     }
13     __syncthreads();
14 }

```

Figure 13: Code: Counting Kernel

```

1  int* counting = (int*)malloc(A_row * sizeof(int));
2  int* d_counting;
3  cudaMalloc(&d_counting, A_row * sizeof(int));
4  cudaMemset(d_counting, 0, A_row * sizeof(int));
5  int set_up_n_blocks = (A_row + N_THREADS_PER_BLOCK - 1) / N_THREADS_PER_BLOCK;
6
7  countingKernel <<< set_up_n_blocks, N_THREADS_PER_BLOCK >>> (A_row, d_counting, d_A_ptr, d_A_ind, d_B_ptr);
8  cudaDeviceSynchronize();
9  cudaMemcpy(counting, d_counting, A_row * sizeof(int), cudaMemcpyDeviceToHost);
10
11 thrust::device_ptr<int> d_counting_ptr(d_counting);
12 thrust::device_vector<int> d_counting_vec(d_counting_ptr, d_counting_ptr + A_row);
13 thrust::device_vector<int> d_order_vec(A_row);
14 thrust::sequence(d_order_vec.begin(), d_order_vec.end());
15 thrust::stable_sort_by_key(d_counting_vec.begin(), d_counting_vec.end(), d_order_vec.begin(), thrust::greater<int>());
16 thrust::device_vector<int> d_operations_vec(A_row + 1);
17 thrust::exclusive_scan(d_counting_vec.begin(), d_counting_vec.end(), d_operations_vec.begin());
18 int tot_operations = d_counting_vec.back() + d_operations_vec[A_row - 1];
19 d_operations_vec[A_row] = tot_operations;

```

Figure 14: Code: Setup Phase

vertex to the final vertex. The vector \hat{V} denotes the value of each path, with each path representing a scaled expansion of an entry of B.

The code implements the expansion phase is shown in Fig. 20 line 10-35. We use one block to expand one row of A and one thread in the block to compute one non-zero entry in the row. The obtained arrays are stored in the shared memory. Then, the results of the expansion phase are presented in Fig. 11 line 37-40.

3.3 Sorting

As the outcome of the Expansion phase, items with the same row index are placed in adjacent positions. However in the segment of data corresponding to the same row index, their column indexes may be duplicated, and they are not arranged in the order from the smallest to the largest. In order to get the COO form of the matrix multiplication result, we need to arrange the column index of each row in the sorting stage.

The algorithm is straightforward, as shown in Fig. 16; it just applies **key_value_sort** to the data segments corresponding to each row, but the implementation is not trivial. The referenced paper underscores that in implementing the ESC algorithm, the Thrust sorting algorithms allocate and free large amounts of temporary global memory each time they are invoked, which represents a non-trivial cost. The inef-

```

parameters:  $A, B$ 
return:  $\hat{C} = (\hat{I}, \hat{J}, \hat{V})$ 

foreach row  $i$  in  $C$  do
  for  $k \in \text{NNZ}(A_{\text{row}_i})$                                 {Note  $A_{i,k}$  is stored in shared memory for reuse}
    for  $j \in \text{NNZ}(B_{\text{row}_k})$ 
       $\hat{I} \leftarrow [\hat{I}, i]$                                 {implicit row index}
       $\hat{J} \leftarrow [\hat{J}, j]$                                 {append column index}
       $\hat{V} \leftarrow [\hat{V}, A_{i,k} * B_{k,j}]$                     {append value}

```

Figure 15: Algorithm: Expansion Phase

efficient sorting performance is a potential bottleneck of the algorithm. Besides, they analyzed some faster global sorting algorithms and exploited the knowledge about the range of input values. However, the results highlight the limited gains achieved by focusing on improving global sorting methods. Therefore, sorting within the GPU's higher-bandwidth shared memory is concentrated on in the modified ESC algorithm to enhance the sorting performance.

The modified ESC algorithm sorts the regions corresponding to each matrix row in parallel within shared memory, which is a feasible and efficient method because: first, the rows of the matrix are independent of each other. Processing one row does not affect the other rows; the global sorting operations over millions of entries are replaced by numerous operations over possibly tens to thousands of entries; sorting the intermediate entries using shared memory reduces global memory operations. In the described implementation, the referenced paper further considered the possible inefficiency of attempting to sort a highly varying workload using a static number of threads. To address this problem, we proportionally scale the number of threads per row of C with the maximum number of entries produced during the expansion. Specifically, when $\text{nnz}(C)$ is less than 32, we will use one thread and bubble sorting network to sort this row. And when $\text{nnz}(C)$ is greater than 32 and less than 768, we will use a block and radix sorting algorithm to sort this row.

```

parameters:  $n$ : number of columns in  $B$ 
              $\hat{C} = (\hat{I}, \hat{J}, \hat{V})$ : column indices,  $\hat{J}$ , reside in shared memory
return:  $\hat{C}, P$                                 { $\hat{J}$  sorted row-wise and permutation vector  $P$ }

foreach row  $i$  in  $C$  do
   $J, V \leftarrow \text{extract}_i \hat{J}, \hat{V}$                                 {extract entries where  $\hat{I} \equiv i$ }
   $m \leftarrow \text{nnz}(J)$                                 { $m$  is the number of expanded entries}
   $J, P \leftarrow \text{key\_value\_sort}(J, [0, m])$                     {key-value sort}

```

Figure 16: Algorithm: Sorting Phase

The code responsible for the sorting phase is shown in Fig. 17 & 20 line 39-50. The sorting operations are applied on the shared arrays obtained from the expansion phase. When $\text{nnz}(C_{\text{row}_i}) \leq 32$, we use the **bubbleSortNetwork** within a thread to sort data segments corresponding to row_i . Otherwise, we use the **cub::BlockRadixSort** within a block to do it. The results of the sorting phase are presented in Fig. 11 line 41-44.

3.4 Contraction

In the contraction phase, we contract the values associated with duplicate column indices in \hat{C} using pairwise addition. In contrast to the predictable nature of the total work required to construct \hat{C} in the expansion phase, the number of duplicates to form any row of C is challenging to compute and often varies significantly between rows in \hat{C} . To address the inefficiency caused by the irregularity of work required, the ESC algorithms try to avoid imbalance by employing the **reduce_by_key** function in Thrust. However, the modified ESC algorithm focuses more on the possible optimizations associated with utilizing shared memory storage following the

```

1  __device__
2  void bubbleSortSwap(int* col_ind, int i, int j) {
3      int temp = col_ind[i];
4      col_ind[i] = col_ind[j];
5      col_ind[j] = temp;
6  }
7
8  __device__
9  void bubbleSortSwap(float* ety_val, int i, int j) {
10     float temp = ety_val[i];
11     ety_val[i] = ety_val[j];
12     ety_val[j] = temp;
13 }
14
15 __device__
16 void bubbleSortNetwork(int* col_ind, float* ety_val, int n) {
17     for (int i = 0; i < n; i++) {
18         for (int j = i + 1; j < n; j++) {
19             if (col_ind[i] > col_ind[j]) {
20                 bubbleSortSwap(col_ind, i, j);
21                 bubbleSortSwap(ety_val, i, j);
22             }
23         }
24     }
25 }

```

Figure 17: Code: Bubble Sort Network

sorting results stored in the shared memory. The pseudo-code of the contraction algorithm is presented in Fig. 18.

parameters: $\hat{C} = (\hat{I}, \hat{J}, \hat{V})$, P (P permutation which sorts \hat{C} row-wise)
return: C

foreach row i in C do

1

$v \leftarrow 0$
 $J, V \leftarrow \text{extract}_i \hat{J}, \hat{V}$
for $j = 0, \dots, \text{nnz}(\hat{C}_{\text{row}_i})$

{initialize output value}
{extract entries where $\hat{I} \equiv i$ }
{Segmented scan by keys, J , over values in V }

2

$v \leftarrow v + V[P_j]$
if $J[j] \neq J[j+1]$

{reduce consecutive values}
{ $J[j+1]$ marks beginning of new nonzero entry}

3

$C_{i,J[j]} \leftarrow v$

{Store accumulated value to output row}

4

$v \leftarrow 0$

{re-initialize output value}

Figure 18: Algorithm: Contraction Phase

The code used to implement the contraction phase is rendered in Fig. 19 & 20 line 52. We apply **contractionOperation** to the shared arrays obtained from the sorting phase. Specifically, we assign one thread to one element and compare the element with the adjacent element. The results of the contraction phases are presented in Fig. 11 line 45-52.

```

1  __device__
2  void contraction0Operation(int* col_ind, float* ety_val, int n) {
3      int tid = threadIdx.x;
4      if (tid < n) {
5          if (tid == 0 || col_ind[tid] != col_ind[tid - 1]) {
6              float res = ety_val[tid];
7              for (int j = tid + 1; j < n && col_ind[j] == col_ind[tid]; j++) {
8                  res += ety_val[j];
9                  col_ind[j] = 0;
10                 ety_val[j] = 0;
11             }
12             ety_val[tid] = res;
13         }
14     }
15 }

```

Figure 19: Code: Contraction Kernel

```

1  __global__
2  void expansionSortingContractionKernel(int* order, int* counting, int* operations,
3
4      int* A_ptr, int* A_ind, float* A_val, int A_row, int A_col, int A_mz,
5      int* B_ptr, int* B_ind, float* B_val, int B_row, int B_col, int B_mz,
6      int* resulting_row, int* resulting_col, float* resulting_val) {
7
8     int bid = blockIdx.x;
9     int tid = threadIdx.x;
10    int block_size = blockDim.x;
11
12    if (bid < A_row) {
13        int row = order[bid];
14        int A_row_start = A_ptr[row];
15        int A_row_end = A_ptr[row + 1];
16        int num_C_row_mz = counting[bid];
17        __shared__ float C_row_val[768];
18        __shared__ int C_row_col[768];
19        __shared__ int C_row_ind;
20        if (tid == 0) {
21            C_row_ind = 0;
22        }
23        __syncthreads();
24
25        for (int entryIdx = tid + A_row_start; entryIdx < A_row_end; entryIdx += block_size) {
26            int k = A_ind[entryIdx];
27            int A_kk = A_val[entryIdx];
28            int B_rowk_start = B_ptr[k];
29            int B_rowk_end = B_ptr[k + 1];
30            for (int i = B_rowk_start; i < B_rowk_end; i++) {
31                int j = B_ind[i];
32                int B_kj = B_val[i];
33                int pos = atomicAdd(&C_row_ind, 1);
34                C_row_val[pos] = A_kk * B_kj;
35                C_row_col[pos] = j;
36            }
37        }
38        __syncthreads();
39
40        if (num_C_row_mz == 32) {
41            if (tid == 0) {
42                bubbleSortNetwork(C_row_col, C_row_val, num_C_row_mz);
43            }
44        }
45        else {
46            typedef cub::BlockRadixSort<int, N_THREADS_PER_BLOCK, 6, float> BlockRadixSort;
47            // Allocate shared memory for BlockRadixSort
48            __shared__ typename BlockRadixSort::TempStorage temp_storage;
49            // Collectively sort the keys
50            BlockRadixSort(temp_storage).Sort(&static_cast<int*>[0], &static_cast<void*>[0], C_row_col + 6 * threadIdx.x, &static_cast<float*>[0], &static_cast<void*>[0], C_row_val + 6 * threadIdx.x));
51        }
52
53        contractionOperation(C_row_col, C_row_val, num_C_row_mz);
54        for (int resIdx = tid; resIdx < num_C_row_mz; resIdx += block_size) {
55            if (C_row_val[resIdx] != 0.f) {
56                resulting_row[operations[bid] + resIdx] = row;
57                resulting_col[operations[bid] + resIdx] = C_row_col[resIdx];
58                resulting_val[operations[bid] + resIdx] = C_row_val[resIdx];
59            }
60        }
61        __syncthreads();
62    }

```

Figure 20: Code: Expansion Sorting and Contraction Kernel

4 COMPUTATIONAL RESULT

We employed two dimensions to evaluate our performance: matrix size and sparse ratio. The sparse ratio represents the ratio of the total number of entries in the matrix to the number of nonzero entries. The performance results are presented in Fig. 21(a). As observed in the figure, the time required to perform matrix multiplication increases with an increase in the matrix size, ranging from 10 to 100. Additionally, the time taken for the multiplication process increases significantly as

the sparse ratio decreases. Fig. 21(b) illustrates our performance during the setup phase. While the patterns of the processing time are not entirely clear with respect to matrix size and sparse ratio, the processing time falls within a narrow range of 2.2×10^{-4} to 2.6×10^{-4} seconds. Therefore, our findings indicate that it takes a comparable amount of time to process the setup phase for matrix sizes ranging from 10 to 100 and sparse ratios ranging from 5 to 23. The performance of the ESC algorithm is presented in Fig. 21(c), where the processing time is plotted against the matrix size and sparse ratio. The figure clearly indicates that the processing time of the ESC algorithm increases with an increase in the matrix size, which aligns with our expectations. Moreover, we observe that the processing time tends to increase as the sparse ratio decreases, and this trend becomes more prominent for larger matrix sizes.

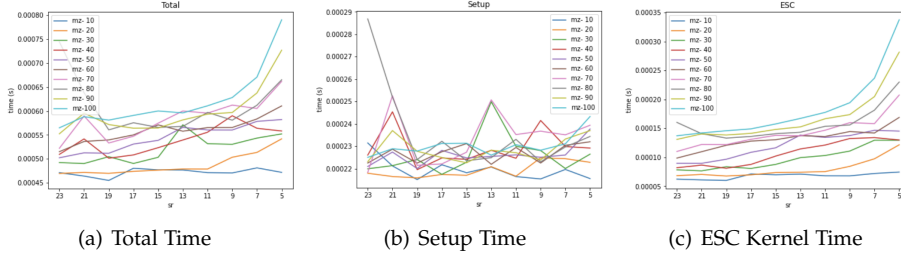


Figure 21: Computational Results: The Modified ESC Algorithm

We conducted a performance comparison of our implementation with two libraries developed by the authors, namely **CUSP** and **CUSParse**, both of which are used for matrix multiplication. The comparison was performed using the same sparse matrices, and the resulting performance is organized in Table. 1. Surprisingly, our implementation outperformed both libraries, with **cusparsSpGEMM_compute** demonstrating better performance than **cusp::multiply**. It is important to note that the total processing time for our implementation ranged from 4×10^{-4} to 8×10^{-4} , whereas CUSParse and CUSP libraries ranged from 8×10^{-4} to 1.5×10^{-3} and 110^{-3} to 6×10^{-3} , respectively.

Matrix Size	The Modified ESC		CUSparse		CUSP
	Total	ESC Kernel	Total	Compute	Total
10	0.000472	0.000068	0.000940	0.000292	0.000881
20	0.000488	0.000080	0.000961	0.000327	0.001626
30	0.000521	0.000098	0.001009	0.000346	0.001343
40	0.000539	0.000107	0.000946	0.000340	0.001445
50	0.000544	0.000120	0.000952	0.000333	0.001607
60	0.000559	0.000132	0.000977	0.000340	0.001860
70	0.000584	0.000143	0.001055	0.000355	0.001586
80	0.000611	0.000157	0.001002	0.000327	0.001491
90	0.000598	0.000168	0.001003	0.000336	0.001815
100	0.000622	0.000184	0.001095	0.000341	0.001785

Table 1: Computational Results

5 CONCLUSION

In the project, we reproduce the modified ESC algorithm from scratch, which exhibits notable speedup compared with the previous version. The setup phase, which contains a reordering scheme, is employed to process the intermediate matrix more effectively, specifically to reduce the load balance between adjacent blocks.

The number of total operations per row of the intermediate matrix is computed in the setup phase, and accordingly, the permutation is made. Besides, the shared memory is thoroughly utilized in the sorting and contraction phase. By parallelly sorting and contracting each part of the intermediate matrix in the shared memory, the time-consuming global operations are reduced, finally leading to performance improvement. In contrast, it is clear that the modified ESC algorithm, which relies on many operations performed in shared memory for speedup, does not behave better in all cases, especially when the shared memory size becomes a bottleneck. The referenced paper mentioned the ecumenical performance with a larger matrix size and small sparse ratio. Due to the computational capacity of the machine, such cases were not included in our tests. However, this will be the direction for our further exploration and enhancement.

6 REFERENCE

- Bell, N., Dalton, S., & Olson, L. N. (2012). Exposing Fine-Grained Parallelism in Algebraic Multigrid Methods. *SIAM Journal on Scientific Computing*, 34(4), C123–C152.
<https://doi.org/10.1137/110838844>
- Bell, N., & Garland, M. (2008). Efficient Sparse Matrix-Vector Multiplication on CUDA. NVIDIA Corporation.
- Chan, T. M. (2007). More Algorithms for All-Pairs Shortest Paths in Weighted Graphs. *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, 590–598.
<https://doi.org/10.1145/1250790.1250877>
- Cohen, E. (1996). On Optimizing Multiplications of Sparse Matrices. *Conference on Integer Programming and Combinatorial Optimization*.
- Davis, T. A. (2018). Graph algorithms via SuiteSparse: GraphBLAS: triangle counting and K-truss. *2018 IEEE High Performance Extreme Computing Conference (HPEC)*, 1–6.
<https://doi.org/10.1109/HPEC.2018.8547538>
- Dalton, S., Olson, L., & Bell, N. (2015). Optimizing Sparse Matrix—Matrix Multiplication for the GPU. *ACM Transactions on Mathematical Software*, 41(4), 1–20.
<https://doi.org/10.1145/2699470>
- Gremse, F., Höfter, A., Schwen, L. O., Kiessling, F., & Naumann, U. (2015). GPU-Accelerated Sparse Matrix-Matrix Multiplication by Iterative Row Merging. *SIAM Journal on Scientific Computing*, 37(1), C54–C71.
<https://doi.org/10.1137/130948811>
- Gao, J., Ji, W., Tan, Z., & Zhao, Y. (2023). A Systematic Survey of General Sparse Matrix-Matrix Multiplication. *ACM Computing Surveys*, 55(12), 1–36.
<https://doi.org/10.1145/3571157>
- Garland, M., Krik, D. B. (2010). Understanding throughput-oriented architectures. *Communications of the ACM*, 53(11), 58–66.
<https://doi.org/10.1145/1839676.1839694>
- Kaplan, H., Sharir, M., & Verbin, E. (2006). Colored intersection searching via sparse rectangular matrix multiplication. *Proceedings of the Twenty-Second Annual Symposium on Computational Geometry*, 52–60.
<https://doi.org/10.1145/1137856.1137866>
- Liu, W., & Vinter, B. (2015). A Framework for General Sparse Matrix-Matrix Multiplication on GPUs and Heterogeneous Processors. *Journal of Parallel and Distributed Computing*, 85, 47–61.
<https://doi.org/10.1016/j.jpdc.2015.06.010>
- Li, K., Yang, W., & Li, K. (2015). Performance Analysis and Optimization for SpMV on GPU Using Probabilistic Modeling. *IEEE Transactions on Parallel and Distributed Systems*, 26(1), 196–205.
<https://doi.org/10.1109/TPDS.2014.2308221>
- Then, M., Kaufmann, M., Chirigati, F., Hoang-Vu, T.-A., Pham, K., Kemper, A., Neumann, T., & Vo, H. T. (2014). The more the merrier: Efficient multi-source graph traversal. *Proceedings of the VLDB Endowment*, 8(4), 449–460.
<https://doi.org/10.14778/2735496.2735507>

7 APPENDIX

This section presented all codes used to implement and test the algorithm, as well as the sample output of the algorithm. It consists of four parts:

- 7.1: implementing and testing the modified ESC algorithm
- 7.2: implementing and testing the `cusparseSpGEMM_compute`
- 7.3: implementing and testing the `cusparse::multiply`
- 7.4: demonstration used in the Sec. 3.

7.1 The Modified ESC Algorithm: Implementation, Test, and Sample Output

```

1  #include <iostream>
2  #include <iomanip>
3  #include <stdio.h>
4  #include <thrust/device_ptr.h>
5  #include <thrust/device_vector.h>
6  #include <thrust/host_vector.h>
7  #include <thrust/sort.h>
8  #include <thrust/copy.h>
9  #include <thrust/scan.h>
10 #include <thrust/sequence.h>
11 #include <thrust/remove.h>
12 #include <thrust/iterator/zip_iterator.h>
13 #include <cub/cub.cuh>
14
15 #include <chrono>
16
17 using namespace std::chrono;
18
19 #define N_WARPS_PER_BLOCK (1 << 2)
20 #define WARP_SIZE (1 << 5)
21 #define N_THREADS_PER_BLOCK (1 << 7)
22
23 void printVector(float* counting, int nrow) {
24     for (int i = 0; i < nrow; i++) {
25         std::cout << counting[i] << " ";
26     }
27     std::cout << "\n";
28 }
29
30 int generateASparseMatrixRandomly(int nrow, int ncol, float** result_matrix, int
    sparse_ratio) {
31     float* A = (float*)malloc(sizeof(float) * nrow * ncol);
32     int nnz = 0;
33     int r;
34     float* cur;
35     for (int i = 0; i < nrow; i++) {
36         for (int j = 0; j < ncol; j++) {
37             r = rand();
38             cur = A + (i * ncol + j);
39             if (r % sparse_ratio == 0) { *cur = 10.0 * (r / (double)RAND_MAX); }
40             else { *cur = 0.0f; }
41             if (*cur != 0.0f) { nnz++; }
42         }
43     }
44     *result_matrix = A;
45     return nnz;
46 }
47
48 void convertToCSRFormat(float* mat, int nrow, int ncol, int nnz, int** ptr, int** indices,
    float** data) {
49     int* row_ptr = (int*)malloc(sizeof(int) * (nrow + 1));
50     int* col_ind = (int*)malloc(sizeof(int) * nnz);
51     float* nz_val = (float*)malloc(sizeof(float) * nnz);

```

```

52 float* cur;
53 int count = 0;
54 for (int i = 0; i < nrow; i++) {
55     row_ptr[i] = count;
56     for (int j = 0; j < ncol; j++) {
57         cur = mat + (i * ncol + j);
58         if (*cur != 0.0f) {
59             col_ind[count] = j;
60             nz_val[count] = *cur;
61             count++;
62         }
63     }
64 }
65 row_ptr[nrow] = count;
66
67 *ptr = row_ptr;
68 *indices = col_ind;
69 *data = nz_val;
70 return;
71 }
72
73 __global__
74 void countingKernel(int n, int* d_counting, int* A_ptr, int* A_ind, int* B_ptr) {
75     int idx = blockIdx.x * blockDim.x + threadIdx.x;
76     __syncthreads();
77     if (idx < n) {
78         int start = A_ptr[idx];
79         int end = A_ptr[idx + 1];
80         for (int i = start; i < end; i++) {
81             int j = A_ind[i];
82             d_counting[idx] = d_counting[idx] + B_ptr[j + 1] - B_ptr[j];
83         }
84     }
85     __syncthreads();
86 }
87
88 __device__
89 void bubbleSortSwap(int* col_ind, int i, int j) {
90     int temp = col_ind[i];
91     col_ind[i] = col_ind[j];
92     col_ind[j] = temp;
93 }
94
95 __device__
96 void bubbleSortSwap(float* ety_val, int i, int j) {
97     float temp = ety_val[i];
98     ety_val[i] = ety_val[j];
99     ety_val[j] = temp;
100 }
101
102 __device__
103 void bubbleSortNetwork(int* col_ind, float* ety_val, int n) {
104     for (int i = 0; i < n; i++) {
105         for (int j = i + 1; j < n; j++) {
106             if (col_ind[i] > col_ind[j]) {
107                 bubbleSortSwap(col_ind, i, j);
108                 bubbleSortSwap(ety_val, i, j);
109             }
110         }
111     }
112 }
113
114 __device__
115 void contractionOperation(int* col_ind, float* ety_val, int n) {
116     int tid = threadIdx.x;
117     if (tid < n) {
118         if (tid == 0 || col_ind[tid] != col_ind[tid - 1]) {
119             float res = ety_val[tid];
120             for (int j = tid + 1; j < n && col_ind[j] == col_ind[tid]; j++) {
121                 res += ety_val[j];
122                 col_ind[j] = 0;

```

```

123         ety_val[j] = 0;
124     }
125     ety_val[tid] = res;
126 }
127 }
128 }
129
130 __global__
131 void expansionSortingContractionKernel(int* order, int* counting, int* operations,
132                                       int* A_ptr, int* A_ind, float* A_val, int A_row,
133                                       int A_col, int A_nnz,
134                                       int* B_ptr, int* B_ind, float* B_val, int B_row,
135                                       int B_col, int B_nnz,
136                                       int* resulting_row, int* resulting_col, float*
137                                       resulting_val) {
138     int bid = blockIdx.x;
139     int tid = threadIdx.x;
140     int block_size = blockDim.x;
141
142     if (bid < A_row) {
143         int row = order[bid];
144         int A_row_start = A_ptr[row];
145         int A_row_end = A_ptr[row + 1];
146         int num_C_row_nnz = counting[bid];
147         __shared__ float C_row_val[768];
148         __shared__ int C_row_col[768];
149         __shared__ int C_row_ind;
150         if (tid == 0) {
151             C_row_ind = 0;
152         }
153         __syncthreads();
154
155         for (int entryIdx = tid + A_row_start; entryIdx < A_row_end; entryIdx +=
156             block_size) {
157             int k = A_ind[entryIdx];
158             int A_ik = A_val[entryIdx];
159             int B_rowk_start = B_ptr[k];
160             int B_rowk_end = B_ptr[k + 1];
161             for (int i = B_rowk_start; i < B_rowk_end; i++) {
162                 int j = B_ind[i];
163                 int B_kj = B_val[i];
164                 int pos = atomicAdd(&C_row_ind, 1);
165                 C_row_val[pos] = A_ik * B_kj;
166                 C_row_col[pos] = j;
167             }
168         }
169
170         __syncthreads();
171
172         if (num_C_row_nnz <= 32) {
173             if (tid == 0) {
174                 bubbleSortNetwork(C_row_col, C_row_val, num_C_row_nnz);
175             }
176         }
177         else {
178             typedef cub::BlockRadixSort<int, N_THREADS_PER_BLOCK, 6, float> BlockRadixSort
179             ;
180             // Allocate shared memory for BlockRadixSort
181             __shared__ typename BlockRadixSort::TempStorage temp_storage;
182             // Collectively sort the keys
183             BlockRadixSort(temp_storage).Sort(*static_cast<int*>[6]>(<static_cast<void*>(<
184                 C_row_col + 6 * threadIdx.x), *static_cast<float*>[6]>(<static_cast<void*>(<
185                 C_row_val + 6 * threadIdx.x)))));
186         }
187
188         contractionOperation(C_row_col, C_row_val, num_C_row_nnz);
189         for (int resIdx = tid; resIdx < num_C_row_nnz; resIdx += block_size) {
190             if (C_row_val[resIdx] != 0.f) {
191                 resulting_row[operations[bid] + resIdx] = row;
192                 resulting_col[operations[bid] + resIdx] = C_row_col[resIdx];
193                 resulting_val[operations[bid] + resIdx] = C_row_val[resIdx];
194             }
195         }
196     }
197 }

```

```

187     }
188 }
189 }
190 __syncthreads();
191 }
192
193 int SPMMM(int* A_ptr, int* A_ind, float* A_val, int A_row, int A_col, int A_nnz,
194          int* B_ptr, int* B_ind, float* B_val, int B_row, int B_col, int B_nnz,
195          int** result_row, int** result_col, float** result_val,
196          float* recorded_time) {
197
198
199     auto tot_start_time = high_resolution_clock::now();
200
201     int* d_A_ptr, * d_A_ind, * d_B_ptr, * d_B_ind;
202     float* d_A_val, * d_B_val;
203     cudaMalloc(&d_A_ptr, (A_row + 1) * sizeof(int));
204     cudaMalloc(&d_A_ind, A_nnz * sizeof(int));
205     cudaMalloc(&d_A_val, A_nnz * sizeof(float));
206     cudaMalloc(&d_B_ptr, (B_row + 1) * sizeof(int));
207     cudaMalloc(&d_B_ind, B_nnz * sizeof(int));
208     cudaMalloc(&d_B_val, B_nnz * sizeof(float));
209     cudaMemcpy(d_A_ptr, A_ptr, (A_row + 1) * sizeof(int), cudaMemcpyHostToDevice);
210     cudaMemcpy(d_A_ind, A_ind, A_nnz * sizeof(int), cudaMemcpyHostToDevice);
211     cudaMemcpy(d_A_val, A_val, A_nnz * sizeof(float), cudaMemcpyHostToDevice);
212     cudaMemcpy(d_B_ptr, B_ptr, (B_row + 1) * sizeof(int), cudaMemcpyHostToDevice);
213     cudaMemcpy(d_B_ind, B_ind, B_nnz * sizeof(int), cudaMemcpyHostToDevice);
214     cudaMemcpy(d_B_val, B_val, B_nnz * sizeof(float), cudaMemcpyHostToDevice);
215
216
217     auto setup_start_time = high_resolution_clock::now();
218     int* counting = (int*)malloc(A_row * sizeof(int));
219     int* d_counting;
220     cudaMalloc(&d_counting, A_row * sizeof(int));
221     cudaMemcpy(d_counting, counting, A_row * sizeof(int), cudaMemcpyHostToDevice);
222     int set_up_n_blocks = (A_row + N_THREADS_PER_BLOCK - 1) / N_THREADS_PER_BLOCK;
223
224     countingKernel <<< set_up_n_blocks, N_THREADS_PER_BLOCK >>> (A_row, d_counting,
225     d_A_ptr, d_A_ind, d_B_ptr);
226     cudaDeviceSynchronize();
227     cudaMemcpy(counting, d_counting, A_row * sizeof(int), cudaMemcpyDeviceToHost);
228
229     thrust::device_ptr<int> d_counting_ptr(d_counting);
230     thrust::device_vector<int> d_counting_vec(d_counting_ptr, d_counting_ptr + A_row);
231     thrust::device_vector<int> d_order_vec(A_row);
232     thrust::sequence(d_order_vec.begin(), d_order_vec.end());
233     thrust::stable_sort_by_key(d_counting_vec.begin(), d_counting_vec.end(), d_order_vec.
234     begin(), thrust::greater<int>());
235     thrust::device_vector<int> d_operations_vec(A_row + 1);
236     thrust::exclusive_scan(d_counting_vec.begin(), d_counting_vec.end(), d_operations_vec.
237     begin());
238     int tot_operations = d_counting_vec.back() + d_operations_vec[A_row - 1];
239     d_operations_vec[A_row] = tot_operations;
240     auto setup_end_time = high_resolution_clock::now();
241
242     int* order = (int*)malloc(A_row * sizeof(int));
243     int* operations = (int*)malloc((A_row + 1) * sizeof(int));
244     thrust::copy(d_order_vec.begin(), d_order_vec.end(), order);
245     thrust::copy(d_operations_vec.begin(), d_operations_vec.end(), operations);
246     thrust::copy(d_counting_vec.begin(), d_counting_vec.end(), counting);
247
248     if (counting[0] > 768) {
249         return 1;
250     }
251
252     int* d_order;
253     int* d_operations;
254     cudaMalloc(&d_order, A_row * sizeof(int));
255     cudaMalloc(&d_operations, (A_row + 1) * sizeof(int));
256     cudaMemcpy(d_order, order, A_row * sizeof(int), cudaMemcpyHostToDevice);

```

```

254 cudaMemcpy(d_operations, operations, (A_row + 1) * sizeof(int), cudaMemcpyHostToDevice
);
255 cudaMemcpy(d_counting, counting, A_row * sizeof(int), cudaMemcpyHostToDevice);
256
257 int* d_resulting_row, * d_resulting_col;
258 float* d_resulting_val;
259 cudaMalloc(&d_resulting_row, tot_operations * sizeof(int));
260 cudaMalloc(&d_resulting_col, tot_operations * sizeof(int));
261 cudaMalloc(&d_resulting_val, tot_operations * sizeof(float));
262
263 int expansion_and_sorting_n_blocks = A_row;
264
265 auto esc_start_time = high_resolution_clock::now();
266 expansionSortingContractionKernel <<< expansion_and_sorting_n_blocks,
    N_THREADS_PER_BLOCK >>> (d_order, d_counting, d_operations,
267
268         d_A_ptr, d_A_ind, d_A_val, A_row, A_col, A_nnz,
269
270         d_B_ptr, d_B_ind, d_B_val, B_row, B_col, B_nnz,
271
272         d_resulting_row, d_resulting_col, d_resulting_val);
273 cudaDeviceSynchronize();
274
275 int* resulting_row = (int*)malloc(tot_operations * sizeof(int));
276 int* resulting_col = (int*)malloc(tot_operations * sizeof(int));
277 float* resulting_val = (float*)malloc(tot_operations * sizeof(float));
278 cudaMemcpy(resulting_row, d_resulting_row, tot_operations * sizeof(int),
    cudaMemcpyDeviceToHost);
279 cudaMemcpy(resulting_col, d_resulting_col, tot_operations * sizeof(int),
    cudaMemcpyDeviceToHost);
280 cudaMemcpy(resulting_val, d_resulting_val, tot_operations * sizeof(float),
    cudaMemcpyDeviceToHost);
281
282 cudaFree(d_A_ptr);
283 cudaFree(d_A_ind);
284 cudaFree(d_A_val);
285 cudaFree(d_B_ptr);
286 cudaFree(d_B_ind);
287 cudaFree(d_B_val);
288 cudaFree(d_counting);
289 cudaFree(d_order);
290 cudaFree(d_operations);
291 cudaFree(d_resulting_row);
292 cudaFree(d_resulting_col);
293 cudaFree(d_resulting_val);
294
295 auto esc_end_time = high_resolution_clock::now();
296
297 free(A_ptr);
298 free(A_ind);
299 free(A_val);
300 free(B_ptr);
301 free(B_ind);
302 free(B_val);
303 free(counting);
304 free(order);
305 free(operations);
306
307 auto tot_end_time = high_resolution_clock::now();
308
309 auto tot_elapsed_time = duration_cast<duration<double>>(tot_end_time - tot_start_time)
;
310 auto setup_elapsed_time = duration_cast<duration<double>>(setup_end_time -
    setup_start_time);
311 auto esc_elapsed_time = duration_cast<duration<double>>(esc_end_time - esc_start_time)
;
312 float tot_elapsed_time_sec = static_cast<float>(tot_elapsed_time.count());
313 float setup_elapsed_time_sec = static_cast<float>(setup_elapsed_time.count());
314 float esc_elapsed_time_sec = static_cast<float>(esc_elapsed_time.count());
315
316 *result_row = resulting_row;

```

```

314     *result_col = resulting_col;
315     *result_val = resulting_val;
316
317     free(resulting_row);
318     free(resulting_col);
319     free(resulting_val);
320
321     recorded_time[0] += tot_elapsed_time_sec;
322     recorded_time[1] += setup_elapsed_time_sec;
323     recorded_time[2] += esc_elapsed_time_sec;
324
325     return 0;
326 }
327
328 int main() {
329
330     int num_matrix_size = 10;
331     int* matrix_size_list = (int*)malloc(num_matrix_size * sizeof(int));
332     for (int i = 0; i < num_matrix_size; i++) {
333         matrix_size_list[i] = 10 * (i + 1);
334     }
335
336     int num_sparse_ratio = 10;
337     int* sparse_ratio_list = (int*)malloc(num_sparse_ratio * sizeof(int));
338     for (int i = 0; i < num_sparse_ratio; i++) {
339         sparse_ratio_list[i] = 23 - 2*i;
340     }
341
342     int repeat_times = 10;
343
344     float* time_tot = (float*)malloc(num_matrix_size * num_sparse_ratio * repeat_times *
345                                     sizeof(float));
346     float* time_setup = (float*)malloc(num_matrix_size * num_sparse_ratio * repeat_times
347                                     * sizeof(float));
348     float* time_esc = (float*)malloc(num_matrix_size * num_sparse_ratio * repeat_times *
349                                     sizeof(float));
350
351     for (int i = 0; i < num_matrix_size; i++) {
352         std::cout << " --- Matrix size: " << matrix_size_list[i] << " --- \n";
353         for (int j = 0; j < num_sparse_ratio; j++) {
354             std::cout << "          *** Sparse ratio: " << sparse_ratio_list[j] << " *** \n";
355             for (int k = 0; k < repeat_times; k++) {
356                 int mz = matrix_size_list[i];
357                 int sr = sparse_ratio_list[j];
358
359                 int A_row = mz;
360                 int A_col = mz;
361                 int B_row = mz;
362                 int B_col = mz;
363                 float* A;
364                 float* B;
365                 int A_nnz = generateASparseMatrixRandomly(A_row, A_col, &A, sr);
366                 int B_nnz = generateASparseMatrixRandomly(B_row, B_col, &B, sr);
367                 int* A_ptr;
368                 int* A_ind;
369                 float* A_val;
370                 convertToCSRFormat(A, A_row, A_col, A_nnz, &A_ptr, &A_ind, &A_val);
371                 int* B_ptr;
372                 int* B_ind;
373                 float* B_val;
374                 convertToCSRFormat(B, B_row, B_col, B_nnz, &B_ptr, &B_ind, &B_val);
375
376                 float* time = (float*)malloc(3 * sizeof(float));
377                 for (int i = 0; i < 3; i++) {
378                     time[i] = 0;
379                 }
380                 int* resulting_row;
381                 int* resulting_col;
382                 float* resulting_val;
383                 int exit_code = SPMMM(A_ptr, A_ind, A_val, A_row, A_col, A_nnz,

```

```

382         B_ptr, B_ind, B_val, B_row, B_col, B_nnz,
383         &resulting_row, &resulting_col, &resulting_val, time);
384
385         if (exit_code == 0) {
386             time1_tot[i * num_sparse_ratio * repeat_times + j * repeat_times + k]
= time[0];
387             time1_setup[i * num_sparse_ratio * repeat_times + j * repeat_times + k
] = time[1];
388             time1_esc[i * num_sparse_ratio * repeat_times + j * repeat_times + k]
= time[2];
389             std::cout << "           Experiment - " << i * num_sparse_ratio *
repeat_times + j * repeat_times + k << " | Time: " << time[0] << " | Setup: " << time
[1] << " | ESC: " << time[2] << " \n";
390         }
391         else {
392             time1_tot[i * num_sparse_ratio * repeat_times + j * repeat_times + k]
= 0;
393             time1_setup[i * num_sparse_ratio * repeat_times + j * repeat_times + k
] = 0;
394             time1_esc[i * num_sparse_ratio * repeat_times + j * repeat_times + k]
= 0;
395             std::cout << "           Experiment - " << i * num_sparse_ratio *
repeat_times + j * repeat_times + k << " | ERROR -> TRY NEXT \n";
396         }
397     }
398 }
399 }
400
401 std::cout << " ##### \n";
402 std::cout << " ### PRINT THE TIME1 ### \n";
403 std::cout << " ##### \n";
404 std::cout << " --- Time1 TOT --- \n";
405 printVector(time1_tot, num_matrix_size * num_sparse_ratio * repeat_times);
406 std::cout << " --- Time1 SETUP --- \n";
407 printVector(time1_setup, num_matrix_size * num_sparse_ratio * repeat_times);
408 std::cout << " --- Time1 ESC --- \n";
409 printVector(time1_esc, num_matrix_size * num_sparse_ratio * repeat_times);
410
411 return 0;
412 }

```

```

1  --- Matrix size: 100 ---
2  *** Sparse ratio: 23 ***
3      Experiment - 900 | Time: 0.000585435 | Setup: 0.000229657 | ESC: 0.000135981
4      Experiment - 901 | Time: 0.000581899 | Setup: 0.000223666 | ESC: 0.000136855
5      Experiment - 902 | Time: 0.000559657 | Setup: 0.000218306 | ESC: 0.000134322
6      Experiment - 903 | Time: 0.000588317 | Setup: 0.00024104 | ESC: 0.000143362
7      Experiment - 904 | Time: 0.00054217 | Setup: 0.000218708 | ESC: 0.000131468
8      Experiment - 905 | Time: 0.000544271 | Setup: 0.00021968 | ESC: 0.000136073
9      Experiment - 906 | Time: 0.000573116 | Setup: 0.000231894 | ESC: 0.000140474
10     Experiment - 907 | Time: 0.000559643 | Setup: 0.000228873 | ESC: 0.000140299
11     Experiment - 908 | Time: 0.000565668 | Setup: 0.000220703 | ESC: 0.000137745
12     Experiment - 909 | Time: 0.000548258 | Setup: 0.000217647 | ESC: 0.000134885
13  *** Sparse ratio: 21 ***
14     Experiment - 910 | Time: 0.000602272 | Setup: 0.000217146 | ESC: 0.0001405
15     Experiment - 911 | Time: 0.000562075 | Setup: 0.000223956 | ESC: 0.000143737
16     Experiment - 912 | Time: 0.0005475 | Setup: 0.000223572 | ESC: 0.000135017
17     Experiment - 913 | Time: 0.000592504 | Setup: 0.000236636 | ESC: 0.000139604
18     Experiment - 914 | Time: 0.000589779 | Setup: 0.000229959 | ESC: 0.00013651
19     Experiment - 915 | Time: 0.000660304 | Setup: 0.000239206 | ESC: 0.000147484
20     Experiment - 916 | Time: 0.000561433 | Setup: 0.000222313 | ESC: 0.00013684
21     Experiment - 917 | Time: 0.000604171 | Setup: 0.000229795 | ESC: 0.000152968
22     Experiment - 918 | Time: 0.000572919 | Setup: 0.00023531 | ESC: 0.000145699
23     Experiment - 919 | Time: 0.000587343 | Setup: 0.000230768 | ESC: 0.00014217

```

Figure 22: Sample Output

7.2 The *cusparseSpGEMM_compute*: Implementation, Test, and Sample Output

```

1  #include <iostream>
2  #include <iomanip>
3  #include <stdio.h>
4  #include <cusparse_v2.h>
5  #include <cuda.h>
6
7  #include <chrono>
8  using namespace std::chrono;
9
10 int generateASparseMatrixRandomly(int nrow, int ncol, float** result_matrix, int
    sparse_ratio) {
11     float* A = (float*)malloc(sizeof(float) * nrow * ncol);
12     int nnz = 0;
13     int r;
14     float* cur;
15     for (int i = 0; i < nrow; i++) {
16         for (int j = 0; j < ncol; j++) {
17             r = rand();
18             cur = A + (i * ncol + j);
19             if (r % sparse_ratio == 0) { *cur = 10.0 * (r / (double)RAND_MAX); }
20             else { *cur = 0.0f; }
21             if (*cur != 0.0f) { nnz++; }
22         }
23     }
24     *result_matrix = A;
25     return nnz;
26 }
27
28 void convertToCSRFormat(float* mat, int nrow, int ncol, int nnz, int** ptr, int** indices,
    float** data) {
29     int* row_ptr = (int*)malloc(sizeof(int) * (nrow + 1));
30     int* col_ind = (int*)malloc(sizeof(int) * nnz);
31     float* nz_val = (float*)malloc(sizeof(float) * nnz);
32     float* cur;
33     int count = 0;
34     for (int i = 0; i < nrow; i++) {
35         row_ptr[i] = count;
36         for (int j = 0; j < ncol; j++) {
37             cur = mat + (i * ncol + j);
38             if (*cur != 0.0f) {
39                 col_ind[count] = j;
40                 nz_val[count] = *cur;
41                 count++;
42             }
43         }
44     }
45     row_ptr[nrow] = count;
46
47     *ptr = row_ptr;
48     *indices = col_ind;
49     *data = nz_val;
50     return;
51 }
52
53 void SPMM_COMP(int* A_ptr, int* A_ind, float* A_val, int A_row, int A_col, int A_nnz,
54     int* B_ptr, int* B_ind, float* B_val, int B_row, int B_col, int B_nnz,
55     int** result_row, int** result_col, float** result_val,
56     float* recorded_time) {
57
58     auto tot_start_time = high_resolution_clock::now();
59
60     float alpha = 1.0f, beta = 0.0f;
61     cusparseOperation_t opA = CUSPARSE_OPERATION_NON_TRANSPOSE;
62     cusparseOperation_t opB = CUSPARSE_OPERATION_NON_TRANSPOSE;
63     cudaDataType computeType = CUDA_R_32F;
64
65     int* dA_csrOffsets, * dA_columns, * dB_csrOffsets, * dB_columns, * dC_csrOffsets, *
        dC_columns;

```



```

66 float* dA_values, * dB_values, * dC_values;
67
68 cudaMalloc((void**)&dA_csrOffsets, sizeof(int) * (A_row + 1));
69 cudaMalloc((void**)&dA_columns, sizeof(int) * A_nnz);
70 cudaMalloc((void**)&dA_values, sizeof(float) * A_nnz);
71 cudaMalloc((void**)&dB_csrOffsets, sizeof(int) * (B_row + 1));
72 cudaMalloc((void**)&dB_columns, sizeof(int) * B_nnz);
73 cudaMalloc((void**)&dB_values, sizeof(float) * B_nnz);
74 cudaMalloc((void**)&dC_csrOffsets, sizeof(int) * (A_row + 1));
75
76 cudaMemcpy(dA_csrOffsets, A_ptr, sizeof(int) * (A_row + 1), cudaMemcpyHostToDevice);
77 cudaMemcpy(dA_columns, A_ind, sizeof(int) * A_nnz, cudaMemcpyHostToDevice);
78 cudaMemcpy(dA_values, A_val, sizeof(float) * A_nnz, cudaMemcpyHostToDevice);
79 cudaMemcpy(dB_csrOffsets, B_ptr, sizeof(int) * (B_row + 1), cudaMemcpyHostToDevice);
80 cudaMemcpy(dB_columns, B_ind, sizeof(int) * B_nnz, cudaMemcpyHostToDevice);
81 cudaMemcpy(dB_values, B_val, sizeof(float) * B_nnz, cudaMemcpyHostToDevice);
82
83 cusparseHandle_t handle = NULL;
84 cusparseSpMatDescr_t matA, matB, matC;
85 void* dBuffer1 = NULL, * dBuffer2 = NULL;
86 size_t bufferSize1 = 0, bufferSize2 = 0;
87 cusparseCreate(&handle);
88
89 cusparseCreateCsr(&matA, A_row, A_col, A_nnz, dA_csrOffsets, dA_columns, dA_values,
90 CUSPARSE_INDEX_32I, CUSPARSE_INDEX_32I, CUSPARSE_INDEX_BASE_ZERO, CUDA_R_32F);
91 cusparseCreateCsr(&matB, B_row, B_col, B_nnz, dB_csrOffsets, dB_columns, dB_values,
92 CUSPARSE_INDEX_32I, CUSPARSE_INDEX_32I, CUSPARSE_INDEX_BASE_ZERO, CUDA_R_32F);
93 cusparseCreateCsr(&matC, A_row, B_col, 0, NULL, NULL, NULL, CUSPARSE_INDEX_32I,
94 CUSPARSE_INDEX_32I, CUSPARSE_INDEX_BASE_ZERO, CUDA_R_32F);
95
96 auto compute_start_time = high_resolution_clock::now();
97 cusparseSpGEMMDescr_t spgemmDesc;
98 cusparseSpGEMM_createDescr(&spgemmDesc);
99 cusparseSpGEMM_workEstimation(handle, opA, opB, &alpha, matA, matB, &beta, matC,
100 computeType, CUSPARSE_SPGEMM_DEFAULT, spgemmDesc, &bufferSize1, NULL);
101 cudaMalloc((void**)&dBuffer1, bufferSize1);
102 cusparseSpGEMM_workEstimation(handle, opA, opB, &alpha, matA, matB, &beta, matC,
103 computeType, CUSPARSE_SPGEMM_DEFAULT, spgemmDesc, &bufferSize1, dBuffer1);
104 cusparseSpGEMM_compute(handle, opA, opB, &alpha, matA, matB, &beta, matC, computeType,
105 CUSPARSE_SPGEMM_DEFAULT, spgemmDesc, &bufferSize2, NULL);
106 cudaMalloc((void**)&dBuffer2, bufferSize2);
107 cusparseSpGEMM_compute(handle, opA, opB, &alpha, matA, matB, &beta, matC, computeType,
108 CUSPARSE_SPGEMM_DEFAULT, spgemmDesc, &bufferSize2, dBuffer2);
109
110 int64_t C_num_rows1, C_num_cols1, C_nnz1;
111 cusparseSpMatGetSize(matC, &C_num_rows1, &C_num_cols1, &C_nnz1);
112 cudaMalloc((void**)&dC_columns, C_nnz1 * sizeof(int));
113 cudaMalloc((void**)&dC_values, C_nnz1 * sizeof(float));
114 cusparseCsrSetPointers(matC, dC_csrOffsets, dC_columns, dC_values);
115 cusparseSpGEMM_copy(handle, opA, opB, &alpha, matA, matB, &beta, matC, computeType,
116 CUSPARSE_SPGEMM_DEFAULT, spgemmDesc);
117
118 cusparseSpGEMM_destroyDescr(spgemmDesc);
119 cusparseDestroySpMat(matA);
120 cusparseDestroySpMat(matB);
121 cusparseDestroySpMat(matC);
122 cusparseDestroy(handle);
123 auto compute_end_time = high_resolution_clock::now();
124
125 int* C_ptr = (int*)malloc(sizeof(int) * (A_row + 1));
126 int* C_ind = (int*)malloc(sizeof(int) * C_nnz1);
127 float* C_val = (float*)malloc(sizeof(float) * C_nnz1);
128
129 cudaMemcpy(C_ptr, dC_csrOffsets, sizeof(int) * (A_row + 1), cudaMemcpyDeviceToHost);
130 cudaMemcpy(C_ind, dC_columns, sizeof(int) * C_nnz1, cudaMemcpyDeviceToHost);
131 cudaMemcpy(C_val, dC_values, sizeof(float) * C_nnz1, cudaMemcpyDeviceToHost);
132
133 cudaFree(dBuffer1);
134 cudaFree(dBuffer2);
135 cudaFree(dA_csrOffsets);
136 cudaFree(dA_columns);

```

```

129 cudaFree(dA_values);
130 cudaFree(dB_csrOffsets);
131 cudaFree(dB_columns);
132 cudaFree(dB_values);
133 cudaFree(dC_csrOffsets);
134 cudaFree(dC_columns);
135 cudaFree(dC_values);
136
137 *result_row = C_ptr;
138 *result_col = C_ind;
139 *result_val = C_val;
140
141 auto tot_end_time = std::chrono::high_resolution_clock::now();
142 auto tot_elapsed_time = duration_cast<duration<double>>(tot_end_time - tot_start_time)
143 ;
144 float tot_elapsed_time_sec = static_cast<float>(tot_elapsed_time.count());
145
146 auto compute_elapsed_time = duration_cast<duration<double>>(compute_end_time -
147 compute_start_time);
148 float compute_elapsed_time_sec = static_cast<float>(compute_elapsed_time.count());
149
150 recorded_time[0] += tot_elapsed_time_sec;
151 recorded_time[1] += compute_elapsed_time_sec;
152 }
153 int main() {
154
155     int num_matrix_size = 10;
156     int* matrix_size_list = (int*)malloc(num_matrix_size * sizeof(int));
157     for (int i = 0; i < num_matrix_size; i++) {
158         matrix_size_list[i] = 10 * (i + 1);
159     }
160
161     int num_sparse_ratio = 10;
162     int* sparse_ratio_list = (int*)malloc(num_sparse_ratio * sizeof(int));
163     for (int i = 0; i < num_sparse_ratio; i++) {
164         sparse_ratio_list[i] = 23 - 2*i;
165     }
166
167     int repeat_times = 10;
168
169     float* time2_tot = (float*)malloc(num_matrix_size * num_sparse_ratio * repeat_times *
170 sizeof(float));
171 float* time2_compute = (float*)malloc(num_matrix_size * num_sparse_ratio *
172 repeat_times * sizeof(float));
173
174 for (int i = 0; i < num_matrix_size; i++) {
175     std::cout << " --- Matrix size: " << matrix_size_list[i] << " --- \n";
176     for (int j = 0; j < num_sparse_ratio; j++) {
177         std::cout << " *** Sparse ratio: " << sparse_ratio_list[j] << " *** \n";
178         for (int k = 0; k < repeat_times; k++) {
179             int mz = matrix_size_list[i];
180             int sr = sparse_ratio_list[j];
181
182             int A_row = mz;
183             int A_col = mz;
184             int B_row = mz;
185             int B_col = mz;
186             float* A;
187             float* B;
188             int A_nnz = generateASparseMatrixRandomly(A_row, A_col, &A, sr);
189             int B_nnz = generateASparseMatrixRandomly(B_row, B_col, &B, sr);
190             int* A_ptr;
191             int* A_ind;
192             float* A_val;
193             convertToCSRFormat(A, A_row, A_col, A_nnz, &A_ptr, &A_ind, &A_val);
194             int* B_ptr;
195             int* B_ind;
196             float* B_val;
197             convertToCSRFormat(B, B_row, B_col, B_nnz, &B_ptr, &B_ind, &B_val);

```

```

196     float* time = (float*)malloc(2 * sizeof(float));
197     for (int i = 0; i < 2; i++) {
198         time[i] = 0;
199     }
200     int* resulting_row;
201     int* resulting_col;
202     float* resulting_val;
203     SPMM_COMP(A_ptr, A_ind, A_val, A_row, A_col, A_nnz,
204               B_ptr, B_ind, B_val, B_row, B_col, B_nnz,
205               &resulting_row, &resulting_col, &resulting_val, time);
206
207     time2_tot[i * num_sparse_ratio * repeat_times + j * repeat_times + k] =
208     time[0];
209     time2_compute[i * num_sparse_ratio * repeat_times + j * repeat_times + k]
210     = time[1];
211     std::cout << "                Experiment - " << i * num_sparse_ratio *
212     repeat_times + j * repeat_times + k << " | Time: " << time[0] << " | Compute: " <<
213     time[1] << " \n";
214 }
215 }

```

```

1  --- Matrix size: 100 ---
2      *** Sparse ratio: 23 ***
3          Experiment - 900 | Time: 0.0011049 | Compute: 0.0004076
4          Experiment - 901 | Time: 0.0010467 | Compute: 0.0003965
5          Experiment - 902 | Time: 0.0012692 | Compute: 0.0002999
6          Experiment - 903 | Time: 0.0009554 | Compute: 0.0002961
7          Experiment - 904 | Time: 0.0008819 | Compute: 0.0002716
8          Experiment - 905 | Time: 0.0008789 | Compute: 0.0002861
9          Experiment - 906 | Time: 0.000916 | Compute: 0.0003309
10         Experiment - 907 | Time: 0.0008767 | Compute: 0.0002917
11         Experiment - 908 | Time: 0.0008867 | Compute: 0.0002781
12         Experiment - 909 | Time: 0.0009151 | Compute: 0.0002915
13     *** Sparse ratio: 21 ***
14         Experiment - 910 | Time: 0.001098 | Compute: 0.0003705
15         Experiment - 911 | Time: 0.0010377 | Compute: 0.0003449
16         Experiment - 912 | Time: 0.0008799 | Compute: 0.0002825
17         Experiment - 913 | Time: 0.000887 | Compute: 0.0002852
18         Experiment - 914 | Time: 0.0010446 | Compute: 0.0003045
19         Experiment - 915 | Time: 0.0011512 | Compute: 0.0004064
20         Experiment - 916 | Time: 0.0011072 | Compute: 0.0003833
21         Experiment - 917 | Time: 0.0009178 | Compute: 0.0003398
22         Experiment - 918 | Time: 0.000936 | Compute: 0.0002983
23         Experiment - 919 | Time: 0.0008782 | Compute: 0.0002889

```

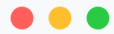
Figure 23: Sample Output

7.3 The `cuspl::multiply`: Implementation, Test, and Sample Output

```

1  #include <iostream>
2  #include <iomanip>
3  #include <stdio.h>
4  #include <thrust/device_ptr.h>
5  #include <thrust/device_vector.h>
6  #include <thrust/host_vector.h>
7  #include <thrust/sort.h>
8  #include <thrust/copy.h>
9  #include <thrust/scan.h>
10 #include <thrust/sequence.h>
11 #include <thrust/remove.h>
12 #include <thrust/iterator/zip_iterator.h>
13 #include <cub/cub.cuh>
14 #include <cuspl/array2d.h>
15 #include <cuspl/csr_matrix.h>
16 #include <cuspl/coo_matrix.h>
17 #include <cuspl/multiply.h>
18 #include <cuspl/gallery/poisson.h>
19 #include <cuspl/gallery/random.h>
20 #include <time.h>
21 #include <chrono>
22
23 using namespace std::chrono;
24
25 #define N_WARPS_PER_BLOCK (1 << 2)
26 #define WARP_SIZE (1 << 5)
27 #define N_THREADS_PER_BLOCK (1 << 7)
28 int main() {
29     for (size_t i = 1; i <= 10; i++)
30     {
31         std::cout << "--- Matrix size: " << i << " -- - " << std::endl;
32         for (size_t j = 5; j < 24; j++)
33         {
34             std::cout << " *** Sparse ratio: " << j << " * *" << std::endl;
35             int mz = i * 10;
36             cuspl::csr_matrix<int, float, cuspl::device_memory> c_a;
37             cuspl::gallery::random(c_a, mz, mz, mz * j / 100);
38             cuspl::csr_matrix<int, float, cuspl::device_memory> c_b;
39             cuspl::gallery::random(c_b, mz, mz, mz * j / 100);
40             cuspl::csr_matrix<int, float, cuspl::device_memory> c_c;
41
42             std::chrono::time_point<std::chrono::system_clock> start, end;
43             start = std::chrono::system_clock::now();
44             cuspl::multiply(c_a, c_b, c_c);
45             end = std::chrono::system_clock::now();
46             std::chrono::duration<double> elapsed_seconds = end - start;
47             std::cout << "matrix size: " << i * 10 << ", cuspl::multiply: " <<
elapsed_seconds.count() << " s" << std::endl;
48         }
49     }
50     return 0;
51 }

```



```
1 --- Matrix size: 10 -- -
2   *** Sparse ratio:5* **
3   cusp::multiply: 0.0001258 s
4   *** Sparse ratio:6* **
5   cusp::multiply: 0.0001457 s
6   *** Sparse ratio:7* **
7   cusp::multiply: 0.0001257 s
8   *** Sparse ratio:8* **
9   cusp::multiply: 0.0001277 s
10  *** Sparse ratio:9* **
11  cusp::multiply: 0.0001045 s
12  *** Sparse ratio:10* **
13  cusp::multiply: 0.0008739 s
```

Figure 24: Sample Output

7.4 The Modified ESC Algorithm: Demonstration

```

1  #include <iostream>
2  #include <iomanip>
3  #include <stdio.h>
4  #include <thrust/device_ptr.h>
5  #include <thrust/device_vector.h>
6  #include <thrust/host_vector.h>
7  #include <thrust/sort.h>
8  #include <thrust/copy.h>
9  #include <thrust/scan.h>
10 #include <thrust/sequence.h>
11 #include <thrust/remove.h>
12 #include <thrust/iterator/zip_iterator.h>
13 #include <cub/cub.cuh>
14
15 #define N_WARPS_PER_BLOCK (1 << 2)
16 #define WARP_SIZE (1 << 5)
17 #define N_THREADS_PER_BLOCK (1 << 7)
18
19 int generateASparseMatrixRandomly(int nrow, int ncol, float** result_matrix, int
    sparse_ratio) {
20     float* A = (float*)malloc(sizeof(float) * nrow * ncol);
21     int nnz = 0;
22     int r;
23     float* cur;
24     for (int i = 0; i < nrow; i++) {
25         for (int j = 0; j < ncol; j++) {
26             r = rand();
27             cur = A + (i * ncol + j);
28             if (r % sparse_ratio == 0) { *cur = 10.0 * (r / (double)RAND_MAX); }
29             else { *cur = 0.0f; }
30             if (*cur != 0.0f) { nnz++; }
31         }
32     }
33     *result_matrix = A;
34     return nnz;
35 }
36
37 void presentAMatrix(float* mat, int nrow, int ncol, int present_row, int present_col) {
38     printf(" --- PRINT THE MATRIX ---- \n");
39     for (int i = 0; i < present_row; i++) {
40         for (int j = 0; j < present_col; j++) {
41             printf("%3.0f", mat[i * ncol + j]);
42         }
43         printf("...\n");
44     }
45     printf("...\n");
46 }
47
48 void convertToCSRFormat(float* mat, int nrow, int ncol, int nnz, int** ptr, int** indices,
    float** data) {
49     int* row_ptr = (int*)malloc(sizeof(int) * (nrow + 1));
50     int* col_ind = (int*)malloc(sizeof(int) * nnz);
51     float* nz_val = (float*)malloc(sizeof(float) * nnz);
52     float* cur;
53     int count = 0;
54     for (int i = 0; i < nrow; i++) {
55         row_ptr[i] = count;
56         for (int j = 0; j < ncol; j++) {
57             cur = mat + (i * ncol + j);
58             if (*cur != 0.0f) {
59                 col_ind[count] = j;
60                 nz_val[count] = *cur;
61                 count++;
62             }
63         }
64     }
65     row_ptr[nrow] = count;
66 }

```

```

67     *ptr = row_ptr;
68     *indices = col_ind;
69     *data = nz_val;
70     return;
71 }
72
73 void presentCSR(int* ptr, int* indices, float* data, int nnz, int nrow) {
74     printf(" --- PRINT THE CSR FORMAT MATRIX ---- \n");
75     printf("ptr - ");
76     for (int i = 0; i <= nrow; i++) {
77         printf("%+5d", ptr[i]);
78     }
79     printf("\n");
80     printf("indices - ");
81     for (int i = 0; i < nnz; i++) {
82         printf("%+5d", indices[i]);
83     }
84     printf("\n");
85     printf("data - ");
86     for (int i = 0; i < nnz; i++) {
87         printf("%+5f", data[i]);
88     }
89     printf("\n");
90 }
91
92 void print(const thrust::device_vector<int>& v)
93 {
94     for (size_t i = 0; i < v.size(); i++)
95         std::cout << " " << v[i];
96     std::cout << "\n";
97 }
98
99 void print(const thrust::host_vector<int>& v)
100 {
101     for (size_t i = 0; i < v.size(); i++)
102         std::cout << " " << v[i];
103     std::cout << "\n";
104 }
105
106 void print(const thrust::device_vector<float>& v)
107 {
108     for (size_t i = 0; i < v.size(); i++)
109         std::cout << " " << std::fixed << std::setprecision(1) << v[i];
110     std::cout << "\n";
111 }
112
113 void print(const thrust::host_vector<float>& v)
114 {
115     for (size_t i = 0; i < v.size(); i++)
116         std::cout << " " << std::fixed << std::setprecision(1) << v[i];
117     std::cout << "\n";
118 }
119
120 void print(thrust::device_vector<int>& v1, thrust::device_vector<int>& v2)
121 {
122     for (size_t i = 0; i < v1.size(); i++)
123         std::cout << " (" << v1[i] << ", " << std::setw(2) << v2[i] << ")";
124     std::cout << "\n";
125 }
126
127 void printVector(int* counting, int nrow) {
128     for (int i = 0; i < nrow; i++) {
129         std::cout << counting[i] << " ";
130     }
131     std::cout << "\n";
132 }
133
134 void printVector(float* counting, int nrow) {
135     for (int i = 0; i < nrow; i++) {
136         std::cout << counting[i] << " ";
137     }

```

```

138     std::cout << "\n";
139 }
140
141 __device__
142 void printVectorDevice(int* vec, int n) {
143     for (int i = 0; i < n; ++i) {
144         printf("%d ", vec[i]);
145     }
146     printf("\n");
147 }
148
149 __device__
150 void printVectorDevice(float* vec, int n) {
151     for (int i = 0; i < n; ++i) {
152         printf("%.2f ", vec[i]);
153     }
154     printf("\n");
155 }
156
157 __global__
158 void countingKernel(int n, int* d_counting, int* A_ptr, int* A_ind, int* B_ptr) {
159     int idx = blockIdx.x * blockDim.x + threadIdx.x;
160     __syncthreads();
161     if (idx < n) {
162         int start = A_ptr[idx];
163         int end = A_ptr[idx + 1];
164         for (int i = start; i < end; i++) {
165             int j = A_ind[i];
166             d_counting[idx] = d_counting[idx] + B_ptr[j + 1] - B_ptr[j];
167         }
168     }
169     __syncthreads();
170 }
171
172 __device__
173 void bubbleSortSwap(int* col_ind, int i, int j) {
174     int temp = col_ind[i];
175     col_ind[i] = col_ind[j];
176     col_ind[j] = temp;
177 }
178
179 __device__
180 void bubbleSortSwap(float* ety_val, int i, int j) {
181     float temp = ety_val[i];
182     ety_val[i] = ety_val[j];
183     ety_val[j] = temp;
184 }
185
186 __device__
187 void bubbleSortNetwork(int* col_ind, float* ety_val, int n) {
188     for (int i = 0; i < n; i++) {
189         for (int j = i + 1; j < n; j++) {
190             if (col_ind[i] > col_ind[j]) {
191                 bubbleSortSwap(col_ind, i, j);
192                 bubbleSortSwap(ety_val, i, j);
193             }
194         }
195     }
196 }
197
198 __device__
199 void contractionOperation(int* col_ind, float* ety_val, int n) {
200     int tid = threadIdx.x;
201     if (tid < n) {
202         if (tid == 0 || col_ind[tid] != col_ind[tid - 1]) {
203             float res = ety_val[tid];
204             for (int j = tid + 1; j < n && col_ind[j] == col_ind[tid]; j++) {
205                 res += ety_val[j];
206                 col_ind[j] = 0;
207                 ety_val[j] = 0;
208             }

```



```

209         ety_val[tid] = res;
210     }
211 }
212 }
213
214 __global__
215 void expansionSortingContractionKernel(int* order, int* counting, int* operations,
216     int* A_ptr, int* A_ind, float* A_val, int A_row, int A_col, int A_nnz,
217     int* B_ptr, int* B_ind, float* B_val, int B_row, int B_col, int B_nnz,
218     int* phase1_row, int* phase1_col, float* phase1_val,
219     int* phase2_row, int* phase2_col, float* phase2_val,
220     int* resulting_row, int* resulting_col, float* resulting_val) {
221     int bid = blockIdx.x;
222     int tid = threadIdx.x;
223     int block_size = blockDim.x;
224
225     if (bid < A_row) {
226         int row = order[bid];
227         int A_row_start = A_ptr[row];
228         int A_row_end = A_ptr[row + 1];
229         int num_C_row_nnz = counting[bid];
230         __shared__ float C_row_val[768];
231         __shared__ int C_row_col[768];
232         __shared__ int C_row_ind;
233         if (tid == 0) {
234             C_row_ind = 0;
235         }
236         __syncthreads();
237
238         for (int entryIdx = tid + A_row_start; entryIdx < A_row_end; entryIdx +=
239             block_size) {
240             int k = A_ind[entryIdx];
241             int A_ik = A_val[entryIdx];
242             int B_rowk_start = B_ptr[k];
243             int B_rowk_end = B_ptr[k + 1];
244             for (int i = B_rowk_start; i < B_rowk_end; i++) {
245                 int j = B_ind[i];
246                 int B_kj = B_val[i];
247                 int pos = atomicAdd(&C_row_ind, 1);
248                 C_row_val[pos] = A_ik * B_kj;
249                 C_row_col[pos] = j;
250             }
251         }
252
253         for (int resIdx = tid; resIdx < num_C_row_nnz; resIdx += block_size) {
254             phase1_row[operations[bid] + resIdx] = row;
255             phase1_col[operations[bid] + resIdx] = C_row_col[resIdx];
256             phase1_val[operations[bid] + resIdx] = C_row_val[resIdx];
257         }
258         __syncthreads();
259
260         if (num_C_row_nnz <= 32) {
261             if (tid == 0) {
262                 bubbleSortNetwork(C_row_col, C_row_val, num_C_row_nnz);
263             }
264         }
265         else {
266             typedef cub::BlockRadixSort<int, N_THREADS_PER_BLOCK, 6, float> BlockRadixSort
267 ;
268             // Allocate shared memory for BlockRadixSort
269             __shared__ typename BlockRadixSort::TempStorage temp_storage;
270             // Collectively sort the keys
271             BlockRadixSort(temp_storage).Sort(*static_cast<int*>[6]>(<static_cast<void*>(<
272 C_row_col + 6 * threadIdx.x)), *static_cast<float*>[6]>(<static_cast<void*>(<
273 + 6 * threadIdx.x)));
274         }
275
276         for (int resIdx = tid; resIdx < num_C_row_nnz; resIdx += block_size) {
277             phase2_row[operations[bid] + resIdx] = row;

```

```

276         phase2_col[operations[bid] + resIdx] = C_row_col[resIdx];
277         phase2_val[operations[bid] + resIdx] = C_row_val[resIdx];
278     }
279
280     contractionOperation(C_row_col, C_row_val, num_C_row_nnz);
281     for (int resIdx = tid; resIdx < num_C_row_nnz; resIdx += block_size) {
282         if (C_row_val[resIdx] != 0.f) {
283             resulting_row[operations[bid] + resIdx] = row;
284             resulting_col[operations[bid] + resIdx] = C_row_col[resIdx];
285             resulting_val[operations[bid] + resIdx] = C_row_val[resIdx];
286         }
287     }
288 }
289 __syncthreads();
290 }
291
292 struct removed_item {
293     template <typename Tuple>
294     __host__ __device__
295     bool operator()(const Tuple& t) {
296         return (thrust::get<2>(t) == 0);
297     }
298 };
299
300 int main() {
301
302     int A_row = 7;
303     int A_col = 7;
304     int B_row = 7;
305     int B_col = 7;
306     int present_row = 7;
307     int present_col = 7;
308     float* A;
309     float* B;
310     int A_nnz = generateASparseMatrixRandomly(A_row, A_col, &A, 4);
311     int B_nnz = generateASparseMatrixRandomly(B_row, B_col, &B, 4);
312     presentAMatrix(A, A_row, A_col, present_row, present_col);
313     presentAMatrix(B, B_row, B_col, present_row, present_col);
314     int* A_ptr;
315     int* A_ind;
316     float* A_val;
317     convertToCSRFormat(A, A_row, A_col, A_nnz, &A_ptr, &A_ind, &A_val);
318     presentCSR(A_ptr, A_ind, A_val, A_nnz, A_row);
319     int* B_ptr;
320     int* B_ind;
321     float* B_val;
322     convertToCSRFormat(B, B_row, B_col, B_nnz, &B_ptr, &B_ind, &B_val);
323     presentCSR(B_ptr, B_ind, B_val, B_nnz, B_row);
324
325     int* d_A_ptr, * d_A_ind, * d_B_ptr, * d_B_ind;
326     float* d_A_val, * d_B_val;
327     cudaMalloc(&d_A_ptr, (A_row + 1) * sizeof(int));
328     cudaMalloc(&d_A_ind, A_nnz * sizeof(int));
329     cudaMalloc(&d_A_val, A_nnz * sizeof(float));
330     cudaMalloc(&d_B_ptr, (B_row + 1) * sizeof(int));
331     cudaMalloc(&d_B_ind, B_nnz * sizeof(int));
332     cudaMalloc(&d_B_val, B_nnz * sizeof(float));
333     cudaMemcpy(d_A_ptr, A_ptr, (A_row + 1) * sizeof(int), cudaMemcpyHostToDevice);
334     cudaMemcpy(d_A_ind, A_ind, A_nnz * sizeof(int), cudaMemcpyHostToDevice);
335     cudaMemcpy(d_A_val, A_val, A_nnz * sizeof(float), cudaMemcpyHostToDevice);
336     cudaMemcpy(d_B_ptr, B_ptr, (B_row + 1) * sizeof(int), cudaMemcpyHostToDevice);
337     cudaMemcpy(d_B_ind, B_ind, B_nnz * sizeof(int), cudaMemcpyHostToDevice);
338     cudaMemcpy(d_B_val, B_val, B_nnz * sizeof(float), cudaMemcpyHostToDevice);
339
340     int* counting = (int*)malloc(A_row * sizeof(int));
341     int* d_counting;
342     cudaMalloc(&d_counting, A_row * sizeof(int));
343     cudaMemset(d_counting, 0, A_row * sizeof(int));
344     int set_up_n_blocks = (A_row + N_THREADS_PER_BLOCK - 1) / N_THREADS_PER_BLOCK;
345

```

```

347 countingKernel <<< set_up_n_blocks, N_THREADS_PER_BLOCK >>> (A_row, d_counting,
348   d_A_ptr, d_A_ind, d_B_ptr);
349 cudaDeviceSynchronize();
350 cudaMemcpy(counting, d_counting, A_row * sizeof(int), cudaMemcpyDeviceToHost);
351 std::cout << " --- SET UP --- " << "\n";
352 std::cout << "original counting: ";
353 printVector(counting, A_row);
354
355 thrust::device_ptr<int> d_counting_ptr(d_counting);
356 thrust::device_vector<int> d_counting_vec(d_counting_ptr, d_counting_ptr + A_row);
357 thrust::device_vector<int> d_order_vec(A_row);
358 thrust::sequence(d_order_vec.begin(), d_order_vec.end());
359 thrust::stable_sort_by_key(d_counting_vec.begin(), d_counting_vec.end(), d_order_vec.
   begin(), thrust::greater<int>());
360 thrust::device_vector<int> d_operations_vec(A_row + 1);
361 thrust::exclusive_scan(d_counting_vec.begin(), d_counting_vec.end(), d_operations_vec.
   begin());
362 int tot_operations = d_counting_vec.back() + d_operations_vec[A_row - 1];
363 d_operations_vec[A_row] = tot_operations;
364
365 std::cout << " --- AFTER REORDER - THRUST VECTOR " << "\n";
366 std::cout << "d_counting_vec: ";
367 print(d_counting_vec);
368 std::cout << "d_order_vec: ";
369 print(d_order_vec);
370 std::cout << "d_operations_vec: ";
371 print(d_operations_vec);
372
373 int* order = (int*)malloc(A_row * sizeof(int));
374 int* operations = (int*)malloc((A_row + 1) * sizeof(int));
375 thrust::copy(d_order_vec.begin(), d_order_vec.end(), order);
376 thrust::copy(d_operations_vec.begin(), d_operations_vec.end(), operations);
377 thrust::copy(d_counting_vec.begin(), d_counting_vec.end(), counting);
378 std::cout << " --- AFTER REORDER - CONVERT TO ptr " << "\n";
379 std::cout << "order: ";
380 printVector(order, A_row);
381 std::cout << "counting: ";
382 printVector(counting, A_row);
383 std::cout << "operations: ";
384 printVector(operations, A_row + 1);
385 int* d_order;
386 int* d_operations;
387 cudaMalloc(&d_order, A_row * sizeof(int));
388 cudaMalloc(&d_operations, (A_row + 1) * sizeof(int));
389 cudaMemcpy(d_order, order, A_row * sizeof(int), cudaMemcpyHostToDevice);
390 cudaMemcpy(d_operations, operations, (A_row + 1) * sizeof(int), cudaMemcpyHostToDevice);
391 cudaMemcpy(d_counting, counting, A_row * sizeof(int), cudaMemcpyHostToDevice);
392
393 int* d_resulting_row, * d_resulting_col;
394 float* d_resulting_val;
395 cudaMalloc(&d_resulting_row, tot_operations * sizeof(int));
396 cudaMalloc(&d_resulting_col, tot_operations * sizeof(int));
397 cudaMalloc(&d_resulting_val, tot_operations * sizeof(float));
398
399 int* d_phase1_row, * d_phase1_col;
400 float* d_phase1_val;
401 cudaMalloc(&d_phase1_row, tot_operations * sizeof(int));
402 cudaMalloc(&d_phase1_col, tot_operations * sizeof(int));
403 cudaMalloc(&d_phase1_val, tot_operations * sizeof(float));
404
405 int* d_phase2_row, * d_phase2_col;
406 float* d_phase2_val;
407 cudaMalloc(&d_phase2_row, tot_operations * sizeof(int));
408 cudaMalloc(&d_phase2_col, tot_operations * sizeof(int));
409 cudaMalloc(&d_phase2_val, tot_operations * sizeof(float));
410
411 int expansion_and_sorting_n_blocks = A_row;
412 expansionSortingContractionKernel <<< expansion_and_sorting_n_blocks,
   N_THREADS_PER_BLOCK >>> (d_order, d_counting, d_operations,

```

```

413     d_A_ptr, d_A_ind, d_A_val, A_row, A_col, A_nnz,
414     d_B_ptr, d_B_ind, d_B_val, B_row, B_col, B_nnz,
415     d_phase1_row, d_phase1_col, d_phase1_val,
416     d_phase2_row, d_phase2_col, d_phase2_val,
417     d_resulting_row, d_resulting_col, d_resulting_val);
418     cudaDeviceSynchronize();
419     int* resulting_row = (int*)malloc(tot_operations * sizeof(int));
420     int* resulting_col = (int*)malloc(tot_operations * sizeof(int));
421     float* resulting_val = (float*)malloc(tot_operations * sizeof(float));
422     cudaMemcpy(resulting_row, d_resulting_row, tot_operations * sizeof(int),
423               cudaMemcpyDeviceToHost);
424     cudaMemcpy(resulting_col, d_resulting_col, tot_operations * sizeof(int),
425               cudaMemcpyDeviceToHost);
426     cudaMemcpy(resulting_val, d_resulting_val, tot_operations * sizeof(float),
427               cudaMemcpyDeviceToHost);
428
429     int* phase1_row = (int*)malloc(tot_operations * sizeof(int));
430     int* phase1_col = (int*)malloc(tot_operations * sizeof(int));
431     float* phase1_val = (float*)malloc(tot_operations * sizeof(float));
432     cudaMemcpy(phase1_row, d_phase1_row, tot_operations * sizeof(int),
433               cudaMemcpyDeviceToHost);
434     cudaMemcpy(phase1_col, d_phase1_col, tot_operations * sizeof(int),
435               cudaMemcpyDeviceToHost);
436     cudaMemcpy(phase1_val, d_phase1_val, tot_operations * sizeof(float),
437               cudaMemcpyDeviceToHost);
438
439     int* phase2_row = (int*)malloc(tot_operations * sizeof(int));
440     int* phase2_col = (int*)malloc(tot_operations * sizeof(int));
441     float* phase2_val = (float*)malloc(tot_operations * sizeof(float));
442     cudaMemcpy(phase2_row, d_phase2_row, tot_operations * sizeof(int),
443               cudaMemcpyDeviceToHost);
444     cudaMemcpy(phase2_col, d_phase2_col, tot_operations * sizeof(int),
445               cudaMemcpyDeviceToHost);
446     cudaMemcpy(phase2_val, d_phase2_val, tot_operations * sizeof(float),
447               cudaMemcpyDeviceToHost);
448
449     std::cout << " --- AFTER EXPANSION and BEFORE SORTING" << "\n";
450     std::cout << "phase1_row: ";
451     printVector(phase1_row, tot_operations);
452     std::cout << "phase1_col: ";
453     printVector(phase1_col, tot_operations);
454     std::cout << "phase1_val: ";
455     printVector(phase1_val, tot_operations);
456
457     std::cout << " --- AFTER SORTING and BEFORE CONTRACTION" << "\n";
458     std::cout << "phase2_row: ";
459     printVector(phase2_row, tot_operations);
460     std::cout << "phase2_col: ";
461     printVector(phase2_col, tot_operations);
462     std::cout << "phase2_val: ";
463     printVector(phase2_val, tot_operations);
464
465     std::cout << " --- FINALLY --- " << "\n";
466     std::cout << "resulting_row: ";
467     printVector(resulting_row, tot_operations);
468     std::cout << "resulting_col: ";
469     printVector(resulting_col, tot_operations);
470     std::cout << "resulting_val: ";
471     printVector(resulting_val, tot_operations);
472
473     std::cout << " --- COO FORMAT --- " << "\n";
474     thrust::host_vector<int> coo_row(resulting_row, resulting_row + tot_operations);
475     thrust::host_vector<int> coo_col(resulting_col, resulting_col + tot_operations);
476     thrust::host_vector<float> coo_val(resulting_val, resulting_val + tot_operations);
477
478     thrust::device_vector<int> d_coo_row = coo_row;
479     thrust::device_vector<int> d_coo_col = coo_col;
480     thrust::device_vector<float> d_coo_val = coo_val;
481
482     typedef thrust::device_vector<int>::iterator IntIterator;
483     typedef thrust::device_vector<float>::iterator FloatIterator;

```

```

475 typedef thrust::tuple<IntIterator, IntIterator, FloatIterator> IteratorTuple;
476 typedef thrust::zip_iterator<IteratorTuple> ZipIterator;
477
478 ZipIterator first = thrust::make_zip_iterator(thrust::make_tuple(d_coo_row.begin(),
479     d_coo_col.begin(), d_coo_val.begin()));
479 ZipIterator last = thrust::make_zip_iterator(thrust::make_tuple(d_coo_row.end(),
480     d_coo_col.end(), d_coo_val.end()));
480
481 ZipIterator new_last = thrust::remove_if(first, last, removed_item());
482
483 d_coo_row.erase(new_last.get_iterator_tuple().get<0>(), d_coo_row.end());
484 d_coo_col.erase(new_last.get_iterator_tuple().get<1>(), d_coo_col.end());
485 d_coo_val.erase(new_last.get_iterator_tuple().get<2>(), d_coo_val.end());
486
487 thrust::host_vector<int> h_coo_row = d_coo_row;
488 thrust::host_vector<int> h_coo_col = d_coo_col;
489 thrust::host_vector<float> h_coo_val = d_coo_val;
490
491 std::cout << "coo_row: ";
492 print(h_coo_row);
493 std::cout << "coo_col: ";
494 print(h_coo_col);
495 std::cout << "coo_val: ";
496 print(h_coo_val);
497 return 0;
498 }

```