

Small Scale Parallel Programming Assignment

Aritra Das Ray s426160

February 2024

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 2 |
| 1.1 | Aim | 2 |
| 1.2 | Objective | 2 |
| 2 | Methodology | 2 |
| 3 | Algorithm and Code Description | 4 |
| 3.1 | Main function | 4 |
| 3.2 | Reading .MTX file and CSR conversion | 5 |
| 3.3 | Generating Vectors based on input columns | 6 |
| 3.4 | Multiplication with OpenMP | 7 |
| 3.5 | Multiplication with CUDA | 8 |
| 4 | Analysis | 9 |
| 4.1 | Functional Test | 9 |
| 4.2 | Performance Tests | 10 |
| A | Appendix | 13 |
| A.1 | C Scripts | 13 |
| A.1.1 | Matrix Vector Multiplication CUDA | 13 |
| A.1.2 | Matrix Vector Multiplication CUDA Performance | 18 |
| A.1.3 | Matrix Vector Multiplication OpenMP | 24 |
| A.1.4 | Matrix Vector Multiplication OpenMP Performance | 29 |

1 Introduction

1.1 Aim

The aim of the work described in this report is to explore Sparse Matrix and their Multiplication with Vectors of different sizes OpenMP and CUDA parallel programming paradigms and analyse their performance.

1.2 Objective

To Achieve our Aims we breakdown the tasks into the following steps:

1. Download matrices in Matrix market Format form <https://sparse.tamu.edu/>
2. Convert the matrices to CSR format
3. Generate vectors having varying columns.
4. Implement matrix vector multiplication in OpenMP and CUDA programming paradigms.
5. Measure time for execution of each operation to analyse performance.

2 Methodology

Sparse Matrix:

In the field of Machine Learning, working with large matrices is common. However, these matrices often contain a significant number of 0s, leading to what is known as a Sparse Matrix. As these 0s do not convey useful information and occupy considerable space, it becomes inefficient to compute them. To address this issue, various storage formats are employed to significantly reduce storage and computing needs. Depending on the distribution and quantity of non-zero entries, different data structures can be utilized for substantial memory savings compared to standard approaches. This trade-off requires more complex access to individual elements and necessitates additional structures for unambiguous recovery of the original matrix.[1]

Sparse Matrix Formats

Sparse Matrix can be stored in several different formats, following methods are the widely used formats:

- Dictionary of Keys (DOK)
- List of Lists
- Coordinated List
- Compressed Sparse Row (CSR)

- Compressed Sparse Column (CSC)

We will utilize the Compressed Sparse Row Method for our project as it allows for effective access and manipulation of matrices.[2]

CSR

The compressed row storage, also known as the compressed sparse row or Yale format, is a method for representing a matrix M using three one-dimensional arrays. These arrays store the nonzero values, row extents, and column indices of the matrix.[2]

The matrix M with dimensions $m \times n$ and nnz non-zero elements is represented using three one-dimensional arrays (V , column index, row index). The array V stores the non-zero values, the column index array stores their respective column positions and both arrays are of equal length. Meanwhile, the row index array has a length of $m+1$ and indicates the starting index of each row in V and the column index. The last element in the row index represents nnz which acts as a placeholder for proper loop iteration.[2]

Parallel Programming

Parallel programming involves executing multiple calculations or processes simultaneously to improve computational efficiency and reduce the time required for running complex tasks. This approach is particularly relevant in the context of CPU (Central Processing Unit) and GPU (Graphics Processing Unit) computing, where different architectures and capabilities are leveraged for optimal performance. Parallel programming with multi-threaded shared memory is a key strategy for maximizing the utilization of these hardware resources.[5]

CPU Parallel Programming

CPUs are general-purpose processors designed to handle a wide range of computing tasks. They are composed of a few cores with high clock speeds capable of executing a variety of operations. Parallel programming on CPUs often involves creating multiple threads within an application, where each thread can run concurrently on different cores, allowing for multitasking within the same application. This is achieved through the use of shared memory, where multiple threads can access and modify the same memory space, enabling efficient communication and data sharing among them.[4]

The main challenge in CPU parallel programming is managing the synchronization between threads to ensure data consistency and prevent race conditions, where two or more threads attempt to modify the same data simultaneously. Tools and libraries such as OpenMP (Open Multi-Processing) and Pthreads (POSIX threads) are commonly used for facilitating multi-threaded programming on CPUs by providing APIs for thread creation, synchronization, and management.[4]

GPU Parallel Programming

GPUs, on the other hand, are specialized processors designed to handle the

massive parallelism required for rendering graphics and video processing. Unlike CPUs, GPUs consist of thousands of smaller cores designed to perform simple calculations on large data sets simultaneously. This makes them highly efficient for algorithms that can be parallelized at a fine-grained level.

Parallel programming on GPUs involves dividing a task into small work items that can be executed simultaneously across many GPU cores. This is achieved through the use of a different programming model and memory hierarchy, including local (private), global, and shared memory spaces. Shared memory on a GPU is a small but fast memory accessible by threads within the same thread block, facilitating fast data exchange and synchronization among these threads.[3]

CUDA (Compute Unified Device Architecture) by NVIDIA and OpenCL (Open Computing Language) are popular frameworks for GPU programming, allowing developers to write programs that execute across both CPUs and GPUs. These frameworks provide the necessary tools for managing memory, defining parallel kernels (functions), and controlling thread execution.[3]

Multi-Threaded Shared Memory Model

The multi-threaded shared memory model in parallel programming exploits the shared memory architecture of both CPUs and GPUs to allow multiple threads to access common memory locations. This model is critical for achieving high performance in computing tasks that require frequent data exchange and synchronization between threads. Efficient use of shared memory can significantly reduce the latency associated with memory accesses and improve the overall speed of parallel algorithms.[9]

However, the efficient design of parallel programs that use shared memory requires careful consideration of memory access patterns, synchronization mechanisms, and the avoidance of bottlenecks that can arise from contention over shared resources. Techniques such as lock-free algorithms, atomic operations, and barrier synchronization are often employed to manage these challenges.[9]

3 Algorithm and Code Description

3.1 Main function

The general architecture of the "main" function is as follows:

- **Initialization:** The function start by defining a Sparse_CSR structure for storing the matrix in Compressed Sparse Row format. It also initialize a seed for random number generation, used later for vector generation.
- **Matrix Reading:** Then it reads a sparse matrix from the .mtx file using read_mtx_file function and stores it in CSR format within csr_matrix Struct.

- **Vector Generation:** Vectors of varying sizes (`vec_k1`, `vec_k2`, `vec_k3`, `vec_k6`) are generated using `vec_gen`. The sizes are based on the number of columns in the CSR matrix and a scaling factor (1, 2, 3, 6), signifying the number of columns of the vector with the third parameter as the seed for randomness.
- **Matrix Vector Multiplication:** The matrix is multiplied to the vectors generated in the previous steps using parallelization with OpenMP and CUDA and the results are stored in `result_k1`, `result_k2`, `result_k3`, `result_k6`. The OpenMP parallelization is achieved using the `multiplication_openmp`, while CUDA parallelization is achieved through `multiplication_cuda`.
- **Printing Results:** Results are printed by looping through the result arrays.
- **Freeing memory:** The memory is freed from `csr_matrix`, generated vectors and results of multiplication to prevent memory leaks.

3.2 Reading .MTX file and CSR conversion

The `read_mtx_file` function is designed to read a sparse matrix from a .mtx file and convert it into Compressed Sparse Row (CSR) format. This process involves several steps, each crucial for efficiently representing sparse matrices in memory. The function's operation can be summarized as follows:

Parameters

- **const char *filename:** A string that represents the path to the .mtx file to be read.
- **Sparse_CSR csr_matrix:** A pointer to a `Sparse_CSR` structure where the matrix in CSR format will be stored.

Process and Algorithms

1. **Opening the File:** Opens the .mtx file for reading. If the file cannot be opened, it reports an error and exits the program.
2. **Skipping Comments:** Reads lines from the file until it encounters a line that does not start with `%`, indicating the end of comments and the start of matrix dimensions.
3. **Reading Matrix Dimensions:** Reads the first non-comment line to extract the number of rows, columns, and non-zero elements in the matrix.
4. **Allocating Memory for COO:** Allocates memory for arrays to temporarily store the matrix in Coordinate (COO) format (row indices, column indices, and values of non-zero elements).

5. **Reading Matrix Entries:** Reads the file's entries (row index, column index, value) into the COO format arrays.
6. **Initializing CSR Structure:** Stores the read dimensions in the `csr_matrix` structure and allocates memory for CSR arrays (`row_ptr`s, `col_indices`, `values`).
7. **Converting to CSR Format:**
 - **Counting Non-Zeros per Row:** Iterates over the COO arrays to count the number of non-zero elements in each row, storing these counts in the `row_ptr`s array.
 - **Computing Row Pointer Indices:** Transforms the counts in `row_ptr`s into starting indices for each row in the CSR arrays.
 - **Filling CSR Arrays:** Iterates over the COO arrays again to fill the `col_indices` and `values` arrays, adjusting the `row_ptr`s indices as needed.
8. **Freeing Temporary Memory:** Releases the memory allocated for the COO format arrays.

Output

The function outputs by filling the passed `Sparse_CSR` structure with the matrix's CSR representation, including:

- The number of rows (`n_row`), columns (`n_cols`), and non-zero elements (`n_nz`).
- Arrays for the row pointers (`row_ptr`s), column indices (`col_indices`), and values of non-zero elements (`values`).

3.3 Generating Vectors based on input columns

The `vec_gen` function generates a matrix or vector with specified dimensions and fills it with random integer values. It's designed to create test vectors or matrices for operations such as matrix-vector multiplication, with the ability to control the randomness for reproducibility through a seed value. Here's a breakdown of its operation:

Parameters

- **int `vec_row`:** The number of rows in the vector or matrix to be generated.
- **int `vec_col`:** The number of columns in the vector or matrix. When `vec_col` is 1, the function generates a traditional vector. Otherwise, it generates a matrix.
- **unsigned int `seed`:** A seed value for the random number generator, ensuring reproducible results across different runs when the same seed is used.

Algorithm

1. **Initialization of Random Number Generator:** Uses the seed parameter to initialize the standard C library's random number generator with `srand`, ensuring that subsequent calls to `rand` produce a sequence of numbers that can be replicated by using the same seed.
2. **Memory Allocation:** Allocates memory for the vector or matrix based on the specified dimensions (`vec_row * vec_col`). It calculates the total number of elements and allocates enough space to store `int` values for each.
3. **Filling with Random Values:** Iterates over each element in the allocated space, assigning it a random integer value between 0 and 9 (inclusive). This is done by using the `rand` function and taking the remainder when divided by 10 (`rand() % 10`), ensuring that all values fall within the specified range.
4. **Error Handling:** Checks if memory allocation failed (`if(vector == NULL)`). If so, it returns `NULL` to indicate an error, allowing the calling function to handle the situation appropriately.

Output

Returns a pointer to the first element of the allocated and filled vector or matrix. This pointer can then be used in further computations or for debugging purposes (e.g., printing the vector/matrix).

3.4 Multiplication with OpenMP

The `multiplication_openmp` function performs matrix-vector multiplication using the Compressed Sparse Row (CSR) format for the matrix and supports parallel computation with OpenMP. This function is designed to work efficiently on multi-core processors by leveraging OpenMP to distribute the workload among multiple threads.

Parameters

- **const Sparse_CSR* csr_matrix:** A pointer to the sparse matrix in CSR format, which includes arrays for row pointers, column indices, and non-zero values, along with the matrix's dimensions.
- **const int* vector:** A pointer to the vector (or matrix when `vec_col > 1`) to be multiplied with the CSR matrix.
- **int vec_col:** The number of columns in the vector. This allows the function to support both traditional vector multiplication (when `vec_col = 1`) and matrix multiplication (when `vec_col > 1`).

Algorithm

1. **Memory Allocation for Result:** Allocates memory for the result of the multiplication. The size is determined by the number of rows in the CSR matrix and the number of columns in the vector.
2. **Parallel Computation:**
 - Retrieves the number of threads to be used from the environment variable `OMP_NUM_THREADS` and parses it to an integer.
 - Uses OpenMP directives (`#pragma omp parallel for`) to parallelize the outer loops over the vector's columns and the CSR matrix's rows.
 - Initializes private variables for each thread to avoid race conditions. These variables include loop indices and temporary variables for intermediate calculations.
 - Each thread calculates a portion of the result matrix, iterating through the non-zero elements of the CSR matrix, multiplying each by the corresponding vector value, and summing these products to form the result elements.
3. **OpenMP Barrier:** Ensures that all threads have completed their portion of the work before proceeding. This is implicitly provided by the end of the `#pragma omp parallel` for region.

Output

Returns a pointer to a dynamically allocated array of doubles that contains the result of the matrix-vector multiplication. The result is stored in row-major order if `vec_col > 1`, treating the vector as a matrix.

3.5 Multiplication with CUDA

The `multiplication_cuda` function performs matrix-vector multiplication leveraging CUDA for GPU acceleration, suitable for sparse matrices in Compressed Sparse Row (CSR) format. It demonstrates efficient GPU utilization for handling large-scale computations more quickly than CPU-based methods.

Parameters

- **const Sparse_CSR* csr_matrix:** A pointer to the sparse matrix in CSR format, which includes arrays for row pointers, column indices, and non-zero values, along with the matrix's dimensions.
- **const int* vector:** A pointer to the vector (or matrix when `vec_col > 1`) to be multiplied with the CSR matrix.
- **int vec_col:** The number of columns in the vector. This allows the function to support both traditional vector multiplication (when `vec_col = 1`) and matrix multiplication (when `vec_col > 1`).

Algorithm

1. **Memory Allocation for Result:** Allocates memory for the result of the multiplication. The size is determined by the number of rows in the CSR matrix and the number of columns in the vector.
2. **Memory Allocation on GPU:** Allocates GPU memory for the matrix's CSR components (`d_row_ptr`, `d_col_ind`, `d_values`), the input vector (`d_vec`), and the result vector (`d_result`).
3. **Data Transfer to GPU:** Copies the CSR matrix components and the input vector from host (CPU) memory to the allocated GPU memory.
4. **Kernel Initialization and Execution:**
 - Initializes CUDA events (start, stop) for timing the kernel execution.
 - Calculates the number of blocks and threads per block needed for the kernel launch, ensuring efficient GPU utilization.
 - Launches the `matrix_vector_multiplication` kernel, passing the necessary parameters for the computation.
 - Checks for kernel launch errors using `cudaGetLastError`.
5. **Data Transfer Back to Host:** Copies the result vector from GPU memory back to host memory.
6. **Timing and Cleanup:**
 - Records the elapsed time for the kernel execution using CUDA events.
 - Frees the allocated GPU memory to avoid memory leaks.

CUDA Kernel: `matrix_vector_multiplication`:

- Executed on the GPU, it calculates the product of the sparse matrix and the vector for each row. For matrices with multiple columns (`vec_col > 1`), it processes each column of the input vector in parallel.
- Each thread computes one or more elements of the result vector, iterating over the non-zero elements of its assigned rows and calculating the dot product.

Output

Returns a pointer to the result vector (`double*`), which contains the outcome of the matrix-vector multiplication. The memory for the result vector is dynamically allocated and should be freed by the caller to prevent memory leaks.

4 Analysis

4.1 Functional Test

For Functional testing, the results obtained from the OpenMP and CUDA code has been verified against an analytical solution. The analytical solution has

been obtained from matrixcalc.org. The analytical Solution is shared in a pdf file along with the code files for verification Based on the figures it can be

```
/var/spool/pbs/mom_priv/jobs/38566.cr2-pbs.SC: line 52:
export: `/var/spool/pbs/aux/38566.cr2-pbs': not a valid
identifier
0
row = 9 , col = 9 , nz = 49
result_k1 : [4.274917, 5.926079, 6.500000, 5.890643,
5.562708, 4.550664, 5.382614, 5.070736, 2.841639, ]
result_k2 : [5.462209, 6.225083, 6.276024, 5.313538,
3.274917, 5.258306, 5.336794, 4.482282, 5.487680, 4.425249,
2.250277, 5.004845, 1.437292, 6.240864, 2.324751, 7.116556,
1.150055, 2.933278, ]
result_k3 : [3.950166, 7.162375, 8.150332, 5.013427,
3.413206, 5.876245, 5.233112, 2.208472, 7.608250, 2.836794,
3.749170, 7.944491, 1.849945, 5.850498, 3.500000, 4.079596,
1.429540, 4.517165, 2.303848, 4.441030, 2.537237, 6.099668,
4.862542, 6.991141, 1.633444, 1.883167, 1.875138, ]
result_k6 : [
3.687292, 6.463040, 7.324751, 5.987957, 5.175249, 6.862542,
5.513427, 3.063261, 7.426079, 2.863372, 3.138289, 7.026300,
1.749723, 7.266888, 3.158084, 5.575028, 2.208472, 6.183278,
2.237126, 7.811047, 2.286960, 8.032116, 7.348838, 8.623755,
4.237403, 1.250277, 2.812708, 4.662375, 6.450166, 4.437569,
5.263150, 2.071013, 2.341916, 5.538344, 5.379983, 4.392026,
2.678710, 4.212209, 1.720515, 3.482835, 6.057586, 2.799834,
3.599668, 6.412375, 3.687292, 7.241141, 7.324751, 4.124585,
2.033499, 1.449889, 1.241694, 2.616833, 2.916667,
1.550111, ]
```

Figure 1: Results from OpenMP code

```
CUDA_VISIBLE_DEVICES=GPU-2b6c369f-e343-b816-fe19-1200242394f1
row = 9 , col = 9 , nz = 49
threadspblock = 1024, numblock = 1Time elapsed for executing with vector columns
1 = 0.043000
threadspblock = 1024, numblock = 1Time elapsed for executing with vector columns
2 = 0.022752
threadspblock = 1024, numblock = 1Time elapsed for executing with vector columns
3 = 0.020400
threadspblock = 1024, numblock = 1Time elapsed for executing with vector columns
6 = 0.020992
result_k1
4.274917, 5.926079, 6.500000, 5.890643, 5.562708, 4.550664, 5.382614, 5.070736,
2.841639,
result_k2
5.462209, 6.225083, 6.276024, 5.313538, 3.274917, 5.258306, 5.336794, 4.482282,
5.487680, 4.425249, 2.250277, 5.004845, 1.437292, 6.240864, 2.324751, 7.116556,
1.150055, 2.933278,
result_k3
3.950166, 7.162375, 8.150332, 5.013427, 3.413206, 5.876245, 5.233112, 2.208472,
7.608250, 2.836794, 3.749170, 7.944491, 1.849945, 5.850498, 3.500000, 4.079596,
1.429540, 4.517165, 2.303848, 4.441030, 2.537237, 6.099668, 4.862542, 6.991141,
1.633444, 1.883167, 1.875138,
result_k6
3.687292, 6.463040, 7.324751, 5.987957, 5.175249, 6.862542, 5.513427, 3.063261,
7.426079, 2.863372, 3.138289, 7.026300, 1.749723, 7.266888, 3.158084, 5.575028,
2.208472, 6.183278, 2.237126, 7.811047, 2.286960, 8.032116, 7.348838, 8.623755,
4.237403, 1.250277, 2.812708, 4.662375, 6.450166, 4.437569, 5.263150, 2.071013,
2.341916, 5.538344, 5.379983, 4.392026, 2.678710, 4.212209, 1.720515, 3.482835,
6.057586, 2.799834, 3.599668, 6.412375, 3.687292, 7.241141, 7.324751, 4.124585,
2.033499, 1.449889, 1.241694, 2.616833, 2.916667, 1.550111,
```

Figure 2: Results from CUDA code

concluded that both the OpenMP and CUDA implementations give the same results. On comparison with the Analytical solution, it has been verified that the solutions are correct.

4.2 Performance Tests

The execution times of computing the matrix multiplication has been captured for 15 matrices. The data obtained has been shared along with the code files. However, Here are some key findings from the data:

1. For certain matrices and vector column configurations, CUDA shows either superior or comparable performance to OpenMP. This can be attributed to CUDA's efficient utilization of GPU resources, which can significantly accelerate computations for parallelizable tasks.
2. As the number of kernels increases from 1 to 8, there's a general trend of decreasing average execution time, suggesting that using more kernels can lead to faster execution times up to a certain point.
3. the reduction in execution time does not always occur linearly with the increase in the number of kernels, indicating that there might be an optimal number of kernels beyond which the performance gain diminishes or becomes inconsistent due to overhead or other factors.

References

- [1] Jason Brownlee. A Gentle Introduction to Sparse Matrices for Machine Learning[Internet]. Guiding Tech Media; 2019 Aug 9 [cited 2024 Feb 25]. Available from: <https://machinelearningmastery.com/sparse-matrices-for-machine-learning/>
- [2] Wikipedia contributors. Sparse matrix [Internet]. Wikipedia, The Free Encyclopedia; 2024 Feb 8, 18:29 [cited 2024 Feb 25]. Available from: https://en.wikipedia.org/w/index.php?title=Sparse_matrix&oldid=1205044030.
- [3] NVIDIA. CUDA Quick Start Guide [Internet]. California, USA: NVIDIA; 2023 Nov 14 [cited 2024 Feb 25]. Available from: <https://docs.nvidia.com/cuda/cuda-quick-start-guide/index.html>
- [4] Wikipedia contributors. OpenMP [Internet]. Wikipedia, The Free Encyclopedia; 2024 Jan 28, 06:11 UTC [cited 2024 Feb 25]. Available from: <https://en.wikipedia.org/w/index.php?title=OpenMP&oldid=1199889263>.
- [5] Wikipedia contributors. Parallel computing [Internet]. Wikipedia, The Free Encyclopedia; 2024 Feb 1, 00:45 UTC [cited 2024 Feb 25]. Available from: https://en.wikipedia.org/w/index.php?title=Parallel_computing&oldid=1201602893.
- [6] Szorbasc. mat_mult.cu in CUDA-CSR-matrix-vector multiplication [Internet]. GitHub; [cited 2024 Feb 25]. Available from: https://github.com/szorbasc/CUDA-CSR-matrix-vector-multiplication/blob/master/mat_mult.cu
- [7] Szorbasc. mat_mult.c in OpenMP-CSR-matrix-vector multiplication [Internet]. GitHub; [cited 2024 Feb 25]. Available from: https://github.com/szorbasc/OpenMP-CSR-matrix-vector-multiplication/blob/master/mat_mult.c
- [8] OpenMP Architecture Review Board. OpenMP Application Programming Interface. [Internet]. Gemini Dr, USA: 2024 Jan 28, 06:11 UTC [cited 2024 Feb 25]. Available from: <https://www.openmp.org/wp-content/uploads/OpenMP-API-Specification-5-2.pdf>
- [9] Herlihy M, Shavit N. The Art of Multiprocessor Programming. San Francisco: Morgan Kaufmann; 2011.[cited 2024 Feb 25]

A Appendix

A.1 C Scripts

A.1.1 Matrix Vector Multiplication CUDA

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <stdbool.h>
4  #include <time.h>
5  #include <string.h>
6  #include <cuda_runtime.h>
7
8  typedef struct Sparse_CSR {
9      int n_row;
10     int n_cols;
11     int n_nz;
12     int* row_ptrs;
13     int* col_indices;
14     double* values;
15 } Sparse_CSR;
16
17 void read_mtx_file(const char *filename, Sparse_CSR *csr_matrix){
18
19     //open file
20     FILE *file = fopen(filename, "r");
21     if(!file) {
22         perror("Unable to open file!");
23         exit(EXIT_FAILURE);
24     }
25
26     //skip comments
27     char line[256];
28     while (fgets(line, sizeof(line), file)) {
29         if (line[0] != '%') break;
30     }
31
32     //initialize variable to store coo matrix from file
33     int coo_rown, coo_coln, coo_nnz;
34
35     //input no of rows, columns and non-zeros
36     sscanf(line, "%d %d %d", &coo_rown, &coo_coln, &coo_nnz);
37     printf("row = %d , col = %d , nz = %d\n", coo_rown, coo_coln,
38           coo_nnz);
39
40     //allocate memory to store the coo values
41     int* coo_rows = (int *)malloc(coo_nnz * sizeof(int));
42     int* coo_cols = (int *)malloc(coo_nnz * sizeof(int));
43     double* coo_val = (double *)malloc(coo_nnz * sizeof(double));
44
45     if (!coo_rows || !coo_cols || !coo_val){
46         perror("Memory Allocation failed!");
47         fclose(file);
48         exit(EXIT_FAILURE);
49     }
50
51     for(int i = 0; i < coo_nnz ; i++){
```

```

51     fscanf(file, "%d %d %lf" , &coo_rows[i], &coo_cols[i], &coo_val
52         [i]);
53 }
54 //variable assignment
55 csr_matrix->n_row = coo_rown;
56 csr_matrix->n_cols = coo_coln;
57 csr_matrix->n_nz = coo_nnz;
58
59
60 // Allocate CSR format arrays
61 csr_matrix->row_ptrs = (int *)malloc((coo_rown + 1) * sizeof(int)
62 );
63 memset(csr_matrix->row_ptrs, 0, (coo_rown + 1) * sizeof(int)); //
64     Initialize to 0
65
66 csr_matrix->col_indices = (int *)malloc(coo_nnz * sizeof(int));
67 csr_matrix->values = (double *)malloc(coo_nnz * sizeof(double));
68
69 // Step 1: Count non-zero elements per row
70 for (int i = 0; i < coo_nnz; i++) {
71     csr_matrix->row_ptrs[coo_rows[i]]++;
72 }
73
74 // Convert counts to starting indices
75 int sum = 0;
76 for (int i = 0; i <= coo_rown; i++) {
77     int temp = csr_matrix->row_ptrs[i];
78     csr_matrix->row_ptrs[i] = sum;
79     sum += temp;
80 }
81
82 // Step 3: Fill CSR format data
83 for (int i = 0; i < coo_nnz; i++) {
84     int row = coo_rows[i];
85     int dest = csr_matrix->row_ptrs[row];
86
87     csr_matrix->col_indices[dest] = coo_cols[i];
88     csr_matrix->values[dest] = coo_val[i];
89
90     csr_matrix->row_ptrs[row]++;
91 }
92
93 //free temporary allocations
94 free(coo_rows);
95 free(coo_cols);
96 free(coo_val);
97 }
98
99 int print_sparse_csr(const Sparse_CSR* csr_matrix) {
100     printf("nz_id\trow\tcol\tmat_val\n");
101     printf("----\n");
102
103     for (size_t i=0; i<csr_matrix->n_row; ++i) {
104         size_t nz_start = csr_matrix->row_ptrs[i];

```

```

105         size_t nz_end = csr_matrix->row_ptrs[i+1];
106         for (size_t nz_id=nz_start; nz_id<nz_end; ++nz_id) {
107             size_t j = csr_matrix->col_indices[nz_id];
108             double val = csr_matrix->values[nz_id];
109             printf("%d\t%d\t%d\t%lf\n",nz_id, i, j, val);
110         }
111     }
112
113     return EXIT_SUCCESS;
114 }
115
116 int free_sparse_csr(Sparse_CSR* csr_matrix) {
117     free(csr_matrix->row_ptrs);
118     free(csr_matrix->col_indices);
119     free(csr_matrix->values);
120
121     return EXIT_SUCCESS;
122 }
123
124 int* vec_gen(int vec_row, int vec_col, unsigned int seed){
125     srand(seed);
126
127     int* vector = (int*)malloc(vec_row * vec_col * sizeof(int));
128
129     if(vector == NULL){
130         return NULL;
131     }
132
133     for (int i = 0; i < vec_row; i++){
134         for (int j = 0; j < vec_col; j++){
135             vector[i * vec_col + j] = rand()%10;
136         }
137     }
138
139     return vector;
140 }
141
142 double* multiplication_cuda(const Sparse_CSR* csr_matrix, const int
    * vector, int vec_col);
143
144 __global__ void matrix_vector_multiplication(const int num_rows,
    const int vec_col, const int *row_ptrs, const int *col_indices,
    const double *mat_val, const int *vec_val, double* res);
145
146
147 int main(int argc, char** argv) {
148     Sparse_CSR csr_matrix;
149     unsigned int seed = 1;
150
151     read_mtx_file( "cage4.mtx" , &csr_matrix);
152
153     //printf("Printing CSR Matrix after storing it");
154     //print_sparse_csr(&csr_matrix);
155
156     int* vec_k1 = vec_gen(csr_matrix.n_cols, 1, seed);
157     int* vec_k2 = vec_gen(csr_matrix.n_cols, 2, seed);
158     int* vec_k3 = vec_gen(csr_matrix.n_cols, 3, seed);

```

```

159     int* vec_k6 = vec_gen(csr_matrix.n_cols, 6, seed);
160
161     double* result_k1 = multiplication_cuda(&csr_matrix, vec_k1, 1);
162     double* result_k2 = multiplication_cuda(&csr_matrix, vec_k2, 2);
163     double* result_k3 = multiplication_cuda(&csr_matrix, vec_k3, 3);
164     double* result_k6 = multiplication_cuda(&csr_matrix, vec_k6, 6);
165
166     printf("result_k1 \n");
167     for(int i = 0; i < (csr_matrix.n_row * 1) ; i++){
168         printf("%lf, " , result_k1[i]);
169     }
170     printf("\n");
171
172     printf("result_k2 \n");
173     for(int i = 0; i < (csr_matrix.n_row * 2) ; i++){
174         printf("%lf, " , result_k2[i]);
175     }
176     printf("\n");
177
178     printf("result_k3 \n");
179     for(int i = 0; i < (csr_matrix.n_row * 3) ; i++){
180         printf("%lf, " , result_k3[i]);
181     }
182     printf("\n");
183
184     printf("result_k6 \n");
185     for(int i = 0; i < (csr_matrix.n_row * 6) ; i++){
186         printf("%lf, " , result_k6[i]);
187     }
188     printf("\n");
189
190     free_sparse_csr(&csr_matrix);
191     free(vec_k1);
192     free(vec_k2);
193     free(vec_k3);
194     free(vec_k6);
195     free(result_k1);
196     free(result_k2);
197     free(result_k3);
198     free(result_k6);
199
200     return EXIT_SUCCESS;
201 }
202
203 double* multiplication_cuda(const Sparse_CSR* csr_matrix, const int
    * vector, int vec_col){
204     double* result = (double*)calloc(csr_matrix->n_row * vec_col ,
        sizeof(double));
205
206     if (result == NULL) {
207         fprintf(stderr, "Failed to allocate memory for result\n")
            ;
208         exit(EXIT_FAILURE);
209     }
210
211     //variable initialization
212     int *d_row_ptr, *d_col_ind, *d_vec;

```



```

213     double *d_values, *d_result;
214
215     // Allocating kernel memory and checking for errors
216     cudaError_t err;
217
218     //allocating kernal memory
219     cudaMalloc(&d_row_ptr, (csr_matrix->n_row+1) * sizeof(int));
220     cudaMalloc(&d_col_ind, csr_matrix->n_nz * sizeof(int));
221     cudaMalloc(&d_values, csr_matrix->n_nz * sizeof(double));
222     cudaMalloc(&d_vec, (csr_matrix->n_cols*vec_col) * sizeof(int));
223     cudaMalloc(&d_result, (csr_matrix->n_row*vec_col) * sizeof(double));
224
225     //moving data from CPU to GPU
226     cudaMemcpy(d_row_ptr, csr_matrix->row_ptrs, (csr_matrix->n_row+1)
227         *sizeof(int), cudaMemcpyHostToDevice);
228     cudaMemcpy(d_col_ind, csr_matrix->col_indices, csr_matrix->n_nz*
229         sizeof(int), cudaMemcpyHostToDevice);
230     cudaMemcpy(d_values, csr_matrix->values, csr_matrix->n_nz*sizeof(
231         double), cudaMemcpyHostToDevice);
232     cudaMemcpy(d_vec, vector, (csr_matrix->n_cols*vec_col)*sizeof(int)
233         ), cudaMemcpyHostToDevice);
234     cudaMemcpy(d_result, result, (csr_matrix->n_row*vec_col)*sizeof(
235         double), cudaMemcpyHostToDevice);
236
237     //cuda event intialization
238     cudaEvent_t start,stop;
239     cudaEventCreate(&start);
240     cudaEventCreate(&stop);
241
242     cudaEventRecord(start);
243     //initializing kernel variables
244     int threadsperblock = 1024;
245     int numblock = ceil(csr_matrix->n_row/threadsperblock);
246
247     if( ceil(csr_matrix->n_row/threadsperblock) == 0) {
248         numblock = 1;
249     } else {
250         numblock = ceil(csr_matrix->n_row/threadsperblock);
251     }
252
253     printf("threadsperblock = %d, numblock = %d", threadsperblock,
254         numblock);
255
256     //launching kernels
257
258     matrix_vector_multiplication<<< numblock , threadsperblock >>>(
259         csr_matrix->n_row, vec_col, d_row_ptr, d_col_ind, d_values,
260         d_vec, d_result);
261     // Check for kernel launch errors
262     err = cudaGetLastError();
263     if (err != cudaSuccess) {
264         fprintf(stderr, "Failed to launch
265             matrix_vector_multiplication kernel: %s\n",
266             cudaGetErrorString(err));
267         exit(EXIT_FAILURE);

```

```

259     }
260
261     cudaMemcpy(result, d_result, (csr_matrix->n_row*vec_col)*sizeof(
        double), cudaMemcpyDeviceToHost);
262
263     cudaEventRecord(stop);
264     cudaEventSynchronize(stop);
265
266     cudaFree(d_row_ptr);
267     cudaFree(d_col_ind);
268     cudaFree(d_values);
269     cudaFree(d_vec);
270     cudaFree(d_result);
271
272     return result;
273 }
274
275
276 --global__ void matrix_vector_multiplication(const int num_rows,
        const int vec_col, const int *row_ptrs, const int *col_indices,
        const double *mat_val, const int *vec_val, double* res){
277     int row = blockIdx.x * blockDim.x + threadIdx.x;
278
279     //printf("row = %d , col = %d, num_rows = %d", row, col, num_rows
        );
280
281     int i, row_start, row_end;
282     double dot;
283
284     if (row < num_rows){
285         for(int k = 0; k<vec_col; k++){
286             dot = 0.0;
287             row_start = row_ptrs[row];
288             row_end = row_ptrs[row+1];
289             for(i = row_start; i<row_end ; i++){
290                 dot += mat_val[i] * vec_val[(col_indices[i]-1)*vec_col+k];
291             }
292             res[row*vec_col+k] = dot;
293         }
294     }
295 }

```

Listing 1: mat_vec_multi_cuda.c

A.1.2 Matrix Vector Multiplication CUDA Performance

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <stdbool.h>
4  #include <time.h>
5  #include <string.h>
6  #include <cuda_runtime.h>
7
8  typedef struct Sparse_CSR {
9      int n_row;
10     int n_cols;

```

```

11     int n_nz;
12     int* row_ptrs;
13     int* col_indices;
14     double* values;
15 } Sparse_CSR;
16
17 void read_mtx_file(const char *filename, Sparse_CSR *csr_matrix){
18
19     //open file
20     FILE *file = fopen(filename, "r");
21     if(!file) {
22         perror("Unable to open file!");
23         exit(EXIT_FAILURE);
24     }
25
26     //skip comments
27     char line[256];
28     while (fgets(line, sizeof(line), file)) {
29         if (line[0] != '%') break;
30     }
31
32     //initialize variable to store coo matrix from file
33     int coo_rown, coo_coln, coo_nnz;
34
35     //input no of rows, columns and non-zeros
36     sscanf(line, "%d %d %d", &coo_rown, &coo_coln, &coo_nnz);
37     printf("row = %d , col = %d , nz = %d\n", coo_rown, coo_coln,
38           coo_nnz);
39
40     //allocate memory to store the coo values
41     int* coo_rows = (int *)malloc(coo_nnz * sizeof(int));
42     int* coo_cols = (int *)malloc(coo_nnz * sizeof(int));
43     double* coo_val = (double *)malloc(coo_nnz * sizeof(double));
44
45     if (!coo_rows || !coo_cols || !coo_val){
46         perror("Memory Allocation failed!");
47         fclose(file);
48         exit(EXIT_FAILURE);
49     }
50
51     for(int i = 0; i < coo_nnz ; i++){
52         fscanf(file, "%d %d %lf" , &coo_rows[i], &coo_cols[i], &coo_val
53               [i]);
54     }
55
56     //variable assignment
57     csr_matrix->n_row = coo_rown;
58     csr_matrix->n_cols = coo_coln;
59     csr_matrix->n_nz = coo_nnz;
60
61     // Allocate CSR format arrays
62     csr_matrix->row_ptrs = (int *)malloc((coo_rown + 1) * sizeof(int)
63 );
64     memset(csr_matrix->row_ptrs, 0, (coo_rown + 1) * sizeof(int)); //
65     Initialize to 0
66

```

```

64     csr_matrix->col_indices = (int *)malloc(coo_nnz * sizeof(int));
65     csr_matrix->values = (double *)malloc(coo_nnz * sizeof(double));
66
67     // Step 1: Count non-zero elements per row
68     for (int i = 0; i < coo_nnz; i++) {
69         csr_matrix->row_ptrs[coo_rows[i]]++;
70     }
71
72     // Convert counts to starting indices
73     int sum = 0;
74     for (int i = 0; i <= coo_rown; i++) {
75         int temp = csr_matrix->row_ptrs[i];
76         csr_matrix->row_ptrs[i] = sum;
77         sum += temp;
78     }
79
80     // Step 3: Fill CSR format data
81     for (int i = 0; i < coo_nnz; i++) {
82         int row = coo_rows[i];
83         int dest = csr_matrix->row_ptrs[row];
84
85         csr_matrix->col_indices[dest] = coo_cols[i]-1;
86         csr_matrix->values[dest] = coo_val[i];
87
88         csr_matrix->row_ptrs[row]++;
89     }
90
91     //free temporary allocations
92     free(coo_rows);
93     free(coo_cols);
94     free(coo_val);
95
96 }
97
98
99 int print_sparse_csr(const Sparse_CSR* csr_matrix) {
100     printf("nz_id\trow\tcol\tmat_val\n");
101     printf("----\n");
102
103     for (size_t i=0; i<csr_matrix->n_row; ++i) {
104         size_t nz_start = csr_matrix->row_ptrs[i];
105         size_t nz_end = csr_matrix->row_ptrs[i+1];
106         for (size_t nz_id=nz_start; nz_id<nz_end; ++nz_id) {
107             size_t j = csr_matrix->col_indices[nz_id];
108             double val = csr_matrix->values[nz_id];
109             printf("%d\t%d\t%d\t%lf\n",nz_id, i, j, val);
110         }
111     }
112
113     return EXIT_SUCCESS;
114 }
115
116 int free_sparse_csr(Sparse_CSR* csr_matrix) {
117     free(csr_matrix->row_ptrs);
118     free(csr_matrix->col_indices);
119     free(csr_matrix->values);
120 }

```

```

121     return EXIT_SUCCESS;
122 }
123
124 int* vec_gen(int vec_row, int vec_col, unsigned int seed){
125     srand(seed);
126
127     int* vector = (int*)malloc(vec_row * vec_col * sizeof(int));
128
129     if(vector == NULL){
130         return NULL;
131     }
132
133     for (int i = 0; i < vec_row; i++){
134         for (int j = 0; j < vec_col; j++){
135             vector[i * vec_col + j] = rand()%10;
136         }
137     }
138
139     return vector;
140 }
141
142 double* multiplication_cuda(const Sparse_CSR* csr_matrix, const int
    * vector, int vec_col);
143
144 __global__ void matrix_vector_multiplication(const int num_rows,
    const int vec_col, const int *row_ptrs, const int *col_indices,
    const double *mat_val, const int *vec_val, double* res);
145
146
147 int main(int argc, char** argv) {
148     Sparse_CSR csr_matrix;
149     unsigned int seed = 1;
150
151     read_mtx_file( "mhd4800a.mtx" , &csr_matrix);
152
153     //printf("Printing CSR Matrix after storing it");
154     //print_sparse_csr(&csr_matrix);
155
156     int* vec_k1 = vec_gen(csr_matrix.n_cols, 1, seed);
157     int* vec_k2 = vec_gen(csr_matrix.n_cols, 2, seed);
158     int* vec_k3 = vec_gen(csr_matrix.n_cols, 3, seed);
159     int* vec_k6 = vec_gen(csr_matrix.n_cols, 6, seed);
160
161     double* result_k1 = multiplication_cuda(&csr_matrix, vec_k1, 1);
162     double* result_k2 = multiplication_cuda(&csr_matrix, vec_k2, 2);
163     double* result_k3 = multiplication_cuda(&csr_matrix, vec_k3, 3);
164     double* result_k6 = multiplication_cuda(&csr_matrix, vec_k6, 6);
165
166
167     free_sparse_csr(&csr_matrix);
168     free(vec_k1);
169     free(vec_k2);
170     free(vec_k3);
171     free(vec_k6);
172     free(result_k1);
173     free(result_k2);
174     free(result_k3);

```

```

175     free(result_k6);
176
177     return EXIT_SUCCESS;
178 }
179
180 double* multiplication_cuda(const Sparse_CSR* csr_matrix, const int
    * vector, int vec_col){
181     double* result = (double*)calloc(csr_matrix->n_row * vec_col ,
        sizeof(double));
182
183     if (result == NULL) {
184         fprintf(stderr, "Failed to allocate memory for result\n")
            ;
185         exit(EXIT_FAILURE);
186     }
187
188     //variable initialization
189     int *d_row_ptr, *d_col_ind, *d_vec;
190     double *d_values, *d_result;
191
192     // Allocating kernel memory and checking for errors
193     cudaError_t err;
194
195     //allocating kernal memory
196     cudaMalloc(&d_row_ptr, (csr_matrix->n_row+1) * sizeof(int));
197     cudaMalloc(&d_col_ind, csr_matrix->n_nz * sizeof(int));
198     cudaMalloc(&d_values, csr_matrix->n_nz * sizeof(double));
199     cudaMalloc(&d_vec, (csr_matrix->n_cols*vec_col) * sizeof(int));
200     cudaMalloc(&d_result, (csr_matrix->n_row*vec_col) * sizeof(double
        ));
201
202     //moving data from CPU to GPU
203     cudaMemcpy(d_row_ptr, csr_matrix->row_ptrs, (csr_matrix->n_row+1)
        *sizeof(int), cudaMemcpyHostToDevice);
204     cudaMemcpy(d_col_ind, csr_matrix->col_indices, csr_matrix->n_nz*
        sizeof(int), cudaMemcpyHostToDevice);
205     cudaMemcpy(d_values, csr_matrix->values, csr_matrix->n_nz*sizeof(
        double), cudaMemcpyHostToDevice);
206     cudaMemcpy(d_vec, vector, (csr_matrix->n_cols*vec_col)*sizeof(int
        ), cudaMemcpyHostToDevice);
207     cudaMemcpy(d_result, result, (csr_matrix->n_row*vec_col)*sizeof(
        double), cudaMemcpyHostToDevice);
208
209     //cuda event initialization
210     cudaEvent_t start,stop;
211     cudaEventCreate(&start);
212     cudaEventCreate(&stop);
213
214     cudaEventRecord(start);
215     //initializing kernel variables
216     int threadsperblock = 1024;
217     int numblock = ceil(csr_matrix->n_row/threadsperblock);
218
219     if( ceil(csr_matrix->n_row/threadsperblock) == 0) {
220         numblock = 1;
221     } else {
222         numblock = ceil(csr_matrix->n_row/threadsperblock);

```

```

223     }
224
225
226     //printf("threadsperblock = %d, numblock = %d", threadsperblock,
227           numblock);
228
229     //launching kernels
230
231     matrix_vector_multiplication<<< numblock , threadsperblock >>>(
232         csr_matrix->n_row, vec_col, d_row_ptr, d_col_ind, d_values,
233         d_vec, d_result);
234     // Check for kernel launch errors
235     err = cudaGetLastError();
236     if (err != cudaSuccess) {
237         fprintf(stderr, "Failed to launch
238             matrix_vector_multiplication kernel: %s\n",
239             cudaGetErrorString(err));
240         exit(EXIT_FAILURE);
241     }
242
243     cudaMemcpy(result, d_result, (csr_matrix->n_row*vec_col)*sizeof(
244         double), cudaMemcpyDeviceToHost);
245
246     cudaEventRecord(stop);
247     cudaEventSynchronize(stop);
248
249     float time_elapsed;
250     cudaEventElapsedTime(&time_elapsed, start, stop);
251     printf("Time elapsed for executing with vector columns %d = %lf\n",
252         vec_col, time_elapsed);
253
254     cudaFree(d_row_ptr);
255     cudaFree(d_col_ind);
256     cudaFree(d_values);
257     cudaFree(d_vec);
258     cudaFree(d_result);
259
260     return result;
261 }
262
263
264
265 __global__ void matrix_vector_multiplication(const int num_rows,
266     const int vec_col, const int *row_ptrs, const int *col_indices,
267     const double *mat_val, const int *vec_val, double* res){
268     int row = blockIdx.x * blockDim.x + threadIdx.x;
269
270     //printf("row = %d , col = %d, num_rows = %d", row, col, num_rows
271           );
272
273     int i, row_start, row_end;
274     double dot;
275
276     if (row < num_rows){
277         for(int k = 0; k<vec_col; k++){
278             dot = 0.0;
279             row_start = row_ptrs[row];
280             row_end = row_ptrs[row+1];

```

```

270         for(i = row_start; i<row_end ; i++){
271             dot += mat_val[i] * vec_val[(col_indices[i])*vec_col+k];
272         }
273         res[row*vec_col+k] = dot;
274     }
275 }
276 }

```

Listing 2: mat_vec_multi_cuda_performance.c

A.1.3 Matrix Vector Multiplication OpenMP

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <stdbool.h>
4  #include <time.h>
5  #include <string.h>
6  #include <omp.h>
7
8  typedef struct Sparse_CSR {
9      int n_row;
10     int n_cols;
11     int n_nz;
12     int* row_ptrs;
13     int* col_indices;
14     double* values;
15 } Sparse_CSR;
16
17 void read_mtx_file(const char *filename, Sparse_CSR *csr_matrix){
18
19     //open file
20     FILE *file = fopen(filename, "r");
21     if(!file) {
22         perror("Unable to open file!");
23         exit(EXIT_FAILURE);
24     }
25
26     //skip comments
27     char line[256];
28     while (fgets(line, sizeof(line), file)) {
29         if (line[0] != '%') break;
30     }
31
32     //initialize variable to store coo matrix from file
33     int coo_rown, coo_coln, coo_nnz;
34
35     //input no of rows, columns and non-zeros
36     sscanf(line, "%d %d %d", &coo_rown, &coo_coln, &coo_nnz);
37     printf("row = %d , col = %d , nz = %d\n", coo_rown, coo_coln,
38           coo_nnz);
39
40     //allocate memory to store the coo values
41     int* coo_rows = (int *)malloc(coo_nnz * sizeof(int));
42     int* coo_cols = (int *)malloc(coo_nnz * sizeof(int));
43     double* coo_val = (double *)malloc(coo_nnz * sizeof(double));

```



```

44 | if (!coo_rows || !coo_cols || !coo_nnz){
45 |     perror("Memory Allocation failed!");
46 |     fclose(file);
47 |     exit(EXIT_FAILURE);
48 | }
49 |
50 | for(int i = 0; i < coo_nnz ; i++){
51 |     fscanf(file, "%d %d %lf" , &coo_rows[i], &coo_cols[i], &coo_val
        [i]);
52 | }
53 |
54 | //variable assignment
55 | csr_matrix->n_row = coo_rown;
56 | csr_matrix->n_cols = coo_coln;
57 | csr_matrix->n_nz = coo_nnz;
58 |
59 |
60 | // Allocate CSR format arrays
61 | csr_matrix->row_ptrs = (int *)malloc((coo_rown + 1) * sizeof(int)
        );
62 | memset(csr_matrix->row_ptrs, 0, (coo_rown + 1) * sizeof(int)); //
        Initialize to 0
63 |
64 | csr_matrix->col_indices = (int *)malloc(coo_nnz * sizeof(int));
65 | csr_matrix->values = (double *)malloc(coo_nnz * sizeof(double));
66 |
67 | // Step 1: Count non-zero elements per row
68 | for (int i = 0; i < coo_nnz; i++) {
69 |     csr_matrix->row_ptrs[coo_rows[i]]++;
70 | }
71 |
72 | // Convert counts to starting indices
73 | int sum = 0;
74 | for (int i = 0; i <= coo_rown; i++) {
75 |     int temp = csr_matrix->row_ptrs[i];
76 |     csr_matrix->row_ptrs[i] = sum;
77 |     sum += temp;
78 | }
79 |
80 | // Step 3: Fill CSR format data
81 | for (int i = 0; i < coo_nnz; i++) {
82 |     int row = coo_rows[i];
83 |     int dest = csr_matrix->row_ptrs[row];
84 |
85 |     csr_matrix->col_indices[dest] = coo_cols[i]-1;
86 |     csr_matrix->values[dest] = coo_val[i];
87 |
88 |     csr_matrix->row_ptrs[row]++;
89 | }
90 |
91 |
92 | //free temporary allocations
93 | free(coo_rows);
94 | free(coo_cols);
95 | free(coo_val);
96 |
97 | }

```

```

98
99 int print_sparse_csr(const Sparse_CSR* csr_matrix) {
100     printf("nz_id\trow\tcol\tmat_val\n");
101     printf("----\n");
102
103     for (size_t i=0; i<csr_matrix->n_row; ++i) {
104         size_t nz_start = csr_matrix->row_ptrs[i];
105         size_t nz_end = csr_matrix->row_ptrs[i+1];
106         for (size_t nz_id=nz_start; nz_id<nz_end; ++nz_id) {
107             size_t j = csr_matrix->col_indices[nz_id];
108             double val = csr_matrix->values[nz_id];
109             printf("%d\t%d\t%d\t%lf\n",nz_id, i, j, val);
110         }
111     }
112
113     return EXIT_SUCCESS;
114 }
115
116 int free_sparse_csr(Sparse_CSR* csr_matrix) {
117     free(csr_matrix->row_ptrs);
118     free(csr_matrix->col_indices);
119     free(csr_matrix->values);
120
121     return EXIT_SUCCESS;
122 }
123
124 int* vec_gen(int vec_row, int vec_col, unsigned int seed){
125     srand(seed);
126
127     int* vector = (int*)malloc(vec_row * vec_col * sizeof(int));
128
129     if(vector == NULL){
130         return NULL;
131     }
132
133     for (int i = 0; i < vec_row; i++){
134         for (int j = 0; j < vec_col; j++){
135             vector[i * vec_col + j] = rand()%10;
136         }
137     }
138
139     return vector;
140 }
141
142 double* multiplication_openmp(const Sparse_CSR* csr_matrix, const
    int* vector, int vec_col);
143
144 int main(int argc, char** argv) {
145     Sparse_CSR csr_matrix;
146     unsigned int seed = 1;
147
148     read_mtx_file("cage4.mtx" , &csr_matrix);
149
150     //printf("Printing CSR Matrix after storing it");
151     //print_sparse_csr(&csr_matrix);
152     //printf("-----");
153

```

```

154     int* vec_k1 = vec_gen(csr_matrix.n_cols, 1, seed);
155     int* vec_k2 = vec_gen(csr_matrix.n_cols, 2, seed);
156     int* vec_k3 = vec_gen(csr_matrix.n_cols, 3, seed);
157     int* vec_k6 = vec_gen(csr_matrix.n_cols, 6, seed);
158
159     double* result_k1 = multiplication_openmp(&csr_matrix, vec_k1, 1)
160     ;
161     double* result_k2 = multiplication_openmp(&csr_matrix, vec_k2, 2)
162     ;
163     double* result_k3 = multiplication_openmp(&csr_matrix, vec_k3, 3)
164     ;
165     double* result_k6 = multiplication_openmp(&csr_matrix, vec_k6, 6)
166     ;
167
168     /*
169     printf("vec_k1 : [");
170     for(int i = 0; i < (csr_matrix.n_cols * 1) ; i++){
171         printf("%d, " , vec_k1[i]);
172     }
173     printf("]\n");
174
175     printf("vec_k2 : [");
176     for(int i = 0; i < (csr_matrix.n_cols * 2) ; i++){
177         printf("%d, " , vec_k2[i]);
178     }
179     printf("]\n");
180
181     printf("vec_k3 : [");
182     for(int i = 0; i < (csr_matrix.n_cols * 3) ; i++){
183         printf("%d, " , vec_k3[i]);
184     }
185     printf("]\n");
186
187     printf("vec_k6 : [");
188     for(int i = 0; i < (csr_matrix.n_cols * 6) ; i++){
189         printf("%d, " , vec_k6[i]);
190     }
191     printf("]\n");
192     */
193     //printf
194     ("-----\n");
195
196     printf("result_k1 : [");
197     for(int i = 0; i < (csr_matrix.n_row * 1) ; i++){
198         printf("%lf, " , result_k1[i]);
199     }
200     printf("]\n");
201
202     printf("result_k2 : [");
203     for(int i = 0; i < (csr_matrix.n_row * 2) ; i++){
204         printf("%lf, " , result_k2[i]);
205     }
206     printf("]\n");
207
208     printf("result_k3 : [");
209     for(int i = 0; i < (csr_matrix.n_row * 3) ; i++){
210         printf("%lf, " , result_k3[i]);
211     }
212     printf("]\n");
213
214     printf("result_k6 : [");
215     for(int i = 0; i < (csr_matrix.n_row * 6) ; i++){
216         printf("%lf, " , result_k6[i]);
217     }
218     printf("]\n");

```

```

205     }
206     printf("]\n");
207
208     printf("result_k6 : [\n");
209     for(int i = 0; i < (csr_matrix.n_row * 6) ; i++){
210         printf("%lf, " , result_k6[i]);
211     }
212     printf("]\n");
213
214     free_sparse_csr(&csr_matrix);
215     free(vec_k1);
216     free(vec_k2);
217     free(vec_k3);
218     free(vec_k6);
219     free(result_k1);
220     free(result_k2);
221     free(result_k3);
222     free(result_k6);
223
224     return EXIT_SUCCESS;
225 }
226
227 double* multiplication_openmp(const Sparse_CSR* csr_matrix, const
    int* vector, int vec_col){
228     double* result = (double *)malloc(csr_matrix->n_row * vec_col *
        sizeof(double));
229
230     if(result == NULL){
231         perror("Failed to allocate memory for result");
232         exit(EXIT_FAILURE);
233     }
234
235     int i, j, k, nz_start, nz_end, nz_id, vec_val;
236     double mat_val, temp;
237     //printf("k\ti\tnz_id\tj\tmat_val\tvec_val\ttemp\n");
238     #pragma omp parallel for private (i, j, k, nz_start, nz_end,
        nz_id, mat_val, vec_val, temp)
239     for(k = 0; k < vec_col; k++){
240         for(i=0 ; i < csr_matrix->n_row; i++){
241             temp = 0.0;
242             nz_start = csr_matrix->row_ptrs[i];
243             nz_end = csr_matrix->row_ptrs[i+1];
244             for(nz_id = nz_start ; nz_id < nz_end ; nz_id++){
245                 j = csr_matrix->col_indices[nz_id];
246                 mat_val = csr_matrix->values[nz_id];
247                 vec_val = vector[(j)*vec_col + k];
248                 temp += (mat_val*vec_val);
249                 //printf("%d\t%d\t%d\t%d\t%lf\t%d\t%lf\n ", k, i , nz_id, j
                    , mat_val, vec_val, temp);
250             }
251             result[i*vec_col+k] = temp;
252             //printf("result[%d] = %lf", (i*vec_col+k), result[i*vec_col+
                k]);
253         }
254     }
255
256     return result;

```

Listing 3: mat_vec_multi_openmp.c

A.1.4 Matrix Vector Multiplication OpenMP Performance

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <stdbool.h>
4  #include <time.h>
5  #include <string.h>
6  #include <omp.h>
7
8  typedef struct Sparse_CSR {
9      int n_row;
10     int n_cols;
11     int n_nz;
12     int* row_ptrs;
13     int* col_indices;
14     double* values;
15 } Sparse_CSR;
16
17 void read_mtx_file(const char *filename, Sparse_CSR *csr_matrix){
18
19     //open file
20     FILE *file = fopen(filename, "r");
21     if(!file) {
22         perror("Unable to open file!");
23         exit(EXIT_FAILURE);
24     }
25
26     //skip comments
27     char line[256];
28     while (fgets(line, sizeof(line), file)) {
29         if (line[0] != '%') break;
30     }
31
32     //initialize variable to store coo matrix from file
33     int coo_rown, coo_coln, coo_nnz;
34
35     //input no of rows, columns and non-zeros
36     sscanf(line, "%d %d %d", &coo_rown, &coo_coln, &coo_nnz);
37     printf("row = %d , col = %d , nz = %d\n", coo_rown, coo_coln,
38           coo_nnz);
39
40     //allocate memory to store the coo values
41     int* coo_rows = (int *)malloc(coo_nnz * sizeof(int));
42     int* coo_cols = (int *)malloc(coo_nnz * sizeof(int));
43     double* coo_val = (double *)malloc(coo_nnz * sizeof(double));
44
45     if (!coo_rows || !coo_cols || !coo_nnz){
46         perror("Memory Allocation failed!");
47         fclose(file);
48         exit(EXIT_FAILURE);
49     }

```

```

50     for(int i = 0; i < coo_nnz ; i++){
51         fscanf(file, "%d %d %lf" , &coo_rows[i], &coo_cols[i], &coo_val
           [i]);
52     }
53
54     //variable assignment
55     csr_matrix->n_row = coo_rown;
56     csr_matrix->n_cols = coo_coln;
57     csr_matrix->n_nz = coo_nnz;
58
59
60     // Allocate CSR format arrays
61     csr_matrix->row_ptrs = (int *)malloc((coo_rown + 1) * sizeof(int)
           );
62     memset(csr_matrix->row_ptrs, 0, (coo_rown + 1) * sizeof(int)); //
           Initialize to 0
63
64     csr_matrix->col_indices = (int *)malloc(coo_nnz * sizeof(int));
65     csr_matrix->values = (double *)malloc(coo_nnz * sizeof(double));
66
67     // Step 1: Count non-zero elements per row
68     for (int i = 0; i < coo_nnz; i++) {
69         csr_matrix->row_ptrs[coo_rows[i]]++;
70     }
71
72     // Convert counts to starting indices
73     int sum = 0;
74     for (int i = 0; i <= coo_rown; i++) {
75         int temp = csr_matrix->row_ptrs[i];
76         csr_matrix->row_ptrs[i] = sum;
77         sum += temp;
78     }
79
80     // Step 3: Fill CSR format data
81     for (int i = 0; i < coo_nnz; i++) {
82         int row = coo_rows[i];
83         int dest = csr_matrix->row_ptrs[row];
84
85         csr_matrix->col_indices[dest] = coo_cols[i]-1;
86         csr_matrix->values[dest] = coo_val[i];
87
88         csr_matrix->row_ptrs[row]++;
89     }
90
91
92     //free temporary allocations
93     free(coo_rows);
94     free(coo_cols);
95     free(coo_val);
96
97 }
98
99 int print_sparse_csr(const Sparse_CSR* csr_matrix) {
100     printf("nz_id\trow\tcol\tmat_val\n");
101     printf("----\n");
102
103     for (size_t i=0; i<csr_matrix->n_row; ++i) {

```

```

104         size_t nz_start = csr_matrix->row_ptrs[i];
105         size_t nz_end = csr_matrix->row_ptrs[i+1];
106         for (size_t nz_id=nz_start; nz_id<nz_end; ++nz_id) {
107             size_t j = csr_matrix->col_indices[nz_id];
108             double val = csr_matrix->values[nz_id];
109             printf("%d\t%d\t%d\t%f\n",nz_id, i, j, val);
110         }
111     }
112
113     return EXIT_SUCCESS;
114 }
115
116 int free_sparse_csr(Sparse_CSR* csr_matrix) {
117     free(csr_matrix->row_ptrs);
118     free(csr_matrix->col_indices);
119     free(csr_matrix->values);
120
121     return EXIT_SUCCESS;
122 }
123
124 int* vec_gen(int vec_row, int vec_col, unsigned int seed){
125     srand(seed);
126
127     int* vector = (int*)malloc(vec_row * vec_col * sizeof(int));
128
129     if(vector == NULL){
130         return NULL;
131     }
132
133     for (int i = 0; i < vec_row; i++){
134         for (int j = 0; j < vec_col; j++){
135             vector[i * vec_col + j] = rand()%10;
136         }
137     }
138
139     return vector;
140 }
141
142 double* multiplication_openmp(const Sparse_CSR* csr_matrix, const
143     int* vector, int vec_col);
144
145 int main(int argc, char** argv) {
146     Sparse_CSR csr_matrix;
147     unsigned int seed = 1;
148
149     read_mtx_file("mhd4800a.mtx" , &csr_matrix);
150
151     //printf("Printing CSR Matrix after storing it");
152     //print_sparse_csr(&csr_matrix);
153     //printf("-----");
154
155     int* vec_k1 = vec_gen(csr_matrix.n_cols, 1, seed);
156     int* vec_k2 = vec_gen(csr_matrix.n_cols, 2, seed);
157     int* vec_k3 = vec_gen(csr_matrix.n_cols, 3, seed);
158     int* vec_k6 = vec_gen(csr_matrix.n_cols, 6, seed);

```

```

159     double* result_k1 = multiplication_openmp(&csr_matrix, vec_k1, 1)
160     ;
161     double* result_k2 = multiplication_openmp(&csr_matrix, vec_k2, 2)
162     ;
163     double* result_k3 = multiplication_openmp(&csr_matrix, vec_k3, 3)
164     ;
165     double* result_k6 = multiplication_openmp(&csr_matrix, vec_k6, 6)
166     ;
167
168     free_sparse_csr(&csr_matrix);
169     free(vec_k1);
170     free(vec_k2);
171     free(vec_k3);
172     free(vec_k6);
173     free(result_k1);
174     free(result_k2);
175     free(result_k3);
176     free(result_k6);
177
178     return EXIT_SUCCESS;
179 }
180
181 double* multiplication_openmp(const Sparse_CSR* csr_matrix, const
182     int* vector, int vec_col){
183     double* result = (double *)malloc(csr_matrix->n_row * vec_col *
184         sizeof(double));
185
186     if(result == NULL){
187         perror("Failed to allocate memory for result");
188         exit(EXIT_FAILURE);
189     }
190
191     int i, j, k, nz_start, nz_end, nz_id, vec_val;
192     double mat_val, temp;
193
194     char* ompThreads = getenv("OMP_NUM_THREADS");
195     int numThreads = atoi(ompThreads);
196     double start, end, execution_time;
197     start = omp_get_wtime();
198
199     #pragma omp parallel for private (i, j, k, nz_start, nz_end,
200         nz_id, mat_val, vec_val, temp)
201     for(k = 0; k<vec_col; k++){
202         for(i=0 ; i < csr_matrix->n_row; i++){
203             temp = 0.0;
204             nz_start = csr_matrix->row_ptrs[i];
205             nz_end = csr_matrix->row_ptrs[i+1];
206             for(nz_id = nz_start ; nz_id < nz_end ; nz_id++){
207                 j = csr_matrix->col_indices[nz_id];
208                 mat_val = csr_matrix->values[nz_id];
209                 vec_val = vector[(j)*vec_col + k];
210                 temp += (mat_val*vec_val);
211             }
212             result[i*vec_col+k] = temp;
213         }
214     }
215     #pragma omp barrier

```



```
209     end = omp_get_wtime();
210     execution_time = end - start;
211     printf("Execution time for %d threads and %d columns = %lf\n", numThreads, vec_col, execution_time);
212
213     return result;
214 }
```

Listing 4: mat_vec_multi_omp_performance.c