

Module 4 Lab

CUDA Tiled Matrix Multiplication

GPU Teaching Kit – Accelerated Computing

OBJECTIVE

Implement a tiled dense matrix multiplication routine using shared memory.

PREREQUISITES

Before starting this lab, make sure that:

- You have completed “Matrix Multiplication” Lab
- You have completed the required module lectures

INSTRUCTIONS

Edit the code in the code tab to perform the following:

- allocate device memory
- copy host memory to device
- initialize thread block and kernel grid dimensions
- invoke CUDA kernel
- copy results from device to host
- deallocate device memory
- implement the matrix-matrix multiplication routine using shared memory and tiling

Instructions about where to place each part of the code is demarcated by the `//@@` comment lines.

LOCAL SETUP INSTRUCTIONS

The most recent version of source code for this lab along with the build-scripts can be found on the [Bitbucket repository](#). A description on how to

use the [CMake](#) tool in along with how to build the labs for local development found in the [README](#) document in the root of the repository.

The executable generated as a result of compiling the lab can be run using the following command:

```
./TiledMatrixMultiplication\_Template -e <expected.raw> \
-i <input0.raw>,<input1.raw> -o <output.raw> -t matrix
```

where <expected.raw> is the expected output, <input0.raw>,<input1.raw> is the input dataset, and <output.raw> is an optional path to store the results. The datasets can be generated using the dataset generator built as part of the compilation process.

QUESTIONS

- (1) How many floating operations are being performed in your matrix multiply kernel? explain.

ANSWER: **One dot-product per output matrix element. $2 * \text{numACols} * \text{numCRows} * \text{numCCols}$**

- (2) How many global memory reads are being performed by your kernel? explain.

ANSWER: **Each thread does $2 * \text{ceil}(\text{numACols} / \text{TILE_WIDTH})$ and there are $\text{numCCols} * \text{numCRows}$ active threads.**

- (3) How many global memory writes are being performed by your kernel? explain.

ANSWER: **Only the output matrix is written. $\text{numCRows} * \text{numCCols}$**

- (4) Describe what further optimizations can be implemented to your kernel to achieve a performance speedup.

ANSWER: **Adjusting kernel launch parameters for optimal occupancy,**

- (5) Compare the implementation difficulty of this kernel compared to the previous MP. What difficulties did you have with this implementation?

ANSWER: **Address index translation from the global space to the shared memory space is a new place to make errors.**

- (6) Suppose you have matrices with dimensions bigger than the max thread dimensions. Sketch an algorithm that would perform matrix multiplication algorithm that would perform the multiplication in this case.

ANSWER: **Each thread could do 4 or 9 or 16 adjacent elements of the matrix (thread coarsening).**

- (7) Suppose you have matrices that would not fit in global memory. Sketch an algorithm that would perform matrix multiplication algorithm that would perform the multiplication out of place.

ANSWER: **Just like the input is tiled for shared memory, the input could be tiled on the host side for kernel execution..**

CODE TEMPLATE

The following code is suggested as a starting point for students. The code handles the import and export as well as the checking of the solution. Students are expected to insert their code in the sections demarcated with `//@@`. Students are expected to leave the other code unchanged. The tutorial page describes the functionality of the `wb*` methods.

```

1  #include <wb.h>
2
3  #define wbCheck(stmt)                                     \
4      do {                                                 \
5          cudaError_t err = stmt;                           \
6          if (err != cudaSuccess) {                         \
7              wbLog(ERROR, "Failed to run stmt ", #stmt);   \
8              wbLog(ERROR, "Got CUDA error ... ", cudaGetErrorString(err)); \
9              return -1;                                     \
10         }                                                 \
11     } while (0)
12
13  // Compute C = A * B
14  __global__ void matrixMultiplyShared(float *A, float *B, float *C,
15                                     int numRows, int numAColumns,
16                                     int numBRows, int numBColumns,
17                                     int numCRows, int numCColumns) {
18      //@@ Insert code to implement matrix multiplication here
19      //@@ You have to use shared memory for this lab
20  }
21
22  int main(int argc, char **argv) {
23      wbArg_t args;
24      float *hostA; // The A matrix
25      float *hostB; // The B matrix
26      float *hostC; // The output C matrix
27      float *deviceA;
28      float *deviceB;
29      float *deviceC;
30      int numRows; // number of rows in the matrix A
31      int numAColumns; // number of columns in the matrix A
32      int numBRows; // number of rows in the matrix B
33      int numBColumns; // number of columns in the matrix B
34      int numCRows; // number of rows in the matrix C (you have to set this)
35      int numCColumns; // number of columns in the matrix C (you have to set
36                      // this)
37
38      args = wbArg_read(argc, argv);
39
40      wbTime_start(Generic, "Importing data and creating memory on host");
41      hostA = (float *)wbImport(wbArg_getInputFile(args, 0), &numRows,
42                               &numAColumns);
43      hostB = (float *)wbImport(wbArg_getInputFile(args, 1), &numBRows,
44                               &numBColumns);
45      //@@ Set numCRows and numCColumns

```

```

46     numCRows    = 0;
47     numCColumns = 0;
48     ///@@ Allocate the hostC matrix
49     wbTime_stop(Generic, "Importing data and creating memory on host");
50
51     wbLog(TRACE, "The dimensions of A are ", numARows, " x ", numAColumns);
52     wbLog(TRACE, "The dimensions of B are ", numBRows, " x ", numBColumns);
53
54     wbTime_start(GPU, "Allocating GPU memory.");
55     ///@@ Allocate GPU memory here
56
57     wbTime_stop(GPU, "Allocating GPU memory.");
58
59     wbTime_start(GPU, "Copying input memory to the GPU.");
60     ///@@ Copy memory to the GPU here
61
62     wbTime_stop(GPU, "Copying input memory to the GPU.");
63
64     ///@@ Initialize the grid and block dimensions here
65
66     wbTime_start(Compute, "Performing CUDA computation");
67     ///@@ Launch the GPU Kernel here
68
69     cudaDeviceSynchronize();
70     wbTime_stop(Compute, "Performing CUDA computation");
71
72     wbTime_start(Copy, "Copying output memory to the CPU");
73     ///@@ Copy the GPU memory back to the CPU here
74
75     wbTime_stop(Copy, "Copying output memory to the CPU");
76
77     wbTime_start(GPU, "Freeing GPU Memory");
78     ///@@ Free the GPU memory here
79
80     wbTime_stop(GPU, "Freeing GPU Memory");
81
82     wbSolution(args, hostC, numCRows, numCColumns);
83
84     free(hostA);
85     free(hostB);
86     free(hostC);
87
88     return 0;
89 }

```

CODE SOLUTION

The following is a possible implementation of the lab. This solution is intended for use only by the teaching staff and should not be distributed to students.

```

1  #include <wb.h>
2
3  #define wbCheck(stmt)                                     \
4      do {                                                 \
5          cudaError_t err = stmt;                           \
6          if (err != cudaSuccess) {                         \
7              wbLog(ERROR, "Failed to run stmt ", #stmt);   \
8              return -1;                                     \
9          }                                                 \
10     } while (0)
11
12 #define TILE_WIDTH 16
13
14 // Compute C = A * B
15 __global__ void matrixMultiply(float *A, float *B, float *C, int numRows,
16                                int numAColumns, int numBRows,
17                                int numBColumns, int numCRows,
18                                int numCColumns) {
19     ///@ Insert code to implement matrix multiplication here
20     __shared__ float ds_M[TILE_WIDTH][TILE_WIDTH];
21     __shared__ float ds_N[TILE_WIDTH][TILE_WIDTH];
22     int bx = blockIdx.x, by = blockIdx.y, tx = threadIdx.x, ty = threadIdx.y,
23         Row = by * TILE_WIDTH + ty, Col = bx * TILE_WIDTH + tx;
24     float Pvalue = 0;
25
26     for (int m = 0; m < (numAColumns - 1) / TILE_WIDTH + 1; ++m) {
27         if (Row < numRows && m * TILE_WIDTH + tx < numAColumns)
28             ds_M[ty][tx] = A[Row * numAColumns + m * TILE_WIDTH + tx];
29         else
30             ds_M[ty][tx] = 0;
31         if (Col < numBColumns && m * TILE_WIDTH + ty < numBRows)
32             ds_N[ty][tx] = B[(m * TILE_WIDTH + ty) * numBColumns + Col];
33         else
34             ds_N[ty][tx] = 0;
35
36         __syncthreads();
37         for (int k = 0; k < TILE_WIDTH; ++k)
38             Pvalue += ds_M[ty][k] * ds_N[k][tx];
39         __syncthreads();
40     }
41     if (Row < numCRows && Col < numCColumns)
42         C[Row * numCColumns + Col] = Pvalue;
43 }
44
45 int main(int argc, char **argv) {
46     wbArg_t args;
47     float *hostA; // The A matrix
48     float *hostB; // The B matrix
49     float *hostC; // The output C matrix
50     float *deviceA;
51     float *deviceB;
52     float *deviceC;
53     int numRows; // number of rows in the matrix A

```

```

54  int numAColumns; // number of columns in the matrix A
55  int numBRows;    // number of rows in the matrix B
56  int numBColumns; // number of columns in the matrix B
57  int numCRows;    // number of rows in the matrix C (you have to set this)
58  int numCColumns; // number of columns in the matrix C (you have to set
59                      // this)
60
61  args = wbArg_read(argc, argv);
62
63  wbTime_start(Generic, "Importing data and creating memory on host");
64  hostA = (float *)wbImport(wbArg_getInputFile(args, 0), &numARows,
65                          &numAColumns);
66  hostB = (float *)wbImport(wbArg_getInputFile(args, 1), &numBRows,
67                          &numBColumns);
68  /// Set numCRows and numCColumns
69  numCRows    = numARows;
70  numCColumns = numBColumns;
71  /// Allocate the hostC matrix
72  hostC = (float *)malloc(sizeof(float) * numCRows * numCColumns);
73  wbTime_stop(Generic, "Importing data and creating memory on host");
74
75  wbLog	TRACE, "The dimensions of A are ", numARows, " x ", numAColumns);
76  wbLog	TRACE, "The dimensions of B are ", numBRows, " x ", numBColumns);
77
78  wbTime_start(GPU, "Allocating GPU memory.");
79  /// Allocate GPU memory here
80  cudaMalloc(&deviceA, sizeof(float) * numARows * numAColumns);
81  cudaMalloc(&deviceB, sizeof(float) * numBRows * numBColumns);
82  cudaMalloc(&deviceC, sizeof(float) * numCRows * numCColumns);
83
84  wbTime_stop(GPU, "Allocating GPU memory.");
85
86  wbTime_start(GPU, "Copying input memory to the GPU.");
87  /// Copy memory to the GPU here
88  cudaMemcpy(deviceA, hostA, sizeof(float) * numARows * numAColumns,
89             cudaMemcpyHostToDevice);
90  cudaMemcpy(deviceB, hostB, sizeof(float) * numBRows * numBColumns,
91             cudaMemcpyHostToDevice);
92
93  wbTime_stop(GPU, "Copying input memory to the GPU.");
94
95  /// Initialize the grid and block dimensions here
96  dim3 dimGrid((numCColumns - 1) / TILE_WIDTH + 1,
97              (numCRows - 1) / TILE_WIDTH + 1, 1);
98  dim3 dimBlock(TILE_WIDTH, TILE_WIDTH, 1);
99
100  wbTime_start(Compute, "Performing CUDA computation");
101  /// Launch the GPU Kernel here
102  matrixMultiply<<<dimGrid, dimBlock>>>(
103      deviceA, deviceB, deviceC, numARows, numAColumns, numBRows,
104      numBColumns, numCRows, numCColumns);
105
106  cudaDeviceSynchronize();

```

```

107     wbTime_stop(Compute, "Performing CUDA computation");
108
109     wbTime_start(Copy, "Copying output memory to the CPU");
110     /// Copy the GPU memory back to the CPU here
111     cudaMemcpy(hostC, deviceC, sizeof(float) * numCRows * numCColumns,
112               cudaMemcpyDeviceToHost);
113
114     wbTime_stop(Copy, "Copying output memory to the CPU");
115
116     wbTime_start(GPU, "Freeing GPU Memory");
117     /// Free the GPU memory here
118     cudaFree(deviceA);
119     cudaFree(deviceB);
120     cudaFree(deviceC);
121
122     wbTime_stop(GPU, "Freeing GPU Memory");
123
124     wbSolution(args, hostC, numCRows, numCColumns);
125
126     free(hostA);
127     free(hostB);
128     free(hostC);
129
130     return 0;
131 }

```