

**ME759**  
**High Performance Computing for Engineering Applications**  
**Tiled Matrix Multiplication**

Edit the source files **matrixmul.cu** and **matrixmul\_kernel.cu** to complete the functionality of the matrix multiplication  $P=M*N$  on the device. The two input matrices have dimensions that allow them to be multiplied. Moreover, no dimension of the two input matrices is larger than 4096. There are several modes of operation for the application. Note that the file interface has been updated to allow the size of the input matrices to be read in.

- a) No arguments: The application will create two randomly sized and initialized matrices such that the matrix operation  $M*N$  is valid, and  $P$  is properly sized to hold the result. After the device multiplication is invoked, it will compute the correct solution matrix using the CPU, and compare that solution with the device-computed solution. If it matches (within a certain tolerance), it will print out "**Test PASSED**" to the screen before exiting.
- b) One argument: The application will use the random initialization to create the input matrices, and write the device-computed output to the file specified by the argument.
- c) Three arguments: The application will read input matrices from provided files. The first argument should be a file containing three integers. The first, second and third integers will be used as **M.height**, **M.width**, and **N.width**. The second and third function arguments will be expected to be files which have exactly enough entries to fill matrices **M** and **N** respectively. No output is written to file.
- d) Four arguments: The application will read its inputs from the files provided by the first three arguments as described above, and write its output to the file provided in the fourth.

Note that if you wish to use the output of one run of the application as an input, you must delete the first line in the output file, which displays the accuracy of the values within the file. The value is not relevant for this application.

Provide a file called **comments.pdf** in which you report on the following:

- 1) For block sizes of 16x16, Set **CUDA\_PROFILE=1** and provide both on the forum (under posting HW5-Tiled Matrix Multiplication Profiling/Timing) and in **comments.pdf** the content of your **cuda\_profile\_0.log** file (or whatever the file is where the profiling information was dumped after one successful run). In **comments.pdf** provide short comments on the results you see reported in the **cuda\_profile\_0.log** file.

**Answer:**

PROFILING DATA FOR BLOCK\_SIZE 16

Inclusive time = 71.395744 ms : Exclusive time = 53.398434 ms GPU computation complete

Start CPU computation

CPU computation complete

Test PASSED

Dimension M[height,width]: 2466 1653

Dimension N[height,width]: 1653 1520

CPU computation time = 20800.273438 ms [pk Gupta3@euler02 Tiled\_matrixmul]\$ cat cuda\_profile\_0.log

# CUDA\_PROFILE\_LOG\_VERSION 2.0

# CUDA\_DEVICE 0 GeForce GTX 480

# CUDA\_CONTEXT 1

# TIMESTAMPFACTOR 132728437adfb7c8

method,gputime,cputime,occupancy

method=[ memcopyHtoD ] gputime=[ 3804.736 ] cputime=[ 3937.817 ]

method=[ memcopyHtoD ] gputime=[ 4095.488 ] cputime=[ 4184.133 ]

method=[ memcopyHtoD ] gputime=[ 2480.640 ] cputime=[ 2560.608 ]

method=[ \_Z15MatrixMulKernel6MatrixS\_S\_ ] gputime=[ 53368.574 ] cputime=[ 6.510 ] occupancy=[ 0.833 ]

method=[ memcopyDtoH ] gputime=[ 6566.496 ] cputime=[ 7215.665 ]

- 2) Just like before, but use block sizes of 32x32.

#### PROFILING DATA FOR BLOCK\_SIZE 32

Inclusive time = 71.401794 ms : Exclusive time = 53.408321 ms GPU computation complete

Start CPU computation

CPU computation complete

Test PASSED

Dimension M[height,width]: 2466 1653

Dimension N[height,width]: 1653 1520

CPU computation time = 29530.380859 ms [pkgupta3@euler02 Tiled\_matrixmul]\$

# CUDA\_PROFILE\_LOG\_VERSION 2.0

# CUDA\_DEVICE 0 GeForce GTX 480

# CUDA\_CONTEXT 1

# TIMESTAMPFACTOR 13272843a7cf1b79

method,gputime,cputime,occupancy

method=[ memcpyHtoD ] gputime=[ 4238.752 ] cputime=[ 4376.221 ]

method=[ memcpyHtoD ] gputime=[ 2518.720 ] cputime=[ 2595.958 ]

method=[ memcpyHtoD ] gputime=[ 3760.512 ] cputime=[ 3817.515 ]

method=[ \_Z15MatrixMulKernel6MatrixS\_S\_ ] gputime=[ 54015.938 ] cputime=[ 7.537 ] occupancy=[ 0.667 ]

method=[ memcpyDtoH ] gputime=[ 6581.600 ] cputime=[ 61381.324 ]

3) Explain why the timing differences between the two cases and why one is superior to the other.

**Answer:** In case of 16 x 16 matrices, the occupancy is 0.833. So, the no. of blocks it is using is 5 out of the 6 and in case of 32 x 32 matrices, the occupancy is 0.667. So, the no. of blocks it is using is 4 out of 6. Hence, 16 x 16 case is superior.

4) Run your application to multiply two matrices of dimension 4096x4096. Report the time required to perform the matrix multiplication on the CPU and then on the GPU. When timing the GPU version, make sure you report both the inclusive and exclusive times.

Provide your timing results both in **comments.pdf** and on the forum.

**Answer:**

Inclusive time = 667.773926 ms : Exclusive time = 582.770691 ms GPU computation complete

Start CPU computation

CPU computation complete

Test PASSED

Dimension M[height,width]: 4096 4096

Dimension N[height,width]: 4096 4096

CPU computation time = 2313680.750000 ms

5) If your input matrices are square, what is the largest dimension that your code can handle on GTX480 after you further modify your code? Provide your timing and profiling results both in **comments.pdf** and on the forum.

**Answer:**

Max dimension is : 11344 x 11344

Timing Results are as follows

Inclusive time = 13328.544922 ms : Exclusive time = 12800.438477 ms GPU computation complete

Start CPU computation

CPU computation complete

Test PASSED

Dimension M[height,width]: 11344 11344

Dimension N[height,width]: 11344 11344

Profiling result:

# CUDA\_PROFILE\_LOG\_VERSION 2.0

# CUDA\_DEVICE 0 GeForce GTX 480

# CUDA\_CONTEXT 1

# TIMESTAMPFACTOR 132c63aafa1f15d8

method,gputime,cputime,occupancy

method=[ memcpyHtoD ] gputime=[ 114319.844 ] cputime=[ 114405.711 ]

method=[ memcpyHtoD ] gputime=[ 114215.523 ] cputime=[ 114247.922 ]

method=[ memcpyHtoD ] gputime=[ 114146.750 ] cputime=[ 114162.203 ]

method=[ \_Z15MatrixMulKernel6MatrixS\_S\_ ] gputime=[ 12799780.000 ] cputime=[ 6.761 ] occupancy=[ 0.833 ]

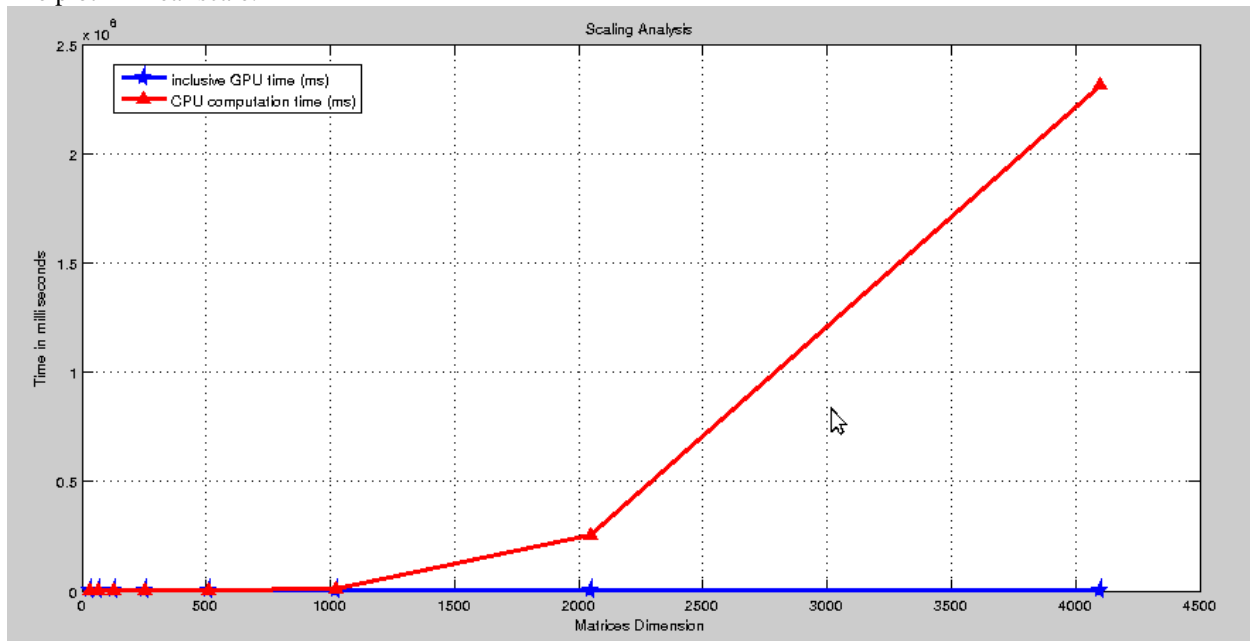
method=[ memcpyDtoH ] gputime=[ 214889.469 ] cputime=[ 215800.094 ]

6) Include in **comments.pdf** and also upload on the forum a plot (**png** format recommended) that shows how both the CPU and *inclusive* GPU times scale with the dimension of the problem. The scaling analysis should be based on matrices of dimension 32x32, 64x64, 128x128,...,4096x4096. When you provide this plot make sure you build the executables in “production” mode (as opposed to “debug” mode).

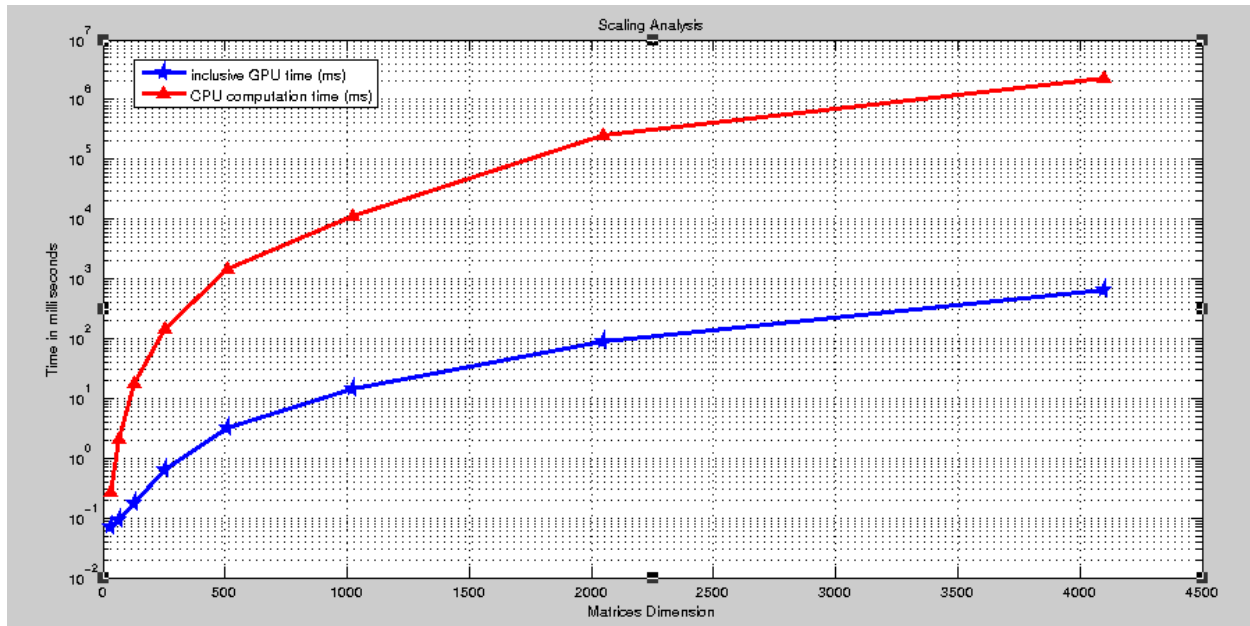
**Answer:**

On the next page

The plot in linear scale:



Plot on log scale:



7) You might notice that if you have very large matrices, your accuracy test (that verifies the result on the GPU against the result on the CPU) might fail. Why do you believe to be the cause of this odd behavior? That is, working ok for smaller matrices but not when dealing with large matrices.

**Answer:**

**The CPU computation time does not gives accurate results for very high matrices some of the timing results are as follows:**

**Note: the CPU results which are in accurate are marked in 'red' and those accurate are marked in 'green'.**

**For 11344 x 11344:**

```
Inclusive time = 13328.544922 ms : Exclusive time = 12800.438477 msGPU computation complete
Start CPU computation
CPU computation complete
Test PASSED
Dimension M[height,width]: 11344 11344
Dimension N[height,width]: 11344 11344
CPU computation time = 0.003200 ms
```

**For 11310 x 11310:**

```
Inclusive time = 13022.593750 ms : Exclusive time = 12472.420898 msGPU computation complete
Start CPU computation
CPU computation complete
Test PASSED
```

Dimension M[height,width]: 11310 11310  
Dimension N[height,width]: 11310 11310  
CPU computation time = 0.014464 ms

For 11300 x 11300:

Inclusive time = 13057.901367 ms : Exclusive time = 12474.960938 msGPU computation complete  
Start CPU computation  
CPU computation complete  
Test PASSED  
Dimension M[height,width]: 11300 11300  
Dimension N[height,width]: 11300 11300  
CPU computation time = 0.014592 ms

For 2048 x 2048:

Inclusive time = 98.589661 ms : Exclusive time = 72.942688 msGPU computation complete  
Start CPU computation  
CPU computation complete  
Test PASSED  
Dimension M[height,width]: 2048 2048  
Dimension N[height,width]: 2048 2048

CPU completion time = 232779.640625 ms

For 4096 x 4096

Inclusive time = 678.817444 ms : Exclusive time = 582.817078 msGPU computation complete  
Start CPU computation  
CPU computation complete  
Test PASSED  
Dimension M[height,width]: 4096 4096  
Dimension N[height,width]: 4096 4096

CPU completion time = 1993869.500000 ms

**The CPU is out of memory for such large calculations and the result is overflow and hence inaccurate.**

### Grading:

Your submission will be graded using the following scheme.

- a) Demo/knowledge: 60%
  - Produces correct result output files for provided inputs
- b) Functionality: 30%
  - Shared memory is used in the kernel to mask global memory access latencies.
- c) Report: 10%
  - Good comments provided in relation to questions 3) and 7) above.