

 (/mp/5983?show=code) Tiled Matrix Multiplication Attempt

Attempt Summary

Submit Attempt for Grading

Remember to answer the questions before clicking.

Dataset Id:	5
Created:	<u>less than a minute ago (2022-03-05 09:06:35 +0000 UTC)</u>
Status:	Correct solution for this dataset.

Timer Output

Kind	Location	Time (ms)	Message
Generic	main.cu::91	13.077868	Importing data and creating memory on host
GPU	main.cu::108	0.840356	Allocating GPU memory.
GPU	main.cu::120	0.119089	Copying input memory to the GPU.
Compute	main.cu::133	0.254237	Performing CUDA computation
Copy	main.cu::141	0.190876	Copying output memory to the CPU
GPU	main.cu::148	0.102842	Freeing GPU Memory

Logger Output

Level	Location	Message
Trace	main::104	The dimensions of A are 200 x 100
Trace	main::105	The dimensions of B are 100 x 256
Trace	main::106	The dimensions of C are 200 x 256

Program Code

```

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 #define TILE_WIDTH 32
14 #define BLOCK_SIZE 8
15 int ceil(int a, int b){
16     return (a + b - 1)/b;
17 }
18
19
20 // Compute C = A * B
21 __global__ void matrixMultiplyShared(float *A, float *B, float *C,
22                                     int numARows, int numAColumns,
23                                     int numBRows, int numBColumns,
24                                     int numCRows, int numCColumns) {
25     //@@ Insert code to implement matrix multiplication here
26     //@@ You have to use shared memory for this MP
27     __shared__ float TileP[TILE_WIDTH][TILE_WIDTH];
28     __shared__ float TileQ[TILE_WIDTH][TILE_WIDTH];
29
30     int block_x = blockIdx.x;
31     int block_y = blockIdx.y;
32     int thread_x = threadIdx.x;

```

```
33  int thread_y = threadIdx.y;
34
35  int row = (block_y * blockDim.y) + thread_y;
36  int col = (block_x * blockDim.x) + thread_x;
37  int next = (numAColumns + BLOCK_SIZE - 1) / BLOCK_SIZE;
38  float hold = 0;
39
40  for (int a = 0; a < next; a++){
41      //load the first matrix tile
42      if((a * BLOCK_SIZE + thread_x) > numAColumns){
43          TileP[thread_y][thread_x] = 0.0;
44      }
45      else{
46          TileP[thread_y][thread_x] = A[row * numAColumns + a * BLOCK_SIZE +
47      }
48      //load the second matrix tile
49      if((a * BLOCK_SIZE + thread_y) >= numBRows){
50          TileQ[thread_y][thread_x] = 0.0;
51      }
52      else{
53          TileQ[thread_y][thread_x] = B[(a * BLOCK_SIZE + thread_y) * numBCol
54      }
55
56      __syncthreads();
57      //perform multiplication calculation
58      for (int b = 0; b < BLOCK_SIZE; b++){
59          hold += TileP[thread_y][b] * TileQ[b][thread_x];
60      }
61      __syncthreads();
62  }
63
64  if (row < numCRows && col < numCColumns){
65      C[row * numCColumns + col] = hold;
66  }
67
68  //__syncthreads();
69
70
71 }
72
73 int main(int argc, char **argv) {
74     wbArg_t args;
75     float *hostA; // The A matrix
76     float *hostB; // The B matrix
77     float *hostC; // The output C matrix
```

```
78 float *deviceA;
79 float *deviceB;
80 float *deviceC;
81 int numRows;    // number of rows in the matrix A
82 int numColumns; // number of columns in the matrix A
83 int numBRows;   // number of rows in the matrix B
84 int numBColumns; // number of columns in the matrix B
85 int numCRows;   // number of rows in the matrix C (you have to set thi
86 int numCColumns; // number of columns in the matrix C (you have to set
87                 // this)
88
89 args = wbArg_read(argc, argv);
90
91 wbTime_start(Generic, "Importing data and creating memory on host");
92 hostA = (float *)wbImport(wbArg_getInputFile(args, 0), &numRows,
93                           &numColumns);
94 hostB = (float *)wbImport(wbArg_getInputFile(args, 1), &numBRows,
95                           &numBColumns);
96 //@@ Set numCRows and numCColumns
97 numCRows = numRows;
98 numCColumns = numBColumns;
99 //@@ Allocate the hostC matrix
100 wbTime_stop(Generic, "Importing data and creating memory on host");
101
102 hostC = (float *) malloc((numCRows * numCColumns) * sizeof(float));
103
104 wbLog	TRACE, "The dimensions of A are ", numRows, " x ", numColumns);
105 wbLog	TRACE, "The dimensions of B are ", numBRows, " x ", numBColumns);
106 wbLog	TRACE, "The dimensions of C are ", numCRows, " x ", numCColumns);
107
108 wbTime_start(GPU, "Allocating GPU memory.");
109 //@@ Allocate GPU memory here
110
111 int size_of_A = numRows * numColumns * sizeof(float);
112 int size_of_B = numBRows * numBColumns * sizeof(float);
113 int size_of_C = numCRows * numCColumns * sizeof(float);
114 cudaMalloc((void **) &deviceA, size_of_A);
115 cudaMalloc((void **) &deviceB, size_of_B);
116 cudaMalloc((void **) &deviceC, size_of_C);
117
118 wbTime_stop(GPU, "Allocating GPU memory.");
119
120 wbTime_start(GPU, "Copying input memory to the GPU.");
121 //@@ Copy memory to the GPU here
122
```

```
123  cudaMemcpy(deviceA, hostA, size_of_A, cudaMemcpyHostToDevice);
124  cudaMemcpy(deviceB, hostB, size_of_B, cudaMemcpyHostToDevice);
125
126  wbTime_stop(GPU, "Copying input memory to the GPU.");
127
128  //@@ Initialize the grid and block dimensions here
129
130  dim3 dimensionBlock(BLOCK_SIZE, BLOCK_SIZE, 1);
131  dim3 dimensionGrid(ceil(numCColumns, BLOCK_SIZE), ceil(numCRows, BLOCK_
132
133  wbTime_start(Compute, "Performing CUDA computation");
134  //@@ Launch the GPU Kernel here
135
136  matrixMultiplyShared<<<dimensionGrid, dimensionBlock>>>(deviceA, device
137
138  cudaDeviceSynchronize();
139  wbTime_stop(Compute, "Performing CUDA computation");
140
141  wbTime_start(Copy, "Copying output memory to the CPU");
142  //@@ Copy the GPU memory back to the CPU here
143
144  cudaMemcpy(hostC, deviceC, size_of_C, cudaMemcpyDeviceToHost);
145
146  wbTime_stop(Copy, "Copying output memory to the CPU");
147
148  wbTime_start(GPU, "Freeing GPU Memory");
149  //@@ Free the GPU memory here
150
151  cudaFree(deviceA);
152  cudaFree(deviceB);
153  cudaFree(deviceC);
154
155  wbTime_stop(GPU, "Freeing GPU Memory");
156
157  wbSolution(args, hostC, numCRows, numCColumns);
158
159  free(hostA);
160  free(hostB);
161  free(hostC);
162
163  return 0;
164 }
165
```

