

Lab 2

George Onwubuya

September 25, 2018

1 Output Files

1.1 Output File 1

```
Setting up the problem...0.020886 s
  A: 1000 x 1000
  B: 1000 x 1000
  C: 1000 x 1000
Allocating device variables...0.176401 s
Copying data from host to device...0.004717 s
Launching kernel...0.022178 s
Copying data from device to host...0.004555 s
Verifying results...TEST PASSED
```

```
Setting up the problem...0.081405 s
  A: 2000 x 2000
  B: 2000 x 2000
  C: 2000 x 2000
Allocating device variables...0.157044 s
Copying data from host to device...0.017920 s
Launching kernel...0.163624 s
Copying data from device to host...0.015868 s
Verifying results...TEST PASSED
```

```
Setting up the problem...0.318687 s
  A: 4000 x 4000
  B: 4000 x 4000
  C: 4000 x 4000
Allocating device variables...0.155963 s
Copying data from host to device...0.071036 s
Launching kernel...1.308468 s
Copying data from device to host...0.049722 s
```

Verifying results...TEST PASSED

Setting up the problem...1.281662 s

A: 8000 x 8000

B: 8000 x 8000

C: 8000 x 8000

Allocating device variables...0.162413 s

Copying data from host to device...0.288608 s

Launching kernel...10.493221 s

Copying data from device to host...0.186492 s

Verifying results...TEST PASSED

Setting up the problem...0.121365 s

A: 4000 x 2000

B: 2000 x 2000

C: 4000 x 2000

Allocating device variables...0.157077 s

Copying data from host to device...0.025333 s

Launching kernel...0.327244 s

Copying data from device to host...0.022017 s

Verifying results...TEST PASSED

1.2 Output File 2

Setting up the problem...0.008318 s

A: 1000 x 500

B: 500 x 500

C: 1000 x 500

Allocating device variables...0.179471 s

Copying data from host to device...0.001774 s

Launching kernel...0.005649 s

Copying data from device to host...0.002503 s

Verifying results...TEST PASSED

Setting up the problem...0.030971 s

A: 2000 x 1000

B: 1000 x 1000

C: 2000 x 1000

Allocating device variables...0.153657 s

Copying data from host to device...0.005067 s

Launching kernel...0.043961 s

```

Copying data from device to host...0.005322 s
Verifying results...TEST PASSED

```

```

Setting up the problem...0.122200 s
  A: 4000 x 2000
  B: 2000 x 2000
  C: 4000 x 2000
Allocating device variables...0.157203 s
Copying data from host to device...0.027051 s
Launching kernel...0.326948 s
Copying data from device to host...0.018884 s
Verifying results...TEST PASSED

```

```

Setting up the problem...0.486641 s
  A: 8000 x 4000
  B: 4000 x 4000
  C: 8000 x 4000
Allocating device variables...0.156578 s
Copying data from host to device...0.098028 s
Launching kernel...2.614291 s
Copying data from device to host...0.102813 s
Verifying results...TEST PASSED

```

2 Performance Analysis

2.1 Square Matrices (n x n)

Execution Time (seconds) for Each Process					
Elements(nxn)	Setting Up	DeviceVar	Kernel	HostToDevice	DeviceToHost
1000	0.020886	0.176401	0.022178	0.004717	0.004555
2000	0.81405	0.157044	0.163624	0.017920	0.015868
4000	0.318687	0.155963	1.308468	0.071036	0.049722
8000	1.281662	0.162413	10.493221	0.288608	0.186492

2.2 Rectangle Matrices ($A = m \times k$ and $B = n \times k$)

Execution Time (seconds) for Each Process					
Elements (m,k,n)	Setting Up	DeviceVar	Kernel	HostToDevice	DeviceToHost
1000, 500, 500	0.008318	0.179471	0.005649	0.001774	0.002503
2000, 1000, 1000	0.030971	0.153657	0.043961	0.005067	0.005322
4000, 2000, 2000	0.122200	0.157203	0.326948	.027051	0.018884
8000, 4000, 4000	0.486641	0.156578	2.614291	0.098028	0.102813

2.3 Comments

For the square matrices each process time followed a similar pattern except for the process of allocating of 'device variables' which shows no correlation or change as the number of elements increase. As the number of elements increase the time taken setting up the problem, allocating device variables, launching the kernel and copying data from the host to the device and vice versa also increase. There is a noticeable direct proportional relationship between these processes and the number of elements. The same conclusion can be made for the rectangular matrices. It can also be observed that the time taken to allocate device variables are approximately the same regardless of the number of elements.

3 Answers

3.1 (i)

The elements in matrix A are loaded m times and in B are loaded n times.

3.2 (ii)

If storing to the global memory is ignored then for each element the global memory is accessed twice. On the first load a floating point operation used to multiply the elements from the corresponding row and column and on the second load a floating-point operation is used to perform addition. There are therefore 2 global memory accesses and 2 floating point operations and therefore the memory access to floating-point compute ratio is 1:1.

4 Main

```

1  /*****
2  *cr
3  *cr          (C) Copyright 2010 The Board of Trustees of the
4  *cr          University of Illinois
5  *cr          All Rights Reserved

```

```

6  *cr
7
8  ↪ *****/
9  #include <stdio.h>
10 #include <stdlib.h>
11 #include "kernel.cu"
12 #include "support.h"
13
14 int main (int argc, char *argv[])
15 {
16
17     Timer timer;
18     cudaError_t cuda_ret;
19
20     // Initialize host variables
21     ↪ -----
22
23     printf("\nSetting up the problem..."); fflush(stdout);
24     startTime(&timer);
25
26     float *A_h, *B_h, *C_h;
27     float *A_d, *B_d, *C_d;
28     size_t A_sz, B_sz, C_sz;
29     unsigned matArow, matAcol;
30     unsigned matBrow, matBcol;
31     dim3 dim_grid, dim_block;
32
33     if (argc == 1) {
34         matArow = 1000;
35         matAcol = matBrow = 1000;
36         matBcol = 1000;
37     } else if (argc == 2) {
38         matArow = atoi(argv[1]);
39         matAcol = matBrow = atoi(argv[1]);
40         matBcol = atoi(argv[1]);
41     } else if (argc == 4) {
42         matArow = atoi(argv[1]);
43         matAcol = matBrow = atoi(argv[2]);
44         matBcol = atoi(argv[3]);
45     } else {
46         printf("\n    Invalid input parameters!"
47             "\n    Usage: ./sgemm                # All matrices
48             ↪ are 1000 x 1000"
49             "\n    Usage: ./sgemm <m>          # All matrices
50             ↪ are m x m"

```

```

48         "\n    Usage: ./sgemm <m> <k> <n>    # A: m x k, B: k
        ↪ x n, C: m x n"
49         "\n");
50         exit(0);
51     }
52
53     A_sz = matArow*matAcol;
54     B_sz = matBrow*matBcol;
55     C_sz = matArow*matBcol;
56
57     A_h = (float*) malloc( sizeof(float)*A_sz );
58     for (unsigned int i=0; i < A_sz; i++) { A_h[i] =
        ↪ (rand()%100)/100.00; }
59
60     B_h = (float*) malloc( sizeof(float)*B_sz );
61     for (unsigned int i=0; i < B_sz; i++) { B_h[i] =
        ↪ (rand()%100)/100.00; }
62
63     C_h = (float*) malloc( sizeof(float)*C_sz );
64
65     stopTime(&timer); printf("%f s\n", elapsedTime(timer));
66     printf("    A: %u x %u\n    B: %u x %u\n    C: %u x %u\n",
        ↪ matArow, matAcol,
67         matBrow, matBcol, matArow, matBcol);
68
69     // Allocate device variables
        ↪ -----
70
71     printf("Allocating device variables..."); fflush(stdout);
72     startTime(&timer);
73
74     //INSERT CODE HERE
75
76     cuda_ret = cudaMalloc((void**)&A_d, sizeof(float)*A_sz );
77     if(cuda_ret != cudaSuccess) FATAL("Unable to allocate
        ↪ device memory");
78
79     cuda_ret = cudaMalloc((void**)&B_d, sizeof(float)*B_sz );
80     if(cuda_ret != cudaSuccess) FATAL("Unable to allocate
        ↪ device memory");
81
82     cuda_ret = cudaMalloc((void**)&C_d, sizeof(float)*C_sz);
83     if(cuda_ret != cudaSuccess) FATAL("Unable to allocate
        ↪ device memory");
84
85     cudaDeviceSynchronize();

```

```

86     stopTime(&timer); printf("%f s\n", elapsedTime(timer));
87
88     // Copy host variables to device
89     ↪ -----
90
91     printf("Copying data from host to device...");
92     ↪ fflush(stdout);
93     startTime(&timer);
94
95     //INSERT CODE HERE
96
97     cuda_ret = cudaMemcpy(A_d, A_h, sizeof(float)*A_sz,
98     ↪ cudaMemcpyHostToDevice);
99     if(cuda_ret != cudaSuccess) FATAL("Unable to copy memory
100     ↪ to device");
101
102     cuda_ret = cudaMemcpy(B_d, B_h, sizeof(float)*B_sz,
103     ↪ cudaMemcpyHostToDevice);
104     if(cuda_ret != cudaSuccess) FATAL("Unable to copy memory
105     ↪ to device");
106
107     cudaDeviceSynchronize();
108     stopTime(&timer); printf("%f s\n", elapsedTime(timer));
109
110     // Launch kernel using standard sgemm interface
111     ↪ -----
112
113     printf("Launching kernel..."); fflush(stdout);
114     startTime(&timer);
115     basicSgemm('N', 'N', matArow, matBcol, matBrow, 1.0f, \
116     ↪ A_d, matArow, B_d, matBrow, 0.0f, C_d, matBrow);
117
118     cuda_ret = cudaDeviceSynchronize();
119     if(cuda_ret != cudaSuccess) FATAL("Unable to launch
120     ↪ kernel");
121     stopTime(&timer); printf("%f s\n", elapsedTime(timer));
122
123     // Copy device variables from host
124     ↪ -----
125
126     printf("Copying data from device to host...");
127     ↪ fflush(stdout);
128     startTime(&timer);
129
130     //INSERT CODE HERE
131
132     cuda_ret = cudaMemcpy(C_h, C_d, sizeof(float)*C_sz,
133     ↪ cudaMemcpyDeviceToHost);

```

```

121         if(cuda_ret != cudaSuccess) FATAL("Unable to copy from
           ↪ device");
122
123     cudaDeviceSynchronize();
124     stopTime(&timer); printf("%f s\n", elapsedTime(timer));
125
126     // Verify correctness
           ↪ -----
127
128     printf("Verifying results..."); fflush(stdout);
129
130     verify(A_h, B_h, C_h, matArow, matAcol, matBcol);
131
132
133     // Free memory
           ↪ -----
134
135     free(A_h);
136     free(B_h);
137     free(C_h);
138
139     //INSERT CODE HERE
140
141     cudaFree(A_d);
142     cudaFree(B_d);
143     cudaFree(C_d);
144
145
146     return 0;
147
148 }

```

5 Kernel

```

1  /*****:{*****/
2  *cr
3  *cr          (C) Copyright 2010 The Board of Trustees of the
4  *cr                      University of Illinois
5  *cr                      All Rights Reserved
6  *cr
7
           ↪ *****/
8
9  #include <stdio.h>
10

```



```

11 __global__ void mysgemm(int m, int n, int k, const float *A,
    ↪ const float *B, float* C) {
12
13
14     ↪ /******
15     *
16     * Compute C = A x B
17     * where A is a (m x k) matrix
18     * where B is a (k x n) matrix
19     * where C is a (m x n) matrix
20     *
21     ↪ *****/
22
23     // INSERT KERNEL CODE HERE
24     int row, col;
25
26     row = blockIdx.y*blockDim.y+threadIdx.y;
27
28     col = blockIdx.x*blockDim.x+threadIdx.x;
29
30     if(( row < m) && (col < n))
31     {
32         float acc = 0;
33
34         for(int index = 0; index < k; index++)
35         {
36             acc = acc + A[row * k + index] * B[index * n + col];
37         }
38
39         C[row * n + col] = acc;
40
41     }
42 }
43
44 void basicSgemv(char transa, char transb, int m, int n, int k,
    ↪ float alpha, const float *A, int lda, const float *B, int
    ↪ ldb, float beta, float *C, int ldc)
45 {
46     if ((transa != 'N') && (transa != 'n')) {
47         printf("unsupported value of 'transa'\n");
48         return;
49     }
50
51     if ((transb != 'N') && (transb != 'n')) {

```

```

52     printf("unsupported value of 'transb'\n");
53     return;
54 }
55
56 if ((alpha - 1.0f > 1e-10) || (alpha - 1.0f < -1e-10)) {
57     printf("unsupported value of alpha\n");
58     return;
59 }
60
61 if ((beta - 0.0f > 1e-10) || (beta - 0.0f < -1e-10)) {
62     printf("unsupported value of beta\n");
63     return;
64 }
65
66 // Initialize thread block and kernel grid dimensions
67 ↪ -----
68
69 const unsigned int BLOCK_SIZE = 16; // Use 16x16 thread
70 ↪ blocks
71
72 //INSERT CODE HERE
73
74 dim3 block(BLOCK_SIZE, BLOCK_SIZE, 1);
75 dim3 grid((n + BLOCK_SIZE - 1)/BLOCK_SIZE, (m + BLOCK_SIZE
76 ↪ -1)/BLOCK_SIZE, 1);
77
78 // Invoke CUDA kernel
79 ↪ -----
80
81 //INSERT CODE HERE
82
83 mysgemm<<< grid, block>>>(m, n, k, A, B, C);
84 }

```