

Computational Science on Many-Core Architectures:

Exercise 5:

Viktor Beck, 11713110

November 2022

Task 1: Performance Modeling: Parameter Identification

To ensure the correctness of the estimated values of this task at least 6 repetitions were made for each measurement and then the median of those was taken as the time (for that particular N).

a)

The PCI Express latency for `cudaMemcpy()` was estimated by simply measuring the time in a program for vector addition for different vector sizes N.

```
1      // repetitions
2      for (int j = 0; j < reps; j++){
3
4          // MEASURING TIME FROM HERE
5          timer.reset();
6
7          cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
8
9          elapsed = timer.get();
10         // TO HERE
11         time_vec.push_back(elapsed);
12
13     }
```

Since the latency is independent from N an estimate can be made from looking at the constant part of the measurements (for small N) which can be seen in the Figure. The estimates for the two GPUs are:

- RTX: $t = 4\mu s$
- K40: $t = 7\mu s$

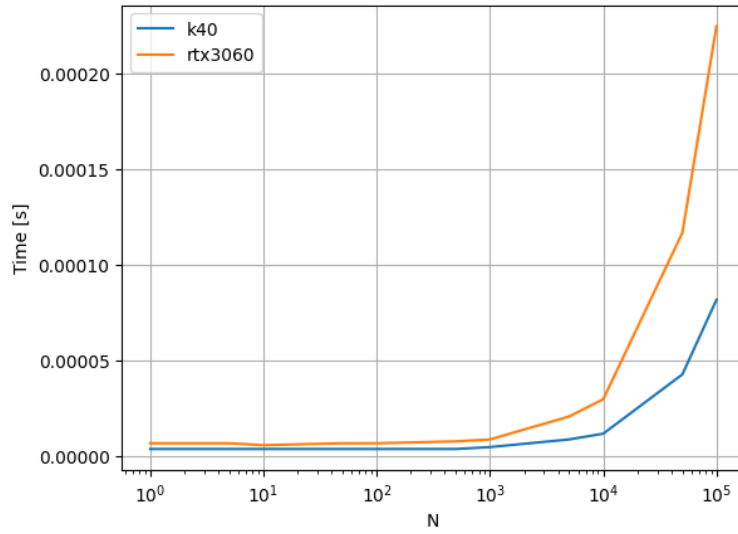


Figure 1: Times for `cudaMemcpy()` for different N

b)

The measurements were taken similar to a) but here we are simply measuring the call of an empty kernel. It is empty such that we only have the call and nothing else in the measurement.

```
1 __global__ void empty(int n, double *x, double *y, double *z){}
```

By looking at the constant part of the plot the following estimates were made:

- RTX: $t = 5\mu s$
- K40: $t = 11\mu s$

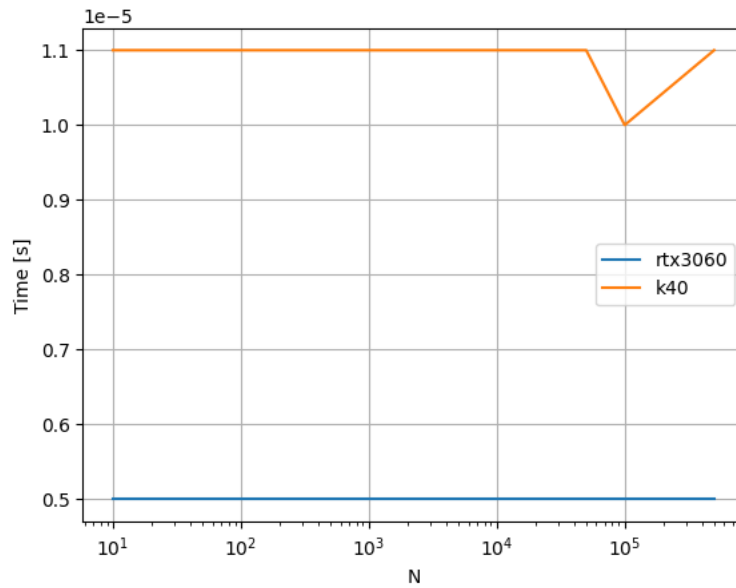


Figure 2: Times for kernel launch (256x256) for different N

c)

In the figure below we can see the peak memory bandwidth for both GPUs measured with the vector addition kernel we have already used in a previous exercise. The highest values measured for the bandwidths are:

- RTX: $B_{peak} = 331.446 \text{ GB/s}$
- K40: $B_{peak} = 167.598 \text{ GB/s}$

The GTX 3060 is allowing a total bandwidth of 360 GB/s, the K40 a bandwidth of 288 GB/s. This is quite close to our own measurement. Also, one can see that we have the typical S-shape of the plot:

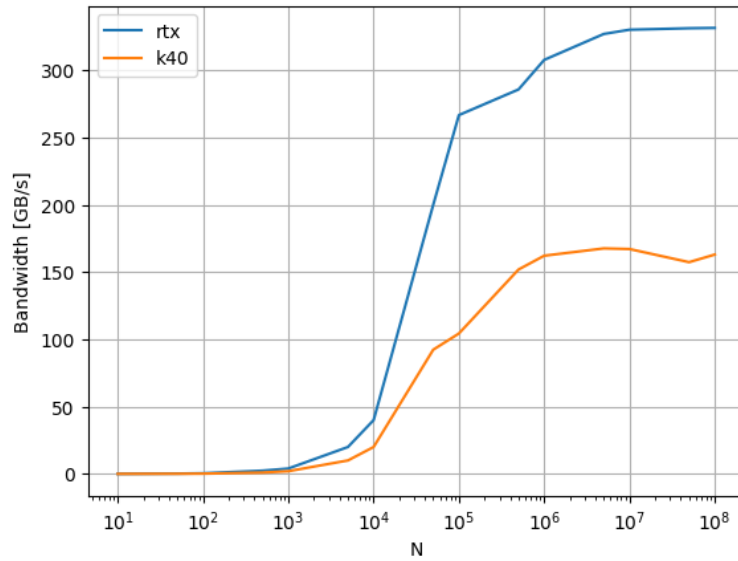


Figure 3: Bandwidth of the vector addition (256x256) for different N

d)

Here only a single thread (blocks x threads = 1x1) is doing the `atomicAdd()` (with an arbitrary input) such that we only have one computation at a time in a for loop which iterates N times.

```

1 __global__ void singleAdd(int n, double *x, double *y)
2 {
3     for (unsigned int i=0; i < n; i++)
4     {
5         atomicAdd(&x[i],1);
6     }
7 }

```

One can see in the figure below the number of `atomicAdd()` per second. The maximum is 1.81362×10^8 for $N = 10^8$.

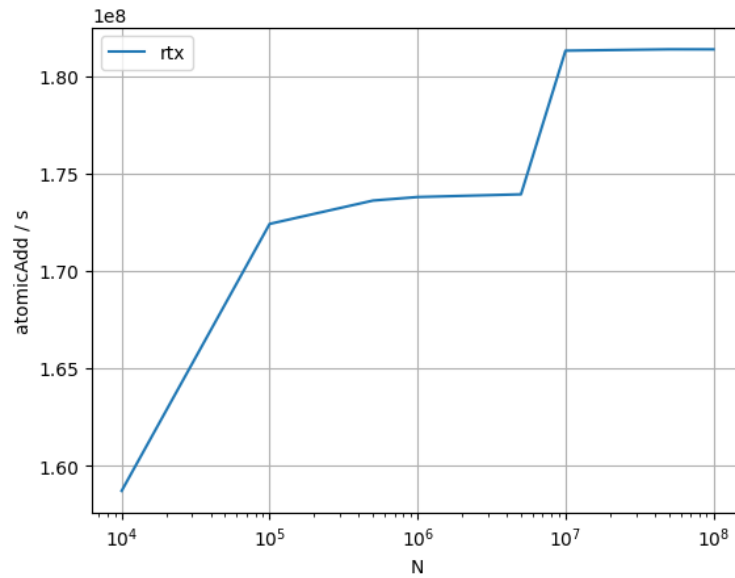


Figure 4: atomicAdd for different N (here N is the number of iterations the single thread has to do)

e)

Here we get the following peak floating point rates per second for the multiply-add (I simply reused the vector addition kernel - see code):

- RTX: 28.2187 GFlops/s
- K40: 12.221 GFlops/s

They were calculated as $2 * 8 * N/t * 10^{-9}$ GFlops/s.

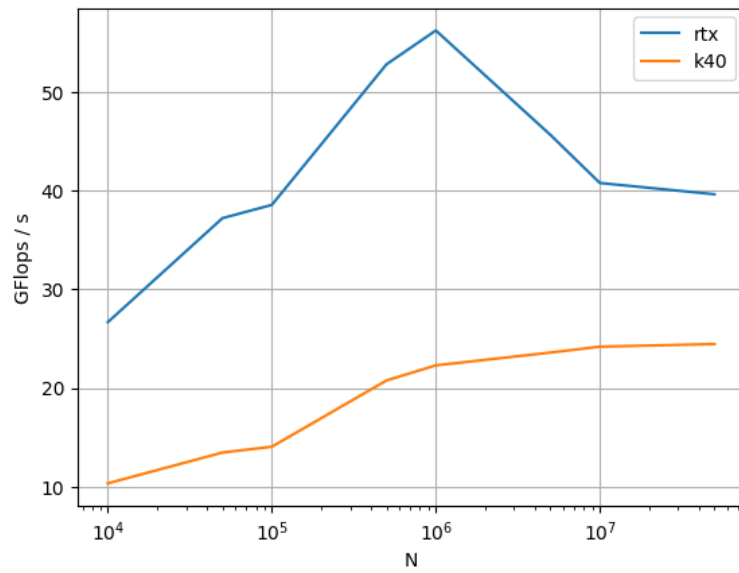


Figure 5: atomicAdd for different N (here N is the number of iterations the single thread has to do)

Appendix

Task 1.a)e)

```
1
2 #include <iostream>
3 #include <stdio.h>
4 #include <vector>
5 #include <algorithm>
6 #include "timer.hpp"
7
8 using Vector = std::vector<double>;
9
10 __global__ void addition(int n, double *x, double *y, double *z)
11 {
12     int id = blockIdx.x*blockDim.x + threadIdx.x;
13     for (size_t i = id; i < n; i += blockDim.x*gridDim.x){
14         //z[i] = x[i] + y[i];
15         z[i] += x[i] + y[i];
16     }
17 }
18
19 int main(void)
20 {
21     double *x, *y, *z, *d_x, *d_y, *d_z;
22     double elapsed;
23
24     Vector time_memcpy;
25     Vector flops;
26
27     Timer timer;
28     int reps = 6; // has to be an even number
29
30     Vector N_vec = {10000,50000,100000,500000,1000000,5000000,10000000,50000000};
31     for (const auto &N : N_vec){
32         Vector time_vec;
33
34         // Allocate host memory and initialize
35         x = (double*)malloc(N*sizeof(double));
36         y = (double*)malloc(N*sizeof(double));
37         z = (double*)malloc(N*sizeof(double));
38
39         for (int i = 0; i < N; i++) {
40             x[i] = i;
41             y[i] = N-1-i;
42             z[i] = 0;
43         }
44
45         // Allocate device memory and copy host data over
46         cudaMalloc(&d_x, N*sizeof(double));
47         cudaMalloc(&d_y, N*sizeof(double));
48         cudaMalloc(&d_z, N*sizeof(double));
49
50         // repetitions
51         for (int j = 0; j < reps; j++){
52
53             // MEASURING TIME FROM HERE
54             timer.reset();
55
56             cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
57
```

```

58         elapsed = timer.get();
59         // TO HERE
60         time_vec.push_back(elapsed);
61
62     }
63
64     cudaMemcpy(d_y, y, N*sizeof(double), cudaMemcpyHostToDevice);
65     cudaMemcpy(d_z, z, N*sizeof(double), cudaMemcpyHostToDevice);
66
67     addition<<<256, 256>>>(N, d_x, d_y, d_z);
68
69     std::sort(time_vec.begin(), time_vec.end());
70     time_memcpy.push_back(time_vec[reps/2]);
71     //bandwidth.push_back(3*(N)*sizeof(double)/time_vec[reps/2]*1e-9);
72     flops.push_back(2 * 8 * N / time_vec[reps/2] * 1e-9);
73
74     //printf("Elapsed: %g\n", time);
75
76     cudaMemcpy(z, d_z, N*sizeof(double), cudaMemcpyDeviceToHost);
77
78     cudaFree(d_x);
79     cudaFree(d_y);
80     cudaFree(d_z);
81     free(x);
82     free(y);
83     free(z);
84 }
85
86 printf("s:\n");
87 for (const auto& value : time_memcpy){
88     std::cout << value << ", " << std::endl;
89 }
90 printf("\nGFlops/s:\n");
91 for (const auto& value : flops){
92     std::cout << value << ", " << std::endl;
93 }
94
95 return EXIT_SUCCESS;
96 }

```

Task 1.b)c)

```

1
2 #include <iostream>
3 #include <stdio.h>
4 #include <vector>
5 #include <algorithm>
6 #include "timer.hpp"
7
8 using Vector = std::vector<double>;
9
10
11 __global__ void addition(int n, double *x, double *y, double *z)
12 {
13     int id = blockIdx.x*blockDim.x + threadIdx.x;
14     for (size_t i = id; i < n; i += blockDim.x*gridDim.x){
15         z[i] = x[i] + y[i];
16         //printf("z: %g\n", z[i*k]);
17     }
18 }

```



```

19
20 __global__ void empty(int n, double *x, double *y, double *z){}
21
22 int main(void)
23 {
24     double *x, *y, *z, *d_x, *d_y, *d_z;
25     double elapsed;
26
27     Vector bandwidth;
28     Vector time_kernel_launch;
29
30     Timer timer;
31     int reps = 10; // has to be an even number
32
33     Vector N_vec = {10,100,1000,10000,50000,100000,500000};
34     for (const auto &N : N_vec){
35         Vector time_vec;
36
37         // Allocate host memory and initialize
38         x = (double*)malloc(N*sizeof(double));
39         y = (double*)malloc(N*sizeof(double));
40         z = (double*)malloc(N*sizeof(double));
41
42         for (int i = 0; i < N; i++) {
43             x[i] = i;
44             y[i] = N-1-i;
45             z[i] = 0;
46         }
47
48         // Allocate device memory and copy host data over
49         cudaMalloc(&d_x, N*sizeof(double));
50         cudaMalloc(&d_y, N*sizeof(double));
51         cudaMalloc(&d_z, N*sizeof(double));
52
53         cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
54         cudaMemcpy(d_y, y, N*sizeof(double), cudaMemcpyHostToDevice);
55         cudaMemcpy(d_z, z, N*sizeof(double), cudaMemcpyHostToDevice);
56
57         // repetitions
58         for (int j = 0; j < reps; j++){
59
60             // MEASURING TIME FROM HERE
61             cudaDeviceSynchronize();
62             timer.reset();
63
64             empty<<<256, 256>>>(N, d_x, d_y, d_z);
65             //addition<<<256, 256>>>(N, d_x, d_y, d_z);
66
67             cudaDeviceSynchronize();
68             elapsed = timer.get();
69             time_vec.push_back(elapsed);
70
71             // TO HERE
72         }
73         std::sort(time_vec.begin(), time_vec.end());
74         time_kernel_launch.push_back(time_vec[reps/2]);
75         bandwidth.push_back(3*(N)*sizeof(double)/time_vec[reps/2]*1e-9);
76         //printf("Elapsed: %g\n", time);
77
78         cudaMemcpy(z, d_z, N*sizeof(double), cudaMemcpyDeviceToHost);
79
80         cudaFree(d_x);

```

```

81         cudaFree(d_y);
82         cudaFree(d_z);
83         free(x);
84         free(y);
85         free(z);
86     }
87
88     printf("s:\n");
89     for (const auto& value : time_kernel_launch){
90         std::cout << value << "," << std::endl;
91     }
92
93     printf("Gb/s:\n");
94     for (const auto& value : bandwidth){
95         std::cout << value << "," << std::endl;
96     }
97
98     return EXIT_SUCCESS;
99 }

```

Task 1.d)

```

1
2 #include <iostream>
3 #include <stdio.h>
4 #include <vector>
5 #include <algorithm>
6 #include "timer.hpp"
7
8 using Vector = std::vector<double>;
9
10 __global__ void singleAdd(int n, double *x, double *y)
11 {
12     for (unsigned int i=0; i < n; i++)
13     {
14         atomicAdd(&x[i],1);
15     }
16 }
17
18 int main(void)
19 {
20     double *x, *y, *d_x, *d_y;
21     double elapsed;
22
23     Vector add_per_second;
24
25     Timer timer;
26     int reps = 6; // has to be an even number
27
28     Vector N_vec = {10000,100000,500000,1000000,5000000,10000000,50000000,100000000};
29     for (const auto &N : N_vec)
30     {
31         Vector time_vec;
32
33         x = (double*)malloc(N*sizeof(double));
34         cudaMalloc(&d_x, N * sizeof(double));
35         cudaMemcpy(d_x, x, N * sizeof(double), cudaMemcpyHostToDevice);
36
37         y = (double*)malloc(N*sizeof(double));
38         cudaMalloc(&d_y, N * sizeof(double));

```

```

39     cudaMemcpy(d_y, y, N * sizeof(double), cudaMemcpyHostToDevice);
40
41     // repetitions
42     for (int j = 0; j < reps; j++)
43     {
44
45         // MEASURING TIME FROM HERE
46         cudaDeviceSynchronize();
47         timer.reset();
48
49         singleAdd<<<1, 1>>>(N, d_x, d_y);
50
51         cudaDeviceSynchronize();
52         elapsed = timer.get();
53         time_vec.push_back(elapsed);
54
55         // TO HERE
56     }
57     std::sort(time_vec.begin(), time_vec.end());
58     add_per_second.push_back(N/time_vec[reps / 2]);
59
60     cudaFree(d_x);
61     free(x);
62 }
63
64 printf("s:\n");
65 for (const auto &value : add_per_second)
66 {
67     std::cout << value << "," << std::endl;
68 }
69
70 return EXIT_SUCCESS;
71 }

```