# 360.252 - Computational Science on Many-Core Architectures

## WS 2020 - Exercise 3

Christian Gollmann, 01435044

Last update: November 7, 2020

# Contents

# 1 Strided Memory Access

I first wrote an implementation working on a small dataset to check if my kernel was operating correctly. Then I moved to the actual timing benchmark.

In both cases there can be seen a clear trend. I cannot really think of a reason why summing every second element is so much slower than say summing every or every sixth element. But in general it definitely makes sense that the runtime goes down when operating on less vector entries.

That the bandwidth goes down so drastically was something I didn't expect, but can also be explained. I therefor refer to an analogy when working on the CPU. There, the smallest amount of loadable memory is the cacheline. So even if only one 8 byte element is needed, the whole (mostly 64byte) cacheline hast to be loaded. That is why from a performance point of view, it is so important to keep an eye on data locality.

In the example here, we destroy the data locality concept by leaving empty spaces in between our vector usage. That's why we are no longer using what we have very efficiently.

In the runtime it can be seen that there are several plateaus, this effect becomes even clearer when turning k up to 127. This must also origin from the GPU's memory layout.

As a recommendation for more complex cases, I would say max out the use of data locality.
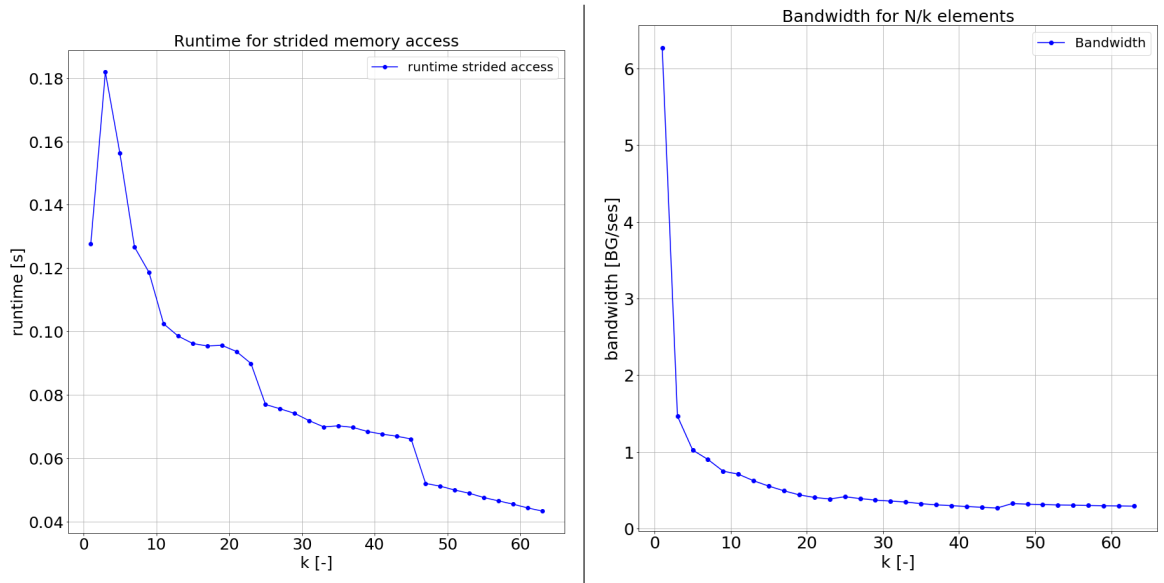


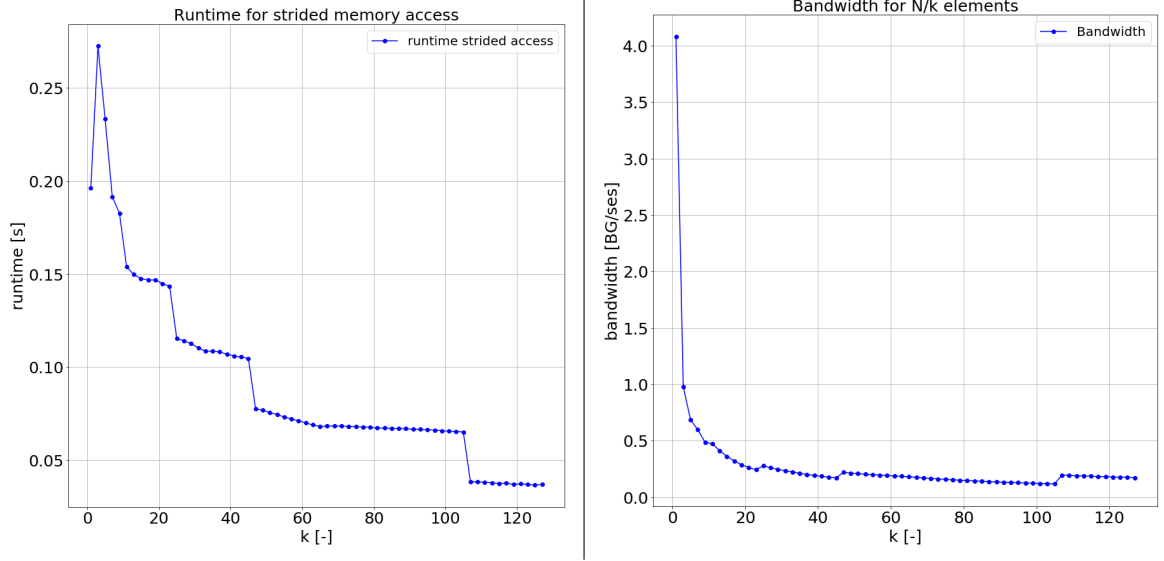Figure 1: Runtime and bandwidth for strided memory access

Figure 2: Runtime and bandwidth for strided memory access

Listing 1: Validation of Code for strided memory access

```cpp
#include <stdio.h>
#include "timer.hpp"
#include <iostream>
#include <algorithm>
#include <vector>


#define MAGIC_NUMBER 10

__global__ void sumVectors(double *x, double *y, double *z, int N, int k)
{
    int thread_id = blockIdx.x * blockDim.x + threadIdx.x;

    for(size_t i = thread_id; i < N/k; i += blockDim.x * gridDim.x)
        z[k*i] = x[k*i] + y[k*i];

}

int main(void)
{
    int N = 100;
    int k = 20;

    double*x, *y, *z, *d_x, *d_y, *d_z;
    Timer timer;

    x = new double[N];
    y = new double[N];
    z = new double[N];


    for (int i = 0; i < N; i++)
    {
        x[i] = 1;
        y[i] = 2;
```

```
36            z[i] = 0;
37        }
38
39
40        cudaMalloc(&d_x, N*sizeof(double));
41        cudaMalloc(&d_y, N*sizeof(double));
42        cudaMalloc(&d_z, N*sizeof(double));
43        cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
44        cudaMemcpy(d_y, y, N*sizeof(double), cudaMemcpyHostToDevice);
45        cudaMemcpy(d_z, z, N*sizeof(double), cudaMemcpyHostToDevice);
46
47        cudaDeviceSynchronize();
48        timer.reset();
49        std::vector<double> timings;
50
51        for(int reps=0; reps < MAGIC_NUMBER; ++reps)
52        {
53            sumVectors<<<256, 256>>>(d_x, d_y, d_z, N, k);
54            cudaDeviceSynchronize();
55            timings.push_back(timer.get());
56        }
57
58        std::sort(timings.begin(), timings.end());
59        double time_elapsed = timings[MAGIC_NUMBER/2];
60
61        cudaMemcpy(z, d_z, N*sizeof(double), cudaMemcpyDeviceToHost);
62
63        printf("Addition took %g seconds", time_elapsed);
64
65        std::cout << std::endl << "z[0] = " << z[0] << std::endl;
66        std::cout << "z[1] = " << z[1] << std::endl;
67        std::cout << "z[k] = " << z[k] << std::endl;
68        std::cout << "z[2*k] = " << z[2*k-1] << std::endl;
69        std::cout << "z[2*k+1] = " << z[2*k-1+1] << std::endl;
70
71        cudaFree(d_x);
72        cudaFree(d_y);
73        cudaFree(d_z);
74        delete x;
75        delete y;
76        delete z;
77
78        return EXIT_SUCCESS;
79    }
```

Listing 2: Actual Code used in benchmark

```
1  #include <stdio.h>
2  #include "timer.hpp"
3  #include <iostream>
4  #include <algorithm>
5  #include <vector>
6
7
8  #define MAGIC_NUMBER 10
9
10 __global__ void sumVectors(double *x, double *y, double *z, int N, int k)
```

```
11  {
12      int thread_id = blockIdx.x * blockDim.x + threadIdx.x;
13
14      for(size_t i = thread_id; i < N/k; i += blockDim.x * gridDim.x)
15          z[k*i] = x[k*i] + y[k*i];
16
17  }
18
19  int main(void)
20  {
21      int N = 100000000;
22
23      double *x, *y, *z, *d_x, *d_y, *d_z;
24      Timer timer;
25
26      x = new double[N];
27      y = new double[N];
28      z = new double[N];
29
30
31      for (int i = 0; i < N; i++)
32      {
33          x[i] = 1;
34          y[i] = 2;
35          z[i] = 0;
36      }
37
38
39      cudaMalloc(&d_x, N*sizeof(double));
40      cudaMalloc(&d_y, N*sizeof(double));
41      cudaMalloc(&d_z, N*sizeof(double));
42      cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
43      cudaMemcpy(d_y, y, N*sizeof(double), cudaMemcpyHostToDevice);
44      cudaMemcpy(d_z, z, N*sizeof(double), cudaMemcpyHostToDevice);
45
46      for(int k = 1; k < 64; k += 2)
47      {
48          cudaDeviceSynchronize();
49          timer.reset();
50          std::vector<double> timings;
51
52          for(int reps=0; reps < MAGIC_NUMBER; ++reps)
53          {
54              sumVectors<<<256, 256>>>(d_x, d_y, d_z, N, k);
55              cudaDeviceSynchronize();
56              timings.push_back(timer.get());
57          }
58
59          std::sort(timings.begin(), timings.end());
60          double time_elapsed = timings[MAGIC_NUMBER/2];
61
62          std::cout << time_elapsed << std::endl;
63      }
64
65      cudaFree(d_x);
66      cudaFree(d_y);
67      cudaFree(d_z);
```

```
68        delete x;
69        delete y;
70        delete z;
71
72        return EXIT_SUCCESS;
73   }
```
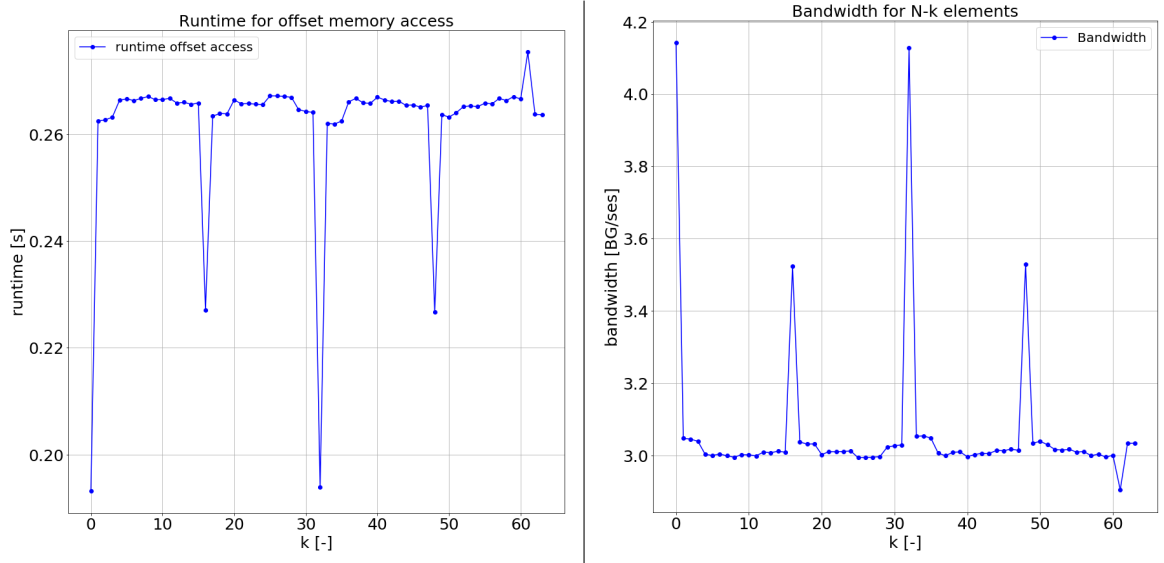
## 2 Offset Memory Access



Figure 3: Runtime and bandwidth for offset memory access

Listing 3: Code for offset access

```cpp
1  #include <stdio.h>
2  #include "timer.hpp"
3  #include <iostream>
4  #include <algorithm>
5  #include <vector>
6
7
8  #define MAGIC_NUMBER 10
9
10 __global__ void sumVectors(double *x, double *y, double *z, int N, int k)
11 {
12     int thread_id = blockIdx.x * blockDim.x + threadIdx.x;
13
14     for(size_t i = thread_id; i < N-k; i += blockDim.x * gridDim.x)
15         z[k+i] = x[k+i] + y[k+i];
16
17 }
18
19 int main(void)
20 {
21     int N = 100000000;
22
23     double *x, *y, *z, *d_x, *d_y, *d_z;
24     Timer timer;
25
26     x = new double[N];
27     y = new double[N];
28     z = new double[N];
29
30
31     for (int i = 0; i < N; i++)
```

```
32        {
33            x[i] = 1;
34            y[i] = 2;
35            z[i] = 0;
36        }
37
38
39        cudaMalloc(&d_x, N*sizeof(double));
40        cudaMalloc(&d_y, N*sizeof(double));
41        cudaMalloc(&d_z, N*sizeof(double));
42        cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
43        cudaMemcpy(d_y, y, N*sizeof(double), cudaMemcpyHostToDevice);
44        cudaMemcpy(d_z, z, N*sizeof(double), cudaMemcpyHostToDevice);
45
46        for(int k = 0; k < 64; k += 1)
47        {
48            cudaDeviceSynchronize();
49            timer.reset();
50            std::vector<double> timings;
51
52            for(int reps=0; reps < MAGIC_NUMBER; ++reps)
53            {
54                sumVectors<<<256, 256>>>(d_x, d_y, d_z, N, k);
55                cudaDeviceSynchronize();
56                timings.push_back(timer.get());
57            }
58
59            std::sort(timings.begin(), timings.end());
60            double time_elapsed = timings[MAGIC_NUMBER/2];
61
62            std::cout << time_elapsed << std::endl;
63        }
64
65        cudaFree(d_x);
66        cudaFree(d_y);
67        cudaFree(d_z);
68        delete x;
69        delete y;
70        delete z;
71
72        return EXIT_SUCCESS;
73 }
```

# 3  Conjugate Gradient - matrix vector product

Here I went with the example from the lecture. But I also took much inspiration from
https://medium.com/analytics-vidhya/sparse-matrix-vector-multiplication-with-cuda-42d191878e8f
which turned out to be a good reference also for later and more sophisticated implimentations.

Listing 4: CUDA kernel for calculating sparse matrix vector multiplication

```
1  __global__ void csr_matvec(int N, int *rowoffsets, int *colindices, double *values,
2  double const *x, double *y)
3  {
4      for (int row = blockDim.x * blockIdx.x + threadIdx.x;
5          row < N;
6          row += gridDim.x * blockDim.x)
7      {
8        double val = 0;
9        for (int jj = rowoffsets[i]; jj < rowoffsets[i+1]; ++jj)
10        {
11            val += values[jj] * x[colindices[jj]];
12        }
13        y[row] = val;
14      }
15  }
```

# 4 Conjugate Gradient - further kernels

Listing 5: kernel for the vector operations in lines 5 and 9

```
1  __global__ void dot_product(double *x, double *y, double *dot, unsigned int n)
2  {
3      unsigned int index = threadIdx.x + blockDim.x*blockIdx.x;
4      unsigned int stride = blockDim.x*gridDim.x;
5
6      __shared__ double cache[256];
7
8      double temp = 0.0;
9      while(index < n){
10         temp += x[index]*y[index];
11
12         index += stride;
13     }
14
15     cache[threadIdx.x] = temp;
16
17     __syncthreads();
18
19     for(int i = blockDim.x/2; i>0; i/=2)
20     {
21         __syncthreads();
22         if(threadIdx.x < i)
23             cache[threadIdx.x] += cache[threadIdx.x + i];
24     }
25
26     if(threadIdx.x == 0){
27         atomicAdd(dot, cache[0]);
28     }
29  }
```

Listing 6: kernel for the vector operations in lines 7, 8 and 12

```
1  __global__ void vector_plus_alpha_vector(double *x, double *y,
2  double *z, double alpha, int N)
3  {
4      int thread_id = blockIdx.x * blockDim.x + threadIdx.x;
5
6      for(size_t i = thread_id; i < N; i += blockDim.x * gridDim.x)
7          z[i] = x[i] + alpha * y[i];
8
9  }
```

# 5 Basic CUDA e)

Now I called the vector addition for a vector with e7 elements with varying grid and block sizes.
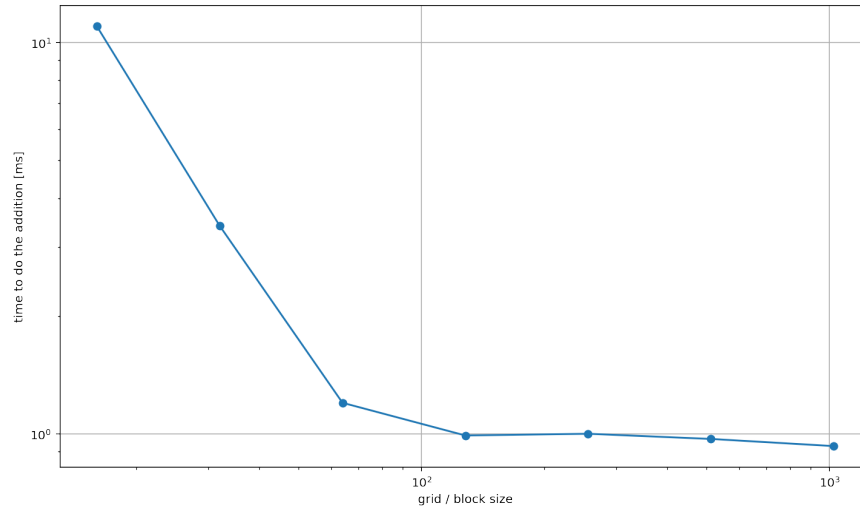It seems that small values (16, 32, 64) lead to a not so good performance.



Figure 4: addition time for vector with e7 elements for different grid / block values

Listing 7: kernel call with different values x

```
1  sumVectors<<<x, x>>>(d_x, d_y, d_z, N);
```

# 6 Dot Product a)

For the Dot Product Exercises I got some inspiration from https://bitbucket.org/jsandham/ algorithms_in_cuda/src/master/dot_product/ and added / changed code as needed. I am aware that my implementations are not the prettiest and might also be error prone. But they get the job done for those specific examples.

Listing 8: Dot Product with two GPU stages

```cpp
#include <stdio.h>
#include <iostream>
#include "timer.hpp"
#include <random>


__global__ void dot_product_first_part(double *x, double *y,
double *temporary, unsigned int n)
{
    unsigned int index = threadIdx.x + blockDim.x*blockIdx.x;
    unsigned int stride = blockDim.x*gridDim.x;

    __shared__ double cache[256];

    double temp = 0.0;
    while(index < n){
        temp += x[index]*y[index];

        index += stride;
    }

    cache[threadIdx.x] = temp;

    __syncthreads();

    for(int i = blockDim.x/2; i>0; i/=2)
    {
        __syncthreads();
        if(threadIdx.x < i)
            cache[threadIdx.x] += cache[threadIdx.x + i];
    }

    if(threadIdx.x == 0){
        temporary[blockIdx.x] = cache[0];
    }
}

__global__ void dot_product_second_part(double* temporary, double* dot)
{
    for(int i = blockDim.x/2; i>0; i/=2)
    {
        __syncthreads();
        if(threadIdx.x < i)
            temporary[threadIdx.x] += temporary[threadIdx.x + i];
    }

    __syncthreads();
```

```
48
49        if(threadIdx.x == 0){
50            *dot = temporary[0];
51        }
52
53  }
54
55  int main()
56  {
57        unsigned int n = 10000;
58        unsigned int x = 256;
59        double *h_prod;
60        double *d_prod;
61        double *h_x, *h_y;
62        double *d_x, *d_y;
63        double *d_temporary;
64        Timer timer;
65
66        h_prod = new double[n];
67        h_x = new double[n];
68        h_y = new double[n];
69
70
71        // fill host array with data
72        for(unsigned int i=0;i<n;i++){
73            h_x[i] = 1;
74            h_y[i] = 2;
75        }
76
77        // start timer
78        timer.reset();
79
80        // allocate memory
81        cudaMalloc(&d_prod, sizeof(double));
82        cudaMalloc(&d_x, n*sizeof(double));
83        cudaMalloc(&d_y, n*sizeof(double));
84        cudaMalloc(&d_temporary, x*sizeof(double));
85
86
87        // copy data to device
88        cudaMemcpy(d_x, h_x, n*sizeof(double), cudaMemcpyHostToDevice);
89        cudaMemcpy(d_y, h_y, n*sizeof(double), cudaMemcpyHostToDevice);
90
91
92        dot_product_first_part<<<x, x>>>(d_x, d_y, d_temporary, n);
93        dot_product_second_part<<<1, x>>>(d_temporary, d_prod);
94
95
96        // copy data back to host
97        cudaMemcpy(h_prod, d_prod, sizeof(double), cudaMemcpyDeviceToHost);
98
99        // get runtime
100       double time_elapsed = timer.get();
101
102
103       // report results
104       std::cout<<"dot product computed on GPU is: "<<*h_prod<<" and took "
```

```
105        << time_elapsed << " s" <<std::endl;
106
107
108        // free memory
109        free(h_prod);
110        free(h_x);
111        free(h_y);
112        cudaFree(d_prod);
113        cudaFree(d_x);
114        cudaFree(d_y);
115
116 }
```

# 7 Dot Product b)

```cpp
1  #include <stdio.h>
2  #include <iostream>
3  #include "timer.hpp"
4  #include <random>
5
6
7  __global__ void dot_product(double *x, double *y, double *temporary, unsigned int n)
8  {
9      unsigned int index = threadIdx.x + blockDim.x*blockIdx.x;
10     unsigned int stride = blockDim.x*gridDim.x;
11
12     __shared__ double cache[256];
13
14     double temp = 0.0;
15     while(index < n){
16         temp += x[index]*y[index];
17
18         index += stride;
19     }
20
21     cache[threadIdx.x] = temp;
22
23     __syncthreads();
24
25     for(int i = blockDim.x/2; i>0; i/=2)
26     {
27         __syncthreads();
28         if(threadIdx.x < i)
29             cache[threadIdx.x] += cache[threadIdx.x + i];
30     }
31
32     if(threadIdx.x == 0){
33         temporary[blockIdx.x] = cache[0];
34     }
35 }
36
37
38
39 int main()
40 {
41     unsigned int n = 10000;
42     unsigned int x = 256;
43     double *h_prod;
44     double *d_prod;
45     double *h_x, *h_y;
46     double *d_x, *d_y;
47     double *d_temporary;
48     double *h_temporary;
49     Timer timer;
50
51     h_prod = new double[n];
52     h_x = new double[n];
53     h_y = new double[n];
```

```
54        h_temporary = new double[x];
55
56
57        // fill host array with data
58        for(unsigned int i=0;i<n;i++){
59            h_x[i] = 1;
60            h_y[i] = 2;
61        }
62
63        // start timer
64        timer.reset();
65
66        // allocate memory
67        cudaMalloc(&d_prod, sizeof(double));
68        cudaMalloc(&d_x, n*sizeof(double));
69        cudaMalloc(&d_y, n*sizeof(double));
70        cudaMalloc(&d_temporary, x*sizeof(double));
71        cudaMemset(d_prod, 0.0, sizeof(double));
72
73
74        // copy data to device
75        cudaMemcpy(d_x, h_x, n*sizeof(double), cudaMemcpyHostToDevice);
76        cudaMemcpy(d_y, h_y, n*sizeof(double), cudaMemcpyHostToDevice);
77
78
79        dot_product<<<x, x>>>(d_x, d_y, d_temporary, n);
80
81
82        // copy data back to host
83        cudaMemcpy(h_temporary, d_temporary, x*sizeof(double), cudaMemcpyDeviceToHost);
84
85        // sum up elements
86        double dot = 0;
87        for(int i = 0; i < x; i++)
88        {
89            dot += h_temporary[i];
90        }
91
92        // get runtime
93        double time_elapsed = timer.get();
94
95
96        // report results
97        std::cout<<"dot product computed on GPU and CPU is: "<<dot<<" and took "
98        << time_elapsed << " s" <<std::endl;
99
100
101       // free memory
102       free(h_prod);
103       free(h_x);
104       free(h_y);
105       cudaFree(d_prod);
106       cudaFree(d_x);
107       cudaFree(d_y);
108
109   }
```

# 8 Dot Product c)

```cpp
1  #include <stdio.h>
2  #include <iostream>
3  #include "timer.hpp"
4  #include <random>
5
6
7  __global__ void dot_product(double *x, double *y, double *dot, unsigned int n)
8  {
9      unsigned int index = threadIdx.x + blockDim.x*blockIdx.x;
10     unsigned int stride = blockDim.x*gridDim.x;
11
12     __shared__ double cache[256];
13
14     double temp = 0.0;
15     while(index < n){
16         temp += x[index]*y[index];
17
18         index += stride;
19     }
20
21     cache[threadIdx.x] = temp;
22
23     __syncthreads();
24
25     for(int i = blockDim.x/2; i>0; i/=2)
26     {
27         __syncthreads();
28         if(threadIdx.x < i)
29             cache[threadIdx.x] += cache[threadIdx.x + i];
30     }
31
32     if(threadIdx.x == 0){
33         atomicAdd(dot, cache[0]);
34     }
35 }
36
37
38
39 int main()
40 {
41     unsigned int n = 10000;
42     double *h_prod;
43     double *d_prod;
44     double *h_x, *h_y;
45     double *d_x, *d_y;
46     Timer timer;
47
48     h_prod = new double[n];
49     h_x = new double[n];
50     h_y = new double[n];
51
52
53     // fill host array with data
```

16

```
54        for(unsigned int i=0;i<n;i++){
55            h_x[i] = 1;
56            h_y[i] = 2;
57        }
58
59        // start timer
60        timer.reset();
61
62        // allocate memory
63        cudaMalloc(&d_prod, sizeof(double));
64        cudaMalloc(&d_x, n*sizeof(double));
65        cudaMalloc(&d_y, n*sizeof(double));
66        cudaMemset(d_prod, 0.0, sizeof(double));
67
68
69        // copy data to device
70        cudaMemcpy(d_x, h_x, n*sizeof(double), cudaMemcpyHostToDevice);
71        cudaMemcpy(d_y, h_y, n*sizeof(double), cudaMemcpyHostToDevice);
72
73
74        dot_product<<<256, 256>>>(d_x, d_y, d_prod, n);
75
76        // copy data back to host
77        cudaMemcpy(h_prod, d_prod, sizeof(double), cudaMemcpyDeviceToHost);
78
79        // get runtime
80        double time_elapsed = timer.get();
81
82
83        // report results
84        std::cout<<"dot product computed on GPU is: "<<*h_prod<<" and took "
85        << time_elapsed << " s" <<std::endl;
86
87
88        // free memory
89        free(h_prod);
90        free(h_x);
91        free(h_y);
92        cudaFree(d_prod);
93        cudaFree(d_x);
94        cudaFree(d_y);
95
96 }
```

# 9 Dot Product d)

It can be seen that at least in my case, there is no difference in runtime for the different methods a to c. This surprises me but maybe is due to how I time the executions. Though, I wanted to take all the steps necessary (also the cudaMalloc and cudaMemcpy) into account. Again there are some similarities in respect to the number of vector entries with Figure 1 and Figure 2. So this behaviour could be memory related.
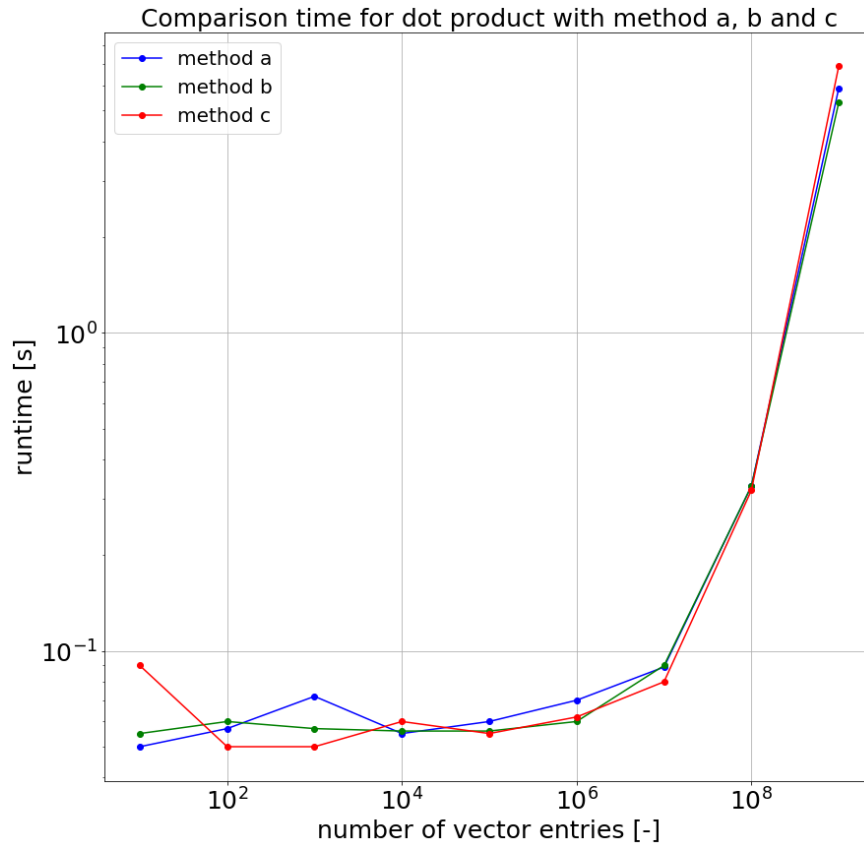


Figure 5: Runtimes for calculating dot product with different methods.