# 360.252 - Computational Science on Many-Core Architectures

## WS 2020 - Exercise 2

Christian Gollmann, 01435044

Last update: November 1, 2020

# Contents

# 1   Basic CUDA a)

Here I have to apologize for the small figure annotations, I only realized how small they are after already having included the figure and I didn't have the figure code anymore.

It can be seen that up to a certain point, the size of allocated memory has no impact on performance, after that the impact scales with the needed memory. This probably is due to the GPU's memory system. There might be an equivalence to the CPU's caches.
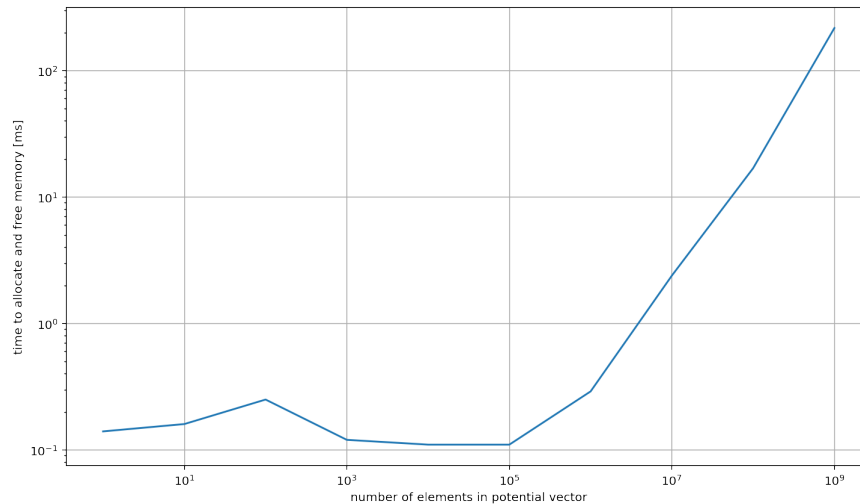


Figure 1: time to allocate and free memory for different sizes N

Listing 1: Code used to generate the output above

```
1   #include <stdio.h>
2   #include "timer.hpp"
3
4   int main(void)
5   {
6       int N = 1000;
7       double *d_x;
8       Timer timer;
9
10      cudaDeviceSynchronize();
11      timer.reset();
12
13      for(int i = 0; i < 100; i++) {
14          cudaMalloc(&d_x, N*sizeof(double));
15          cudaFree(d_x);
16          cudaDeviceSynchronize();
17      }
18
19      cudaDeviceSynchronize();
20      double time_elapsed = timer.get();
21
22      printf("Allocation and Freeing took %g seconds", time_elapsed/100);
23
24      return EXIT_SUCCESS;
25  }
```

## 2 Basic CUDA b)

Here, clearly copying each entry individually is the slowest as I expected it to be.

| Initializing vectors for 1M elements | | |
|---|---|---|
| Method | time to initialize [ms] | bandwidth [MB/s] |
| Initialize in kernel | 13.8 | 579.7 |
| Copy the data | 3 | 2666.7 |
| Copy each individual entry | 7742 | 1 |

Listing 2: Initialize directly within a dedicated CUDA kernel

```
1  #include <stdio.h>
2  #include "timer.hpp"
3
4  __global__
5  void init(int N)
6  {
7      double *x, *y;
8
9      x = new double[N];
10     y = new double[N];
11
12     for (int i = 0; i < N; i++)
13     {
14         x[i] = i;
15         y[i] = N-1-i;
16     }
17 }
18
19 int main(void)
20 {
21     int M = 1;
22     int N = 1000000;
23     Timer timer;
24
25     cudaDeviceSynchronize();
26     timer.reset();
27
28     init<<<(M+255)/256, 256>>>(N);
29
30     cudaDeviceSynchronize();
31     double time_elapsed = timer.get();
32
33     printf("Initializing in kernel took %g seconds", time_elapsed);
34
35     return EXIT_SUCCESS;
36 }
```

Listing 3: Copy the data

```
1  #include <stdio.h>
2  #include "timer.hpp"
3
4  int main(void)
5  {
6      int N = 1000000;
```

```
7
8        double *x, *y, *d_x, *d_y;
9        Timer timer;
10
11       x = new double[N];
12       y = new double[N];
13
14       for (int i = 0; i < N; i++)
15       {
16           x[i] = i;
17           y[i] = N-1-i;
18       }
19
20       cudaDeviceSynchronize();
21       timer.reset();
22
23       cudaMalloc(&d_x, N*sizeof(double));
24       cudaMalloc(&d_y, N*sizeof(double));
25       cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
26       cudaMemcpy(d_y, y, N*sizeof(double), cudaMemcpyHostToDevice);
27
28       cudaDeviceSynchronize();
29       double time_elapsed = timer.get();
30
31       printf("Initializing outside took %g seconds", time_elapsed);
32
33       cudaFree(d_x);
34       cudaFree(d_y);
35       delete x;
36       delete y;
37
38       return EXIT_SUCCESS;
39  }
```

Listing 4: Copy each individual entry

```
1   #include <stdio.h>
2   #include "timer.hpp"
3
4   int main(void)
5   {
6        int N = 1000000;
7
8        double *x, *y, *d_x, *d_y;
9        Timer timer;
10
11       x = new double[N];
12       y = new double[N];
13
14       for (int i = 0; i < N; i++)
15       {
16           x[i] = i;
17           y[i] = N-1-i;
18       }
19
20       cudaDeviceSynchronize();
21       timer.reset();
```

```
22
23        cudaMalloc(&d_x, N*sizeof(double));
24        cudaMalloc(&d_y, N*sizeof(double));
25        for (int i = 0; i < N; i++)
26        {
27             cudaMemcpy(d_x+i, x+i, sizeof(double), cudaMemcpyHostToDevice);
28             cudaMemcpy(d_y+i, y+i, sizeof(double), cudaMemcpyHostToDevice);
29        }
30
31        cudaDeviceSynchronize();
32        double time_elapsed = timer.get();
33
34        printf("Initialising and copying by piece took %g seconds", time_elapsed);
35
36        cudaFree(d_x);
37        cudaFree(d_y);
38        delete x;
39        delete y;
40
41        return EXIT_SUCCESS;
42 }
```

# 3 Basic CUDA c)

Listing 5: CUDA kernel that sums two vectors

```cpp
1  #include <stdio.h>
2  #include "timer.hpp"
3  #include <iostream>
4
5  __global__ void sumVectors(double *x, double *y, double *z, int N)
6  {
7      int thread_id = blockIdx.x * blockDim.x + threadIdx.x;
8
9      for(size_t i = thread_id; i < N; i += blockDim.x * gridDim.x)
10         z[i] = x[i] + y[i];
11
12 }
13
14 int main(void)
15 {
16     int N = 100;
17
18     double *x, *y, *z, *d_x, *d_y, *d_z;
19     Timer timer;
20
21     x = new double[N];
22     y = new double[N];
23     z = new double[N];
24
25
26     for (int i = 0; i < N; i++)
27     {
28         x[i] = i;
29         y[i] = N-1-i;
30         z[i] = 0;
31     }
32
33
34     cudaMalloc(&d_x, N*sizeof(double));
35     cudaMalloc(&d_y, N*sizeof(double));
36     cudaMalloc(&d_z, N*sizeof(double));
37     cudaMemcpy(d_x, x, N*sizeof(double), cudaMemcpyHostToDevice);
38     cudaMemcpy(d_y, y, N*sizeof(double), cudaMemcpyHostToDevice);
39     cudaMemcpy(d_z, z, N*sizeof(double), cudaMemcpyHostToDevice);
40
41     cudaDeviceSynchronize();
42     timer.reset();
43
44     sumVectors<<<(N+255)/256, 256>>>(d_x, d_y, d_z, N);
45
46     cudaDeviceSynchronize();
47     double time_elapsed = timer.get();
48
49     cudaMemcpy(z, d_z, N*sizeof(double), cudaMemcpyDeviceToHost);
50
51     printf("Addition took %g seconds", time_elapsed);
52
53     std::cout << std::endl << "z[0] = " << z[0] << std::endl;
```

```
54
55        cudaFree ( d_x ) ;
56        cudaFree ( d_y ) ;
57        cudaFree ( d_z ) ;
58        delete  x ;
59        delete  y ;
60        delete  z ;
61
62        return  EXIT_SUCCESS ;
63    }
```

## 4  Basic CUDA d)

For small values of N ($<$e5) it doesn't make much difference but after hitting a certain threshold, execution time seems to increase exponentially. Consider though that I always call the kernel with different values of N but still those are interesting results. There's much resemblance to the plot in Figure 1.
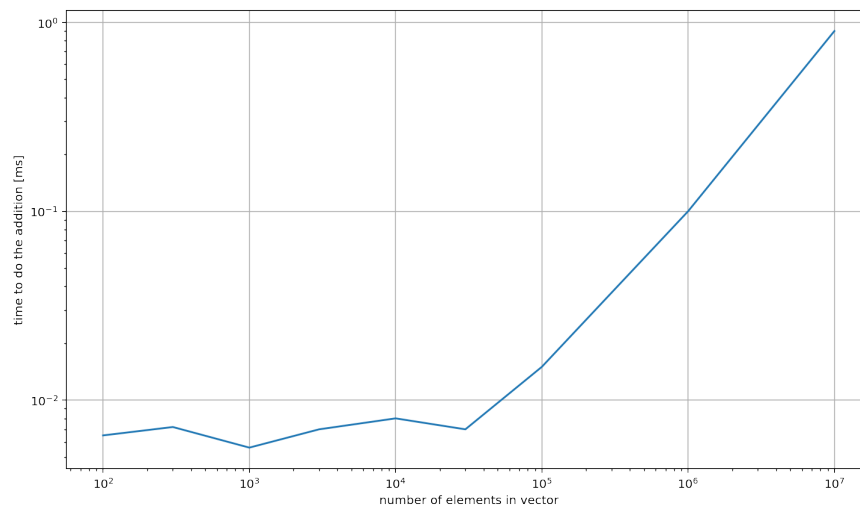


Figure 2: addition time for different values of N

Listing 6: addition benchmark

```
1   #include <stdio.h>
2   #include "timer.hpp"
3   #include <iostream>
4
5   __global__ void sumVectors(double *x, double *y, double *z, int N)
6   {
7       int thread_id = blockIdx.x * blockDim.x + threadIdx.x;
8
9       for(size_t i = thread_id; i < N; i += blockDim.x * gridDim.x)
10          z[i] = x[i] + y[i];
11
12  }
13
14  int main(void)
15  {
16      int N = 100;
17
18      double *x, *y, *z, *d_x, *d_y, *d_z;
19      Timer timer;
20
21      x = new double[N];
22      y = new double[N];
23      z = new double[N];
24
25
26      for (int i = 0; i < N; i++)
27      {
```

```
28              x [ i ] = i ;
29              y [ i ] = N−1−i ;
30              z [ i ] = 0;
31          }
32
33
34          cudaMalloc(&d_x , N∗sizeof (double ));
35          cudaMalloc(&d_y , N∗sizeof (double ));
36          cudaMalloc(&d_z , N∗sizeof (double ));
37          cudaMemcpy( d_x , x , N∗sizeof (double ) , cudaMemcpyHostToDevice );
38          cudaMemcpy( d_y , y , N∗sizeof (double ) , cudaMemcpyHostToDevice );
39          cudaMemcpy( d_z , z , N∗sizeof (double ) , cudaMemcpyHostToDevice );
40
41          cudaDeviceSynchronize ( );
42          timer . reset ( );
43          for ( int i = 0; i < 100; i++) {
44              sumVectors<<<(N+255)/256 , 256>>>(d_x , d_y , d_z , N);
45              cudaDeviceSynchronize ( );
46          }
47
48          cudaDeviceSynchronize ( );
49          double time_elapsed = timer . get ( );
50
51          cudaMemcpy( z , d_z , N∗sizeof (double ) , cudaMemcpyDeviceToHost );
52
53          printf (" Addition took %g seconds " , time_elapsed /100);
54
55          std :: cout << std :: endl << " z [ 0 ] = " << z [ 0 ] << std :: endl ;
56
57          cudaFree ( d_x );
58          cudaFree ( d_y );
59          cudaFree ( d_z );
60          delete x ;
61          delete y ;
62          delete z ;
63
64          return EXIT_SUCCESS;
65  }
```

# 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.
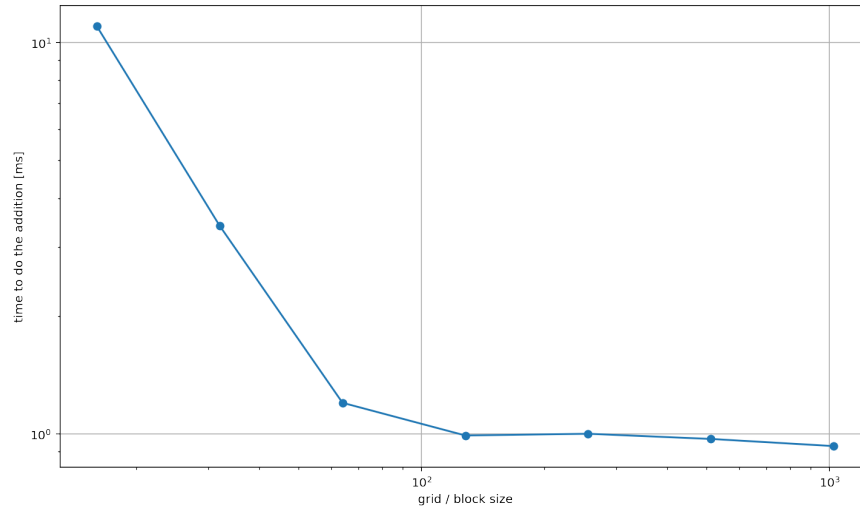


Figure 3: addition time for vector with e7 elements for different grid / block values

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/ algo-rithms_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

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

```cpp
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)

Listing 9: Dot Product with GPU and CPU combined

```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];
```

13

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


__global__ void dot_product(double *x, double *y, double *dot, 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){
        atomicAdd(dot, cache[0]);
    }
}



int main()
{
    unsigned int n = 10000;
    double *h_prod;
    double *d_prod;
    double *h_x, *h_y;
    double *d_x, *d_y;
    Timer timer;

    h_prod = new double[n];
    h_x = new double[n];
    h_y = new double[n];


    // fill host array with data
```

```
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.
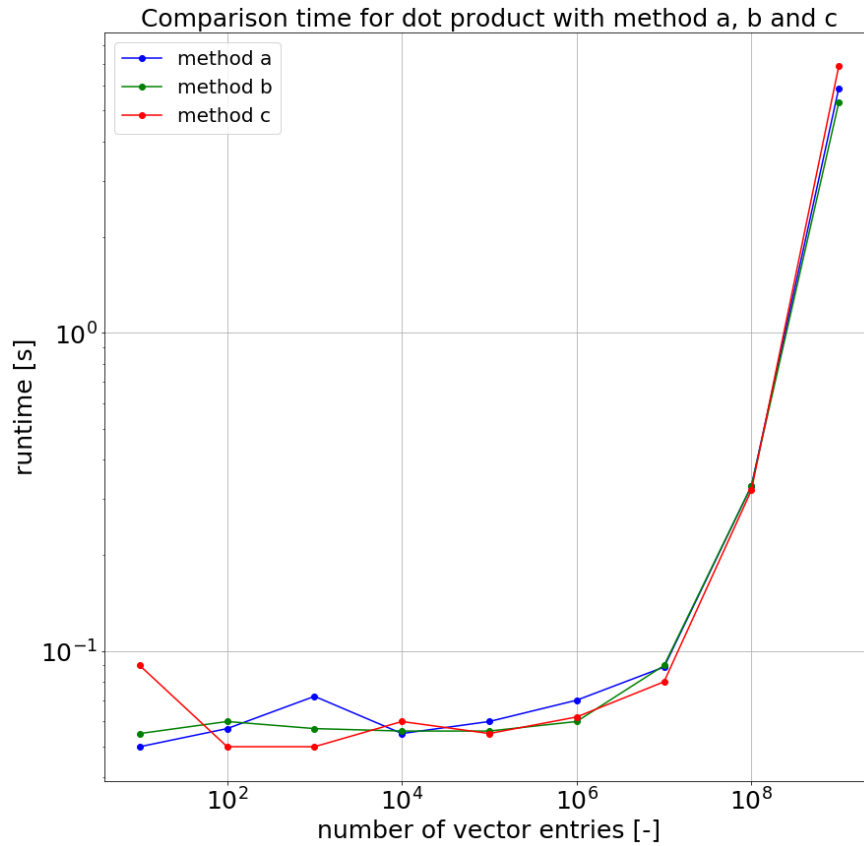


Figure 4: Runtimes for calculating dot product with different methods.