# 360.252 - Computational Science on Many-Core Architectures

## WS 2020 - Exercise 3

Christian Gollmann, 01435044

Last update: November 9, 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. 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 therefore 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 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.

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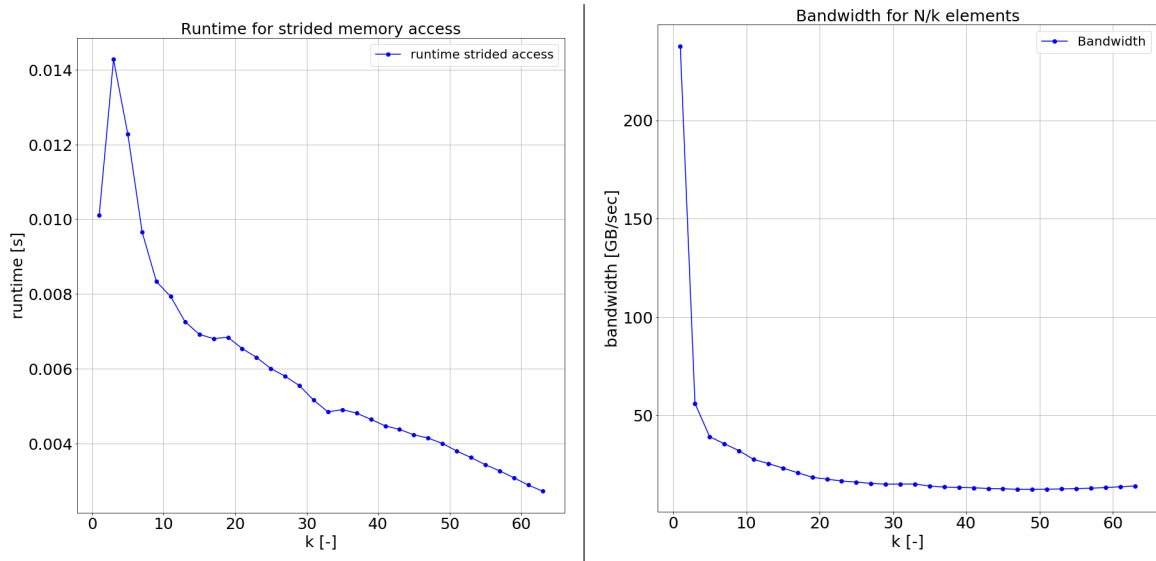


Figure 1: Runtime and bandwidth for strided memory access

```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 = 100;
22      int k = 20;
23
24      double *x, *y, *z, *d_x, *d_y, *d_z;
25      Timer timer;
26
27      x = new double[N];
28      y = new double[N];
29      z = new double[N];
30
31
32      for (int i = 0; i < N; i++)
33      {
34          x[i] = 1;
35          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      }
```

```cpp
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

```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 = 1; k < 64; k += 2)
47      {
48          cudaDeviceSynchronize();
49          std::vector<double> timings;
50
51          for(int reps=0; reps < MAGIC_NUMBER; ++reps)
52          {
53              timer.reset();
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

Two things can be recognized from the following plot:

1. The runtime, and so the bandwidth, are "good" whenever k equals a multiple of sixteen (0, 16, 32, ...).

2. A slow runtime does not really depend on k, eg. it does not matter whether we skip 3 or 15 entries.

I strongly suspect that the GPU's cacheline size is 8*16 = 128 bytes. If, let's say, we skip the first 5 entries, we don't use the cacheline at its most efficiency and therefore have some kind of cache misses. Those cache misses then influence the achievable bandwidth. I'm not completely sure how all of this ties together, but I'm positive that it origins from the cacheline concept.

For future implementations, I suggest always using a "full" cacheline if possible, thereby referring to the principle of data locality from before.

As an update, after completing the exercises, I found out more about Global Memory Coalescing on https://developer.nvidia.com/blog/how-access-global-memory-efficiently-cuda-c-kernels/. There I read that during execution, threads are grouped into warps whose warp size is 32 threads. On the webpage there is also more or less the same example we are supposed to do here. I was happy to see, that I was able to reproduce the results and that my assumptions more or less aligned with what was explained there. As I don't intend to just reproduce the content there, I decided to stick with my initial assumptions and give this finding now as a reference.
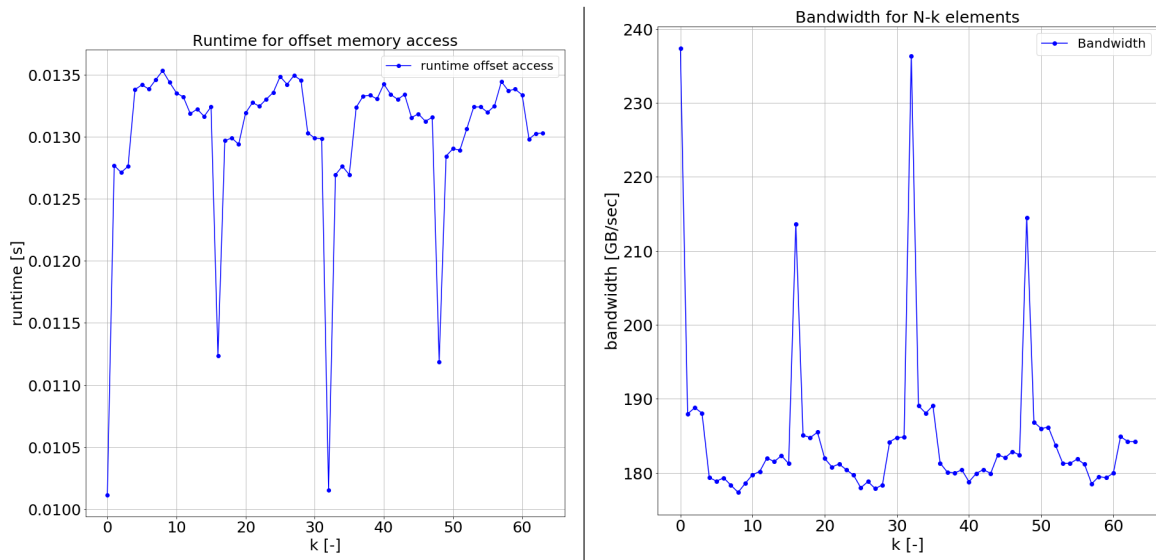
Figure 2: Runtime and bandwidth for offset memory access

Listing 3: Code for offset access

```
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          std::vector<double> timings;
50
51          for(int reps=0; reps < MAGIC_NUMBER; ++reps)
52          {
53              timer.reset();
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 and corrected it slightly because I think there was a math typo in line 9. 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 implementations.

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[row]; jj < rowoffsets[row+1]; ++jj)
10       {
11           val += values[jj] * x[colindices[jj]];
12       }
13       y[row] = val;
14     }
15 }
```

# 4 Conjugate Gradient - further kernels

For the following two kernels, I could build upon code I had already written for previous exercises and tweaked it slightly.

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  }
```