VIENNA UNIVERSITY OF TECHNOLOGY

360.252 COMPUTATIONAL SCIENCE ON MANY CORE ARCHITECTURES

INSTITUTE FOR MICROELECTRONICS

# Exercise 4

*Authors:*
Camilo TELLO FACHIN
12127084

*Supervisor:*
Dipl.-Ing. Dr.techn. Karl RUPP

November 15, 2022

TECHNISCHE
UNIVERSITÄT
WIEN
Vienna University of Technology

# Abstract

Here documented the results of exercise 4.

# Contents

# 1   Dot Product with Warp Shuffles (3/4 Points)

## 1.1   Task 1a

In The listing below see the code part for the Cuda Kernel for Exercise 4.1.a

C++ Cuda Code for 1a Kernel

```
1   // result  = (sum, abssum, squares, zero)
2    __global__  void cuda_1a(int N, double *x, double *sum, double *abssum, double *squares, double *zeros)
3   {
4      __shared__  double shared_mem_sum[512];
5      __shared__  double shared_mem_abssum[512];
6      __shared__  double shared_mem_squares[512];
7      __shared__  double shared_mem_zeros[512];
8
9     double sum_thr = 0;
10    double abssum_thr = 0;
11    double squares_thr  = 0;
12    double zeros_thr  = 0;
13
14    for (int  i = blockIdx.x * blockDim.x + threadIdx.x;  i < N; i += blockDim.x * gridDim.x) {
15      sum_thr   += x[i];
16      abssum_thr += abs(x[i]);
17      squares_thr  += pow(x[i],2);
18      zeros_thr  += (double)x[i]==0;
19    }
20
21    shared_mem_sum[threadIdx.x] = sum_thr;
22    shared_mem_abssum[threadIdx.x] = abssum_thr;
23    shared_mem_squares[threadIdx.x]  = squares_thr;
24    shared_mem_zeros[threadIdx.x]  = zeros_thr;
25
26    for (int  k = blockDim.x / 2; k > 0; k /= 2) {
27      __syncthreads ();
28      if (threadIdx.x < k) {
29        shared_mem_sum[threadIdx.x]  += shared_mem_sum[threadIdx.x + k];
30        shared_mem_abssum[threadIdx.x]  += shared_mem_abssum[threadIdx.x + k];
31        shared_mem_squares[threadIdx.x]  += shared_mem_squares[threadIdx.x + k];
32        shared_mem_zeros[threadIdx.x]  += shared_mem_zeros[threadIdx.x + k];
33      }
34    }
35
36    if (threadIdx.x == 0) {
37      atomicAdd(sum, shared_mem_sum[0]);
38      atomicAdd(abssum, shared_mem_abssum[0]);
39      atomicAdd(squares, shared_mem_squares[0]);
40      atomicAdd(zeros, shared_mem_zeros[0]);
41
42    }
43 }
```

ipe | INSTITUTE FOR MICROELECTRONICS

## 1.2 Task 1b

In the b part of Exercise 1, the warpshuffles had to be made use of instead of the shared memory. See below the listing with the code snippet for the CUDA Kernel with warpshuffles.

C++ Cuda Code for 1b Kernel

```
1  __global__ void cuda_1b(int N, double *x, double *sum, double *abssum, double *squares, double *zeros)
2  {
3      double sum_thr = 0;
4      double abssum_thr = 0;
5      double squares_thr = 0;
6      double zeros_thr = 0;
7
8      for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < N; i += blockDim.x * gridDim.x) {
9          sum_thr  += x[i];
10         abssum_thr += abs(x[i]);
11         squares_thr += pow(x[i],2);
12         zeros_thr += (double)x[i]==0;
13     }
14
15     for (int i=warpSize/2; i>0; i=i/2){
16         sum_thr += __shfl_down_sync(−1, sum_thr, i);
17         abssum_thr += __shfl_down_sync(−1, abssum_thr, i);
18         squares_thr += __shfl_down_sync(−1, squares_thr, i);
19         zeros_thr += __shfl_down_sync(−1, zeros_thr, i);
20     }
21
22     if ((threadIdx.x &(warpSize−1))== 0) {
23     atomicAdd(sum, sum_thr);
24     atomicAdd(abssum, abssum_thr);
25     atomicAdd(squares, squares_thr);
26     atomicAdd(zeros, zeros_thr);
27     }
28  }
```

## 1.3 Task 1c

In exercise 4.1.b, the same had to be done as before, but with warp shuffles only with one thread per entry. This can be done by changing the arguments of the Kernel invocation.

C++ Cuda Code for 1c Kernel Invocation

```
1      cuda_1c<<<((N+255)/256), 256>>>(N, cuda_x, cuda_sum, cuda_abssum, cuda_squares, cuda_zeros);
```

## 2  Multiple Dot Products (4 Points total)

### 2.1  Task 2a

In this Section it was asked specifically to "write a Kernel", so that is what you will find below, just a Kernel without any useful implementation!

C++ Cuda Code Kernel for 2a

```
1   __global__  void cuda_2a(int N, double *x, double *y, double *result)
2   {
3     __shared__ double shared_m1[256];
4     __shared__ double shared_m2[256];
5     __shared__ double shared_m3[256];
6     __shared__ double shared_m4[256];
7     __shared__ double shared_m5[256];
8     __shared__ double shared_m6[256];
9     __shared__ double shared_m7[256];
10    __shared__ double shared_m8[256];
11
12    double dot1=0, dot2=0, dot3=0, dot4=0, dot5=0, dot6=0, dot7=0, dot8=0;
13    for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < N; i += blockDim.x * gridDim.x) {
14      dot1 += x[i] * y[i];
15      dot2 += x[i] * y[i+N];
16      dot3 += x[i] * y[i+2*N];
17      dot4 += x[i] * y[i+3*N];
18      dot5 += x[i] * y[i+4*N];
19      dot6 += x[i] * y[i+5*N];
20      dot7 += x[i] * y[i+6*N];
21      dot8 += x[i] * y[i+7*N];
22    }
23
24    shared_m1[threadIdx.x] = dot1;
25    shared_m2[threadIdx.x] = dot2;
26    shared_m3[threadIdx.x] = dot3;
27    shared_m4[threadIdx.x] = dot4;
28    shared_m5[threadIdx.x] = dot5;
29    shared_m6[threadIdx.x] = dot6;
30    shared_m7[threadIdx.x] = dot7;
31    shared_m8[threadIdx.x] = dot8;
32
33    for (int k = blockDim.x / 2; k > 0; k /= 2) {
34      __syncthreads();
35      if (threadIdx.x < k) {
36        shared_m1[threadIdx.x] += shared_m1[threadIdx.x + k];
37        shared_m2[threadIdx.x] += shared_m2[threadIdx.x + k];
38        shared_m3[threadIdx.x] += shared_m3[threadIdx.x + k];
39        shared_m4[threadIdx.x] += shared_m4[threadIdx.x + k];
40        shared_m5[threadIdx.x] += shared_m5[threadIdx.x + k];
41        shared_m6[threadIdx.x] += shared_m6[threadIdx.x + k];
42        shared_m7[threadIdx.x] += shared_m7[threadIdx.x + k];
43        shared_m8[threadIdx.x] += shared_m8[threadIdx.x + k];
```

```
44        }
45      }
46
47      if (threadIdx.x == 0){
48        atomicAdd(&result[0], shared_m1[0]);
49        atomicAdd(&result[1], shared_m2[0]);
50        atomicAdd(&result[2], shared_m3[0]);
51        atomicAdd(&result[3], shared_m4[0]);
52        atomicAdd(&result[4], shared_m5[0]);
53        atomicAdd(&result[5], shared_m6[0]);
54        atomicAdd(&result[6], shared_m7[0]);
55        atomicAdd(&result[7], shared_m8[0]);
56      }
57    }
```

## 2.2   Task 2b

"Add a loop around that kernel to call it K/8 times to correctly compute the K dot products" This was
achieved by adjusting the algorithm in the following way:

C++ Cuda Code adaptation 2b

```
1   void Loop_function(int N, int K, double *x, double **y, double **results)
2   {
3       int h= K/8;
4       for (; h>0; h−−)
5       {
6           cuda_task2<<<256, 256>>>(N, x, y[h−1], results[h−1]);
7       }
8   }
9
10  float avg_vec(std::vector<float> timings_vec)
11  {
12      float avg = std::accumulate(timings_vec.begin(), timings_vec.end(), 0.0);
13      avg /= timings_vec.size();
14      return avg;
15  }
16
17  int main(void)
18  {
19
20    Timer timer;
21
22    for (int N = 1000; N < 1000001; N *= 10)
23    {
24        for (size_t K = 8; K < 33; K += 8)
25        {
26
27        std::vector<float> timings;
28
29        //
30        // allocate host memory:
31        //
32
33        double *x = (double*)malloc(sizeof(double) * N);
34        double *results   = (double*)malloc(sizeof(double) * K);
35        double **results3 = (double**)malloc(sizeof(double) * K/8);
36        for (size_t i=0; i<K/8; ++i) {
37           results3[i] = (double*)malloc(sizeof(double) * 8);
38        }
39
40        double **v = (double**)malloc(sizeof(double*) * K/8);
41        for (size_t i=0; i<K/8; ++i) {
42          v[i] = (double*)malloc(sizeof(double) * N*8);
43        }
44
45        //
```

```
46        // allocate  device  memory
47        //
48
49        double *cuda_x; cudaMalloc((&cuda_x), sizeof(double)*N);
50        double **cuda_v = (double**)malloc(sizeof(double*) * K/8);
51        for ( size_t  i=0; i<K/8; ++i) {
52          cudaMalloc( (void **)(&cuda_v[i]),  sizeof(double)*N*8);
53        }
54
55        double **cuda_results3 = (double**)malloc(sizeof(double*) * K/8);
56        for ( size_t  i=0; i<K/8; ++i) {
57          cudaMalloc( (void **)(&cuda_results3[i]),  sizeof(double)*8);
58        }
59
60        //
61        //  initialize  v
62        //
63
64        std :: fill (x, x + N, 1.0);
65        for ( size_t  i=0; i<K/8; ++i) {
66          for ( size_t  j=0; j<N*8; ++j) {
67            v[i][j] = 1 + rand() / (1.1 * RAND_MAX);
68          }
69        }
70
71        //
72        // Reference  calculation  on CPU:
73        //
74
75        for ( size_t  i=0; i<K; ++i) {results[i]=0;}
76        // dot product
77        for ( size_t  i=0; i<K/8; ++i) {
78          for ( size_t  j=0; j<N*8; ++j) {
79            results [i*8+j/N] += x[j%N] * v[i][j];
80          }
81        }
82
83        //
84        // Copy data to GPU
85        //
86
87        cudaMemcpy(cuda_x, x, sizeof(double)*N, cudaMemcpyHostToDevice);
88        for ( size_t  i=0; i<K/8; ++i) {
89          cudaMemcpy(cuda_v[i], v[i],  sizeof(double)*N*8, cudaMemcpyHostToDevice);
90        }
91
92        for (int  i = 0; i  < 11; i++)
93        {
94          //
95          // CUDA implementation
96          //
```

```
97
98          CUDA_ERRCHK(cudaDeviceSynchronize());
99          timer.reset();
100         Loop_function(N, K, cuda_x, cuda_v, cuda_results3);
101         CUDA_ERRCHK(cudaDeviceSynchronize());
102         float elapsed_time_0 = timer.get();
103         if (i > 0)
104         {
105             timings.push_back(elapsed_time_0);
106         }
107
108     }
109
110     std::cout << N << "," << K << "," << avg_vec(timings) << std::endl;
111
112     //
113     // Copy data to Host
114     //
115
116     for ( size_t i=0; i<K/8; ++i) {
117       cudaMemcpy(results3[i], cuda_results3[i], sizeof(double)*8, cudaMemcpyDeviceToHost);
118     }
```

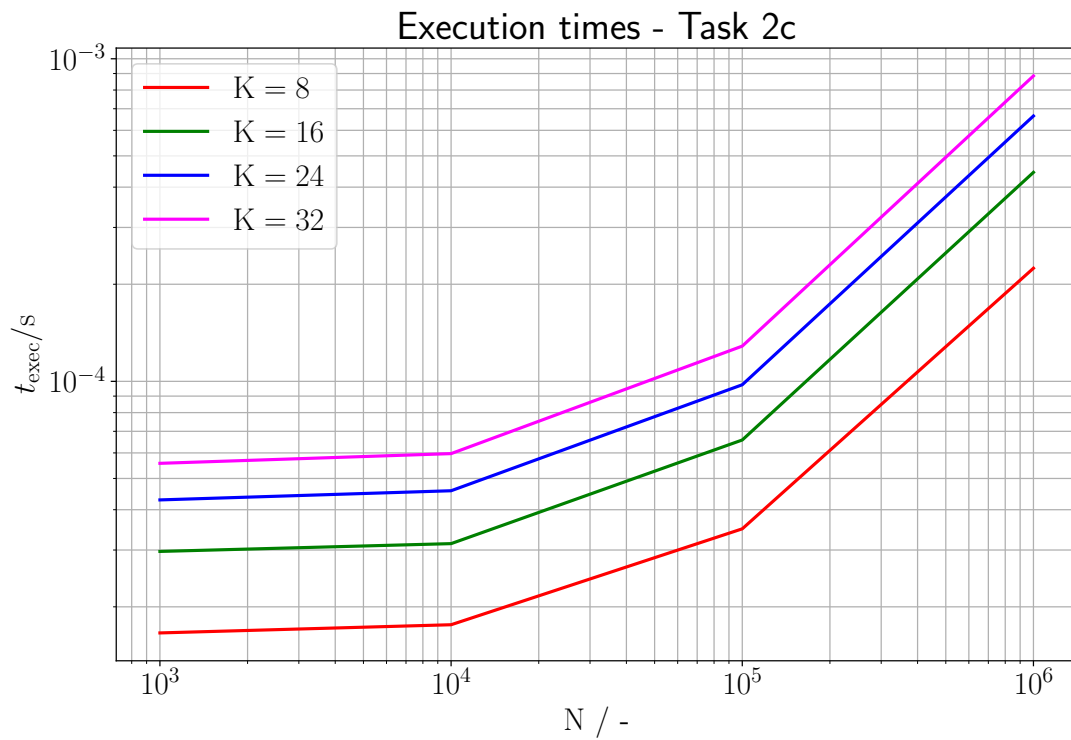## 2.3  Task 2c

In figure 1 below, the plot for task 2c.



Figure 1: Execution Times of concurrent dot products for different K and N

## 2.4  Task 2d

By checking if k mod(8) is 0, if not find a scheme to additionally start another kernel for a single dot product to balance out, otherwise proceed as before.