

Computational Science on Many-Core Architectures Exercise 4

Example 1 Multiple Dot Products (4 Points total)

a)

First I tried to implement this with a given y -matrix, but it does not work that well and I got a SEGMENTATION error.

Listing 1: kernel for 1a)

```

1  #include <stdio.h>
2  # include "timer.hpp"
3  #include <vector>
4
5  __global__ void dot_pro(int N, double *x, double *y0, double *y1, double *y2,
6  double *y3, double *y4, double *y5, double *y6, double *y7, double *dot)
7  {
8
9      unsigned int ind = threadIdx.x + blockDim.x*blockIdx.x;
10     unsigned int str = blockDim.x*gridDim.x;
11
12     __shared__ double cache0[256];
13     __shared__ double cache1[256];
14     __shared__ double cache2[256];
15     __shared__ double cache3[256];
16     __shared__ double cache4[256];
17     __shared__ double cache5[256];
18     __shared__ double cache6[256];
19     __shared__ double cache7[256];
20
21     double tmpsum0 = 0.0;
22     double tmpsum1 = 0.0;
23     double tmpsum2 = 0.0;
24     double tmpsum3 = 0.0;
25     double tmpsum4 = 0.0;
26     double tmpsum5 = 0.0;
27     double tmpsum6 = 0.0;
28     double tmpsum7 = 0.0;
29
30     double val = x[0];
31
32     while(ind < N)
33     {
34         for (int i = blockDim.x * blockIdx.x + threadIdx.x; i < N; i += blockDim.x *
35         {
36             val = x[i];
37             tmpsum0 += val * y0[i];
38             tmpsum1 += val * y1[i];
39             tmpsum2 += val * y2[i];

```

```
40         tmpsum3 += val * y3[i];
41         tmpsum4 += val * y4[i];
42         tmpsum5 += val * y5[i];
43         tmpsum6 += val * y6[i];
44         tmpsum7 += val * y7[i];
45     }
46
47     ind += str;
48 }
49
50     cache0[threadIdx.x] = tmpsum0;
51     cache1[threadIdx.x] = tmpsum1;
52     cache2[threadIdx.x] = tmpsum2;
53     cache3[threadIdx.x] = tmpsum3;
54     cache4[threadIdx.x] = tmpsum4;
55     cache5[threadIdx.x] = tmpsum5;
56     cache6[threadIdx.x] = tmpsum6;
57     cache7[threadIdx.x] = tmpsum7;
58
59     __syncthreads();
60
61     for(int i = blockDim.x/2; i>0; i/=2)
62     {
63         __syncthreads();
64         if(threadIdx.x < i)
65         {
66             cache0[threadIdx.x] += cache0[threadIdx.x + i];
67             cache1[threadIdx.x] += cache1[threadIdx.x + i];
68             cache2[threadIdx.x] += cache2[threadIdx.x + i];
69             cache3[threadIdx.x] += cache3[threadIdx.x + i];
70             cache4[threadIdx.x] += cache4[threadIdx.x + i];
71             cache5[threadIdx.x] += cache5[threadIdx.x + i];
72             cache6[threadIdx.x] += cache6[threadIdx.x + i];
73             cache7[threadIdx.x] += cache7[threadIdx.x + i];
74         }
75     }
76
77     if(threadIdx.x == 0)
78     {
79         atomicAdd(dot + 0, cache0[0]);
80         atomicAdd(dot + 1, cache1[0]);
81         atomicAdd(dot + 2, cache2[0]);
82         atomicAdd(dot + 3, cache3[0]);
83         atomicAdd(dot + 4, cache4[0]);
84         atomicAdd(dot + 5, cache5[0]);
85         atomicAdd(dot + 6, cache6[0]);
86         atomicAdd(dot + 7, cache7[0]);
87     }
88 }
```

b)

Listing 2: kernel for 1a)

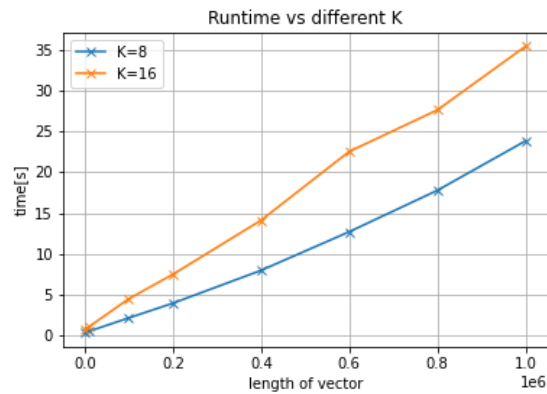
```

1  #include <stdio.h>
2  #include "timer.hpp"
3  #include <vector>
4  for (int g = 0; g < anz; g++)
5  {
6      for (int i=0; i<K/8; ++i)
7  Bandwidth_offset {
8      cudaDeviceSynchronize();
9      cudaMemcpy(d_y0, y[i*8+0], sizeof(double)*N, cudaMemcpyHostToDevice);
10     cudaMemcpy(d_y1, y[i*8+1], sizeof(double)*N, cudaMemcpyHostToDevice);
11     cudaMemcpy(d_y2, y[i*8+2], sizeof(double)*N, cudaMemcpyHostToDevice);
12     cudaMemcpy(d_y3, y[i*8+3], sizeof(double)*N, cudaMemcpyHostToDevice);
13     cudaMemcpy(d_y4, y[i*8+4], sizeof(double)*N, cudaMemcpyHostToDevice);
14     cudaMemcpy(d_y5, y[i*8+5], sizeof(double)*N, cudaMemcpyHostToDevice);
15     cudaMemcpy(d_y6, y[i*8+6], sizeof(double)*N, cudaMemcpyHostToDevice);
16     cudaMemcpy(d_y7, y[i*8+7], sizeof(double)*N, cudaMemcpyHostToDevice);
17     dot_pro<<<s, s>>>(N, d_x, d_y0, d_y1, d_y2, d_y3, d_y4
18     , d_y5, d_y6, d_y7, d_res_cblas);
19     cudaMemcpy(resultslarge+i*8, d_res_cblas
20     , sizeof(double)*8, cudaMemcpyDeviceToHost);
21     for (int j = 0; j < 8; j++)
22     {
23         res_cblas[j] = 0;
24     }
25     cudaMemcpy(d_res_cblas, res_cblas, sizeof(double)*8, cudaMemcpyHostToDevice)
26     }
27 }
28 printf("Dot product took %g seconds", 1000*timer.get()/anz);

```

c)

I got a error message that I did not understand and tht's why I could not ran the programm with the other K -values.



Inconsistency detected by ld.so: dl-fini.c: 87: _dl_fini: Assertion `ns != LM_ID_BASE || i == nloaded' failed!

d)

On approach could be to store the an arbitrary number of vectors in a matirx and send it to the kernel and then in kernel calculate the dot-product with a dynamic and static arrays.

Pipelined CG (4 Points total)

Different Kernels

Listing 3: kernel line 7 to 9)

```

1  __global__ void x_plus_a_p(int N, double *x, double *p, double *r, double *Ap,
2  double alpha, double beta)
3  {
4      for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < N;
5          i += blockDim.x * gridDim.x)
6      {
7          x[i] = x[i] + alpha * p[i];
8          r[i] = r[i] - alpha * Ap[i];
9          p[i] = r[i] + beta * p[i];
10     }
11 }
```

Listing 4: kernel line 10 to 12)

```

1  __global__ void diff_dot_prod(int N, double *Ap, double *p, double *r,
2  double *ApAp, double *pAp, double *rr, int *csr_rowoffsets,
3  int *csr_colindices, double *csr_values)
4  {
5      __shared__ double cache0[512];
6      __shared__ double cache1[512];
7      __shared__ double cache2[512];
8
9      double tmpsum0 = 0;
10     double tmpsum1 = 0;
11     double tmpsum2 = 0;
12
13     for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < N; i += blockDim.x *
14         {
15         double sum = 0;
16         for (int k = csr_rowoffsets[i]; k < csr_rowoffsets[i + 1]; k++)
17         {
18             sum += csr_values[k] * p[csr_colindices[k]];
19         }
20         Ap[i] = sum;
21
22         tmpsum0 += Ap[i] * Ap[i];
23         tmpsum1 += p[i] * Ap[i];
24         tmpsum2 += r[i] * r[i];
25     }
26
27     cache0[threadIdx.x] = tmpsum0;
28     cache1[threadIdx.x] = tmpsum1;
29     cache2[threadIdx.x] = tmpsum2;
30
31
32
```

```

33     for (int k = blockDim.x / 2; k > 0; k /= 2)
34     {
35         __syncthreads();
36         if (threadIdx.x < k)
37         {
38             cache0[threadIdx.x] += cache0[threadIdx.x + k];
39             cache1[threadIdx.x] += cache1[threadIdx.x + k];
40             cache2[threadIdx.x] += cache2[threadIdx.x + k];
41         }
42     }
43
44     if (threadIdx.x == 0)
45     {
46         atomicAdd(ApAp, cache0[0]);
47         atomicAdd(pAp, cache1[0]);
48         atomicAdd(rr, cache2[0]);
49     }
50 }

```

Listing 5: major loop for iteration

```

1  while (1)
2  {
3      x_plus_a_p<<<512,512>>>(N, cuda_solution, cuda_p, cuda_r, cuda_Ap,
4      alpha, beta);
5
6
7      cudaMemcpy(cuda_ApAp, &zero, sizeof(double), cudaMemcpyHostToDevice);
8      cudaMemcpy(cuda_pAp, &zero, sizeof(double), cudaMemcpyHostToDevice);
9      cudaMemcpy(cuda_rr, &zero, sizeof(double), cudaMemcpyHostToDevice);
10     diff_dot_prod<<<512,512>>>(N, cuda_Ap, cuda_p, cuda_r, cuda_ApAp,
11     cuda_pAp, cuda_rr, csr_rowoffsets, csr_colindices, csr_values);
12
13     cudaMemcpy(&beta, cuda_ApAp, sizeof(double), cudaMemcpyDeviceToHost);
14     cudaMemcpy(&alpha, cuda_pAp, sizeof(double), cudaMemcpyDeviceToHost);
15     cudaMemcpy(&residual_norm_squared, cuda_rr, sizeof(double),
16     cudaMemcpyDeviceToHost);
17
18
19     // line convergence check:
20     if (std::sqrt(residual_norm_squared / initial_residual_squared) < 1e-6)
21     {
22         break;
23     }
24
25     // line 13:
26     alpha = residual_norm_squared / alpha;
27
28     // line 14:
29     beta = alpha * alpha * beta / residual_norm_squared - 1;
30

```

```

31     if (iters > 10000)
32     break; // solver didn't converge
33     ++iters;
34 }

```

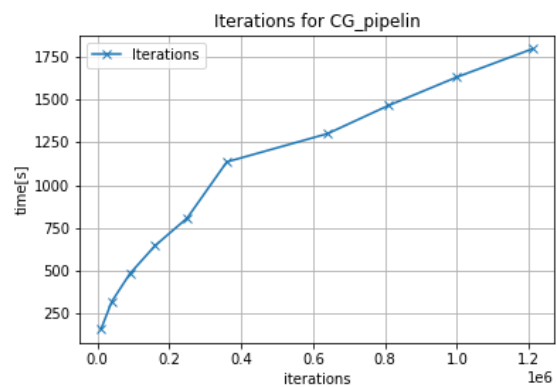
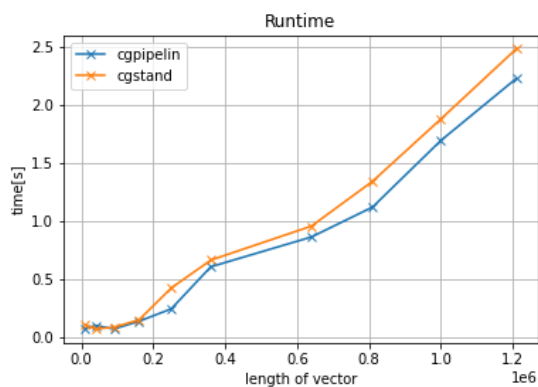
Listing 6: main

```

1  for(int i = 100; i < 1200; i = i + 100)
2  {
3      solve_system(i);
4      std::cout << "i:" << i << std::endl;
5      std::cout << "" << std::endl;
6      std::cout << "" << std::endl;
7  }

```

Benchmark



I expected that the pipeline CG should perform better and also that the number of iterations should be reduced but the number of iterations are exactly the same. Maybe I have to work on the second kernel in the iteration.