

VIENNA UNIVERSITY OF TECHNOLOGY

360.252 COMPUTATIONAL SCIENCE ON MANY CORE ARCHITECTURES

INSTITUTE FOR MICROELECTRONICS

Exercise 7

Authors:

Camilo TELLO FACHIN
12127084

Supervisor:

Dipl.-Ing. Dr.techn. Karl RUPP

December 11, 2022



TECHNISCHE
UNIVERSITÄT
WIEN

Vienna University of Technology

Abstract

Here documented the results of Exercise 7.

Contents

1	Pipelined Conjugate Gradients (5/5 Points)	1
1.1	Implement the Algorithm (3/3 Points)	1
1.1.1	Blue Font Algorithm Part	1
1.1.2	Red Font Algorithm Part	2
1.1.3	Bringing Everything together	3
1.2	Comparison of Pipelined CG and Classical CG on both GPU's (2/2 Points) . . .	4
1.2.1	CG Algorithm Performance on RTX3060	4
1.2.2	CG Algorithm Performance on Tesla K40	5

1 Pipelined Conjugate Gradients (5/5 Points)

1.1 Implement the Algorithm (3/3 Points)

In the listings below the Kernels for the blue and red colored font from lecture slides 7a. The task was to only execute 2 Kernel call's per iteration, with these two Kernels, this can be achieved.

1.1.1 Blue Font Algorithm Part

Kernel for Blue Font Algorithm Part

```

1  __global__ void cuda_blue(int N, double *x, double *p, double *Ap, double *r, double *r_ip, double Alpha,
2  double Beta)
3  {
4  __shared__ double shared_memory[512];
5  double partial_dot_product = 0;
6
7  for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < N; i += blockDim.x * gridDim.x){
8      double p_thread = p[i];
9      double Ap_thread = Ap[i];
10     double r_thread = r[i] - Alpha * Ap_thread;
11     x[i] += Alpha * p_thread;
12     r[i] = r_thread;
13     p[i] = r_thread + Beta * p_thread;
14     partial_dot_product += r_thread * r_thread;
15 }
16 shared_memory[threadIdx.x] = partial_dot_product ;
17 for (int j = blockDim.x / 2; j > 0; j /= 2) {
18     __syncthreads();
19     if (threadIdx.x < j) {
20         shared_memory[threadIdx.x] += shared_memory[threadIdx.x + j];
21     }
22     if (threadIdx.x == 0) r_ip[blockIdx.x] = shared_memory[0];
23 }
24 }
```

1.1.2 Red Font Algorithm Part

Kernels for Red Algorithm Part

```

1  __global__ void cuda_blue(int N, double *x, double *p, double *Ap, double *r, double *r_ip, double Alpha,
2      double Beta)
3  {
4      __shared__ double shared_memory[512];
5      double partial_dot_product = 0;
6      for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < N; i += blockDim.x * gridDim.x){
7          double p_thread = p[i];
8          double Ap_thread = Ap[i];
9          double r_thread = r[i] - Alpha * Ap_thread;
10         x[i] += Alpha * p_thread;
11         r[i] = r_thread;
12         p[i] = r_thread + Beta * p_thread;
13         partial_dot_product += r_thread * r_thread;
14     }
15     shared_memory[threadIdx.x] = partial_dot_product;
16     for (int j = blockDim.x / 2; j > 0; j /= 2) {
17         __syncthreads();
18         if (threadIdx.x < j) {
19             shared_memory[threadIdx.x] += shared_memory[threadIdx.x + j];
20         }
21     }
22     if (threadIdx.x == 0) r_ip[blockIdx.x] = shared_memory[0];
23 }

```

1.1.3 Bringing Everything together

In the function `conjugate_gradient_pipelined()` we use both the Kernels from before and iterate till convergence. In order for this algorithm to work, it needs initial α_0 , β_0 and also an A_{p0} , this is done before the actual iterations in the `while()` loop. One can use the old functions from the last exercise like the dot product and the matrix vector product to obtain the two values and the vector. In the listing below shown the `while()` loop that iterates till convergence.

CG while() loop

```

1  int  iters  = 0;
2  cudaDeviceSynchronize();
3  timer.reset();
4  while (1) {
5      // Line 2–4 and partial of line 6, The blue colored part in the algorithm:
6      cuda_blue<<<blocks_Inch,thrds_block>>>(N, cuda_solution, cuda_p, cuda_Ap, cuda_r, cuda_ip_rr, alpha,
          beta);
7      // cudaMemcpy for blockwise inner product <r,i, r_i> for CPU calculations of alpha and beta
8      cudaMemcpy(bwip_rr, cuda_ip_rr, sizeof(double) * blocks_Inch, cudaMemcpyDeviceToHost);
9      // Line 5 and 6, the red colored part in the algorithm
10     cuda_red<<<blocks_Inch,thrds_block>>>(N, csr_rowoffsets, csr_colindices, csr_values, cuda_p, cuda_Ap,
        cuda_ip_ApAp, cuda_ip_pAp);
11     // cudaMemcpy for blockwise inner product <Ap,i, Ap,i> and <p,i, Ap,i> for CPU calculations of alpha
        and beta
12     cudaMemcpy(bwip_pAp, cuda_ip_pAp, sizeof(double) * blocks_Inch, cudaMemcpyDeviceToHost);
13     cudaMemcpy(bwip_ApAp, cuda_ip_ApAp, sizeof(double) * blocks_Inch, cudaMemcpyDeviceToHost);
14     // CPU summation of blockwise inner products.
15     ip_rr = bwip_rr[0];
16     ip_ApAp = bwip_ApAp[0];
17     ip_pAp = bwip_pAp[0];
18     for (size_t i=1; i<blocks_Inch; ++i)
19     {
20         ip_rr += bwip_rr[i];
21         ip_ApAp += bwip_ApAp[i];
22         ip_pAp += bwip_pAp[i];
23     }
24     //Check if convergence criterion is fulfilled .
25     if (std::sqrt( ip_rr / initial_residual_squared ) < 1e-6) {
26         break;
27     }
28     // Computation of alpha and beta for next while iteration .
29     alpha = ip_rr / ip_pAp;
30     beta = ( alpha*alpha*ip_ApAp - ip_rr ) / ip_rr ;
31
32     if ( iters > 10000)
33         break; // solver didn't converge
34     ++iters;
35 }

```

1.2 Comparison of Pipelined CG and Classical CG on both GPU's (2/2 Points)

1.2.1 CG Algorithm Performance on RTX3060

For this specific implementation, the runtime for the pipelined CG algorithm on the RTX3060 outperforms the classic one only for \sqrt{N} larger than approximately 300. After this point, the improvements of the pipelined algorithm appear to be only minor.

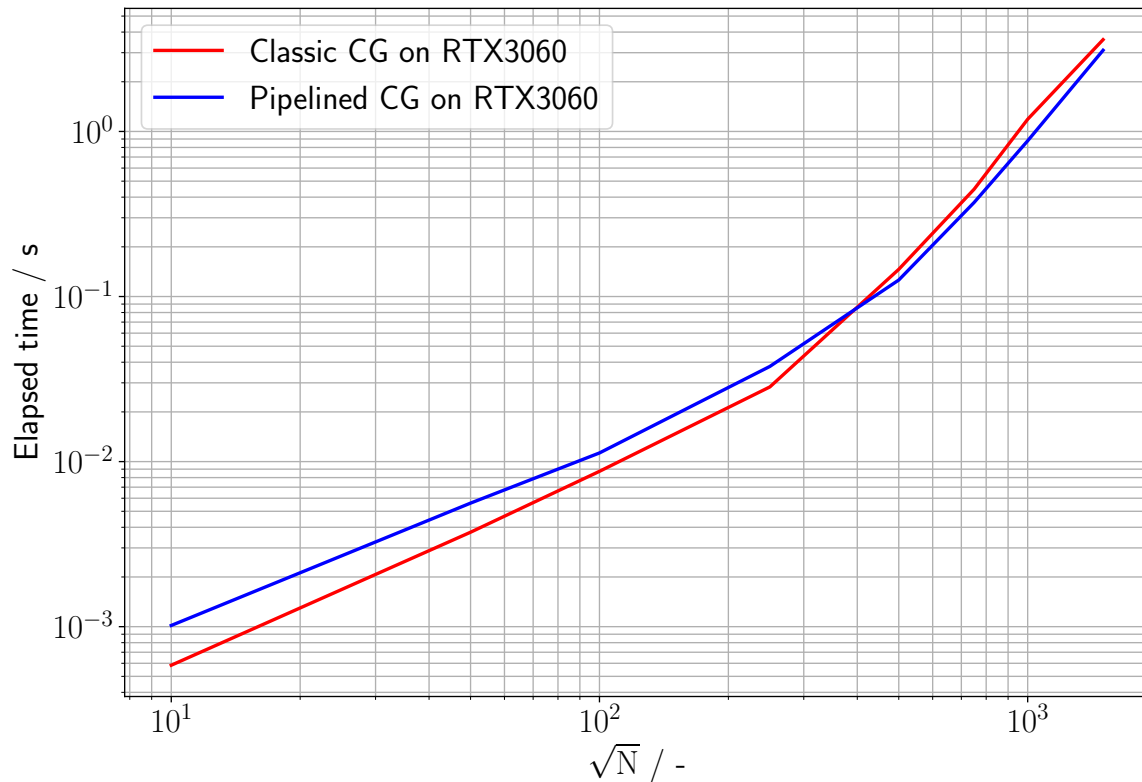


Figure 1: Elapsed time for Classical and Pipelined CG Algorithm on the RTX3060 GPU

1.2.2 CG Algorithm Performance on Tesla K40

For this specific implementation, the runtime for the pipelined CG algorithm on the Tesla K40 outperforms the classic one only for \sqrt{N} larger than approximately 350. After this point, the improvements of the pipelined algorithm appear to be only minor.

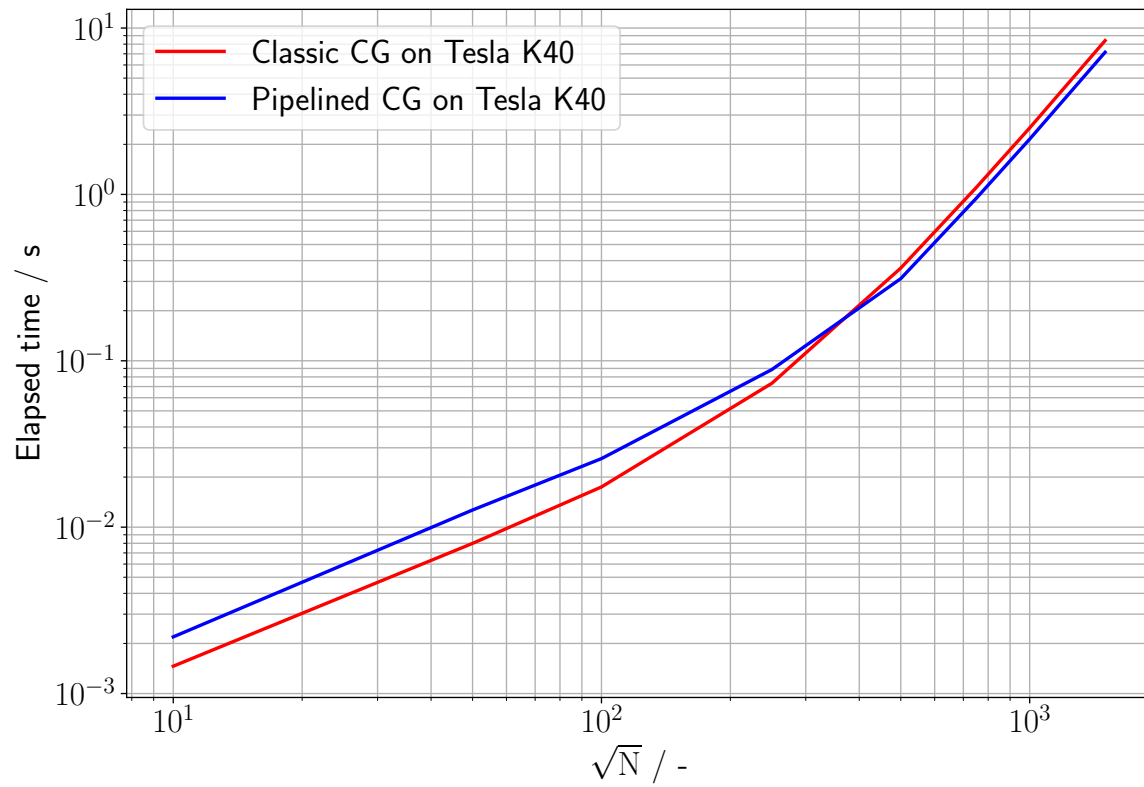


Figure 2: Elapsed time for Classical and Pipelined CG Algorithm on the Tesla K40 GPU