

Programming Massive-Parallel Processor - Exercise 1



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Dominik Marino - 2468378
21. November 2020

Inhaltsverzeichnis

1	Cuda Devices	2
2	CPU Multiplication	2
2.1	Allocation Matrices	2
2.2	Fill Matrices	2
2.3	Matrix CPU Multiplication	2
2.4	Testing CPU	3
3	GPU Multiplication	3
3.1	Set CUDA Device	3
3.2	Allocate GPU Memory	3
3.3	Upload 2 GPU	4
3.4	Matrix GPU Multiplication	4
3.5	Download 2 GPU	5
3.6	Testing GPU	5
3.7	Compare CPU Matrix with GPU Matrix	5
3.8	Print computation time from CPU and GPU	5
4	utils	6
4.1	Print Matrix	6
5	Question	6

1 Cuda Devices

```
1 void print_cuda_devices() {
2     int nDevice;
3     cudaGetDeviceCount(&nDevice);
4
5     for (int i = 0; i < nDevice; i++) {
6         cudaDeviceProp prop;
7         cudaGetDeviceProperties(&prop, i);
8
9         std::cout << "Compute capability (Major/Minor): \t" << prop.major << ", " << prop.minor << std::endl;
10        std::cout << "Multiprocessor count: \t\t\t" << prop.multiProcessorCount << std::endl; //return in int
11        std::cout << "GPU clock rate: \t\t\t" << prop.clockRate << " kHz (" << prop.clockRate / 1000000.0 << "GHz)" << std::endl;
12        std::cout << "Total global memory: \t\t\t" << prop.totalGlobalMem << " bytes (" << prop.totalGlobalMem / 1024.0 << " MiB)" << std::endl; //return in bytes
13        std::cout << "L2 chache size: \t\t\t" << prop.l2CacheSize << " bytes (" << prop.l2CacheSize / 1048576.0 << " KiB)" << std::endl; //return in bytes
14        std::cout << "\n" << std::endl;
15    }
16 }
```

Listing 1: Print Cuda Devices (matmul.cc)

2 CPU Multiplication

2.1 Allocation Matrices

```
1 // TODO: Allocate CPU matrices (see matrix.cc)
2 // Matrix sizes:
3 // Input matrices:
4 // Matrix M: pmpp::M_WIDTH, pmpp::M_HEIGHT
5 // Matrix N: pmpp::N_WIDTH, pmpp::N_HEIGHT
6 // Output matrices:
7 // Matrix P: pmpp::P_WIDTH, pmpp::P_HEIGHT
8
9 CPUMatrix inMatrixCPU1 = matrix_alloc_cpu(pmpp::M_HEIGHT, pmpp::M_WIDTH);
10 CPUMatrix inMatrixCPU2 = matrix_alloc_cpu(pmpp::N_WIDTH, pmpp::N_WIDTH);
11
12 CPUMatrix outMatrixCPU = matrix_alloc_cpu(pmpp::P_WIDTH, pmpp::P_HEIGHT);
```

Listing 2: Allocate CPU matrices (matmul.cc)

2.2 Fill Matrices

```
1 // TODO: Fill the CPU input matrices with the provided test values (pmpp::fill(CPUMatrix &m, CPUMatrix &n))
2 pmpp::fill(inMatrixCPU1, inMatrixCPU2);
```

Listing 3: using fill method for CPU input matrices (matmul.cc)

2.3 Matrix CPU Multiplication

```
1 // TODO (Task 5): Start CPU timing here!
2 timer_tp start_cpu = timer_now();
3
4 // TODO: Run your implementation on the CPU (see mul_cpu.cc)
5 matrix_mul_cpu(inMatrixCPU1, inMatrixCPU2, outMatrixCPU);
6
7 // TODO (Task 5): Stop CPU timing here!
8 timer_tp stop_cpu = timer_now();
9 float cpu_comp_time = timer_elapsed(start_cpu, stop_cpu);
```

Listing 4: Start and Stop timer and run the implementation (matmul.cc)

```

1 void matrix_mul_cpu(const CPUMatrix &m, const CPUMatrix &n, CPUMatrix &p) {
2
3     for (int i = 0; i < m.height; i++) {
4         for (int j = 0; j < n.width; j++) {
5             float sum = 0.0;
6             for (int k = 0; k < n.height; k++) {
7                 sum = sum + (m.elements[i * m.width + k] * n.elements[k * n.width + j]);
8             }
9             p.elements[i * p.width + j] = sum;
10        }
11    }
12 }

```

Listing 5: Matrix multiplication (mul_cpu.cc)

2.4 Testing CPU

```

1 // TODO: Check your matrix for correctness (pmpp::test_cpu(const CPUMatrix &p))
2 pmpp::test_cpu(outMatrixCPU);

```

Listing 6: Test Output Matrix (matmul.cc)

3 GPU Multiplication

3.1 Set CUDA Device

```

1 // TODO: Set CUDA device
2 int nDevice = 0;
3 cudaGetDeviceCount(&nDevice);
4 int userDeviceInput = nDevice - 1;
5
6 if (userDeviceInput < nDevice) {
7     cudaSetDevice(userDeviceInput);
8 } else {
9     printf("error: invalid device choosen\n");
10 }

```

Listing 7: Set Cuda Device and check if the choosen device is invalid (matmul.cc)

3.2 Allocate GPU Memory

```

1 // TODO: Allocate GPU matrices (see matrix.cc)
2 GPUMatrix inMatrixGPU1 = matrix_alloc_gpu(pmpp::M_HEIGHT, pmpp::M_WIDTH);
3 GPUMatrix inMatrixGPU2 = matrix_alloc_gpu(pmpp::N_WIDTH, pmpp::N_WIDTH);
4
5 GPUMatrix outMatrixGPU = matrix_alloc_gpu(pmpp::P_WIDTH, pmpp::P_HEIGHT);

```

Listing 8: Allocate GPU Matrices (matmul.cc)

```

1 GPUMatrix matrix_alloc_gpu(int width, int height)
2 {
3     GPUMatrix m;
4     m.width = width;
5     m.height = height;
6
7     cudaMallocPitch((void**)&m.elements, &m.pitch, width * sizeof(float), m.height);
8
9     return m;
10 }

```

Listing 9: Matrix alloc GPU (matmul.cc)

3.3 Upload 2 GPU

```
1 // TODO: Upload the CPU input matrices to the GPU (see matrix.cc)
2 matrix_upload(inMatrixCPU1, inMatrixGPU1);
3 matrix_upload(inMatrixCPU2, inMatrixGPU2);
```

Listing 10: Upload Matrix 2 GPU (matmul.cc)

```
1 void matrix_upload(const CPUMatrix &src, GPUMatrix &dst)
2 {
3     int size = src.height * src.width;
4     cudaMemcpy(dst.elements, src.elements, size * sizeof(float), cudaMemcpyHostToDevice);
5 }
```

Listing 11: Upload Matrix (matrix.cc)

3.4 Matrix GPU Multiplication

```
1 // TODO (Task 5): Start GPU timing here!
2 // TODO (Task 5): Start GPU timing here!
3 cudaEvent_t evStart, evStop;
4 cudaEventCreate(&evStart);
5 cudaEventCreate(&evStop);
6 cudaEventRecord(evStart, 0);
7
8 // TODO: Run your implementation on the GPU (see mul_gpu.cu)
9 matrix_mul_gpu(inMatrixGPU1, inMatrixGPU2, outMatrixGPU);
10
11 // TODO (Task 5): Stop GPU timing here!
12 cudaEventRecord(evStop, 0);
13 cudaEventSynchronize(evStop);
14 float elapsedTime_ms;
15
16 cudaEventElapsedTime(&elapsedTime_ms, evStart, evStop);
17 printf("CUDA processing took: %f ms\n", elapsedTime_ms);
18 cudaEventDestroy(evStart);
19 cudaEventDestroy(evStop);
```

Listing 12: Matrix Multiplication with timer for calculation (matmul.cc)

```
1 // TODO (Task 4): Implement matrix multiplication CUDA kernel
2 __global__ void multiplication_kernel(GPUMatrix m, GPUMatrix n, GPUMatrix out) {
3     int tx = threadIdx.x;
4     int ty = threadIdx.y;
5     float pValue = 0;
6
7     for (int k = 0; k < m.width; k++) {
8         float Melement = m.elements[ty * m.width + k];
9         float Nelement = n.elements[k * n.width + tx];
10
11         pValue += Melement * Nelement;
12     }
13     int out_index = ty * out.width + tx;
14     out.elements[ty * out.width + tx] = pValue;
15 }
16
17 void matrix_mul_gpu(const GPUMatrix &m, const GPUMatrix &n, GPUMatrix &p)
18 {
19     // TODO (Task 4): Determine execution configuration and call CUDA kernel
20     dim3 grid(1, 1); //div_up(m.height, n.width);
21     dim3 dimBlock(m.width, n.height);
22
23     multiplication_kernel<<<grid, dimBlock >>>(m, n, p);
24
25     cudaError_t err = cudaGetLastError();
26
27     if (err != cudaSuccess)
28         printf("Error: %s\n", cudaGetErrorString(err));
29
30     cudaDeviceSynchronize();
31 }
```

Listing 13: Kernel and Set up (mul_gpu.cu)

3.5 Download 2 CPU

```
1 // TODO: Download the GPU output matrix to the CPU (see matrix.cc)
2 CPMatrix dOutputMatrixfromGPU = matrix_alloc_cpu(pmpp::P_WIDTH, pmpp::P_HEIGHT);
3 matrix_download(outMatrixGPU, dOutputMatrixfromGPU);
```

Listing 14: download Matrix (matmul.cc)

```
1 void matrix_download(const GPUMatrix &src, CPMatrix &dst)
2 {
3     int size = src.height * src.width;
4     cudaMemcpy(dst.elements, src.elements, size * sizeof(float), cudaMemcpyDeviceToHost);
5 }
```

Listing 15: Download Matrix from GPU (matrix.cc)

3.6 Testing GPU

```
1 // TODO: Check your downloaded matrix for correctness (pmpp::test_gpu(const CPMatrix &p))
2 pmpp::test_gpu(dOutputMatrixfromGPU);
```

Listing 16: Testing GPU Matrix (matmul.cc)

3.7 Compare CPU Matrix with GPU Matrix

```
1 // TODO: Compare CPU result with GPU result (see matrix.cc)
2 matrix_compare_cpu(outMatrixCPU, outputMatrixGPU);
```

Listing 17: Compare CPU Matrix with GPU Matrix (matmul.cc)

```
1 void matrix_compare_cpu(const CPMatrix &a, const CPMatrix &b)
2 {
3     int j = 0;
4     for (int i = 0; i < a.height * a.width; i++) {
5         if (a.elements[i] == b.elements[i]) {
6             j++;
7         }
8     }
9     if (j == (a.height * a.width)) {
10         std::cout << "CPU Matrix and GPU Matrix are equal" << std::endl;
11     } else {
12         std::cout << "Matrices are not equal" << std::endl;
13     }
14 }
```

Listing 18: Compare CPU Matrix with GPU Matrix (matmul.cc)

3.8 Print computation time from CPU and GPU

```
1 //print computation time from cpu and gpu
2 print_time(cpu_comp_time, elapsedTime_ms);
```

Listing 19: computation time from CPU and GPU (matmul.cc)

```
1 void print_time(float cpu_comp_time, float gpu_comp_time) {
2     std::cout << "\n" << std::endl;
3     std::cout << "Computing Time for CPU: " << cpu_comp_time << std::endl;
4     std::cout << "Computing Time for GPU: " << gpu_comp_time << std::endl;
5
6     if (cpu_comp_time < gpu_comp_time) {
7         std::cout << "CPU is faster" << std::endl;
8     } else if (gpu_comp_time < cpu_comp_time) {
9         std::cout << "GPU is faster" << std::endl;
10    } else {
11        std::cout << "Computing time is equal" << std::endl;
12    }
13 }
14 }
15 }
```

Listing 20: Compare computation time (timer.h)

4 utils

4.1 Print Matrix

```
1 void print_matrix(CPUMatrix &matrix) {  
2     int j = 0;  
3     for (int i = 0; i < matrix.height * matrix.width; ++i) {  
4         std::cout << matrix.elements[i] << ' ';  
5         j++;  
6         if (j == matrix.width) {  
7             std::cout << std::endl;  
8             j = 0;  
9         }  
10    }  
11 }
```

Listing 21: Print Matrix in console (matmul.cc)

5 Question

- Where do the differences come from?
 - The CPU calculation is slower than the GPU calculation. If we look up to the CPU implementation, then we see that the algorithm is $O(n^3)$. In the case of matrix multiplication we can parallelize the computations, because the GPU have much more threads. In our task has every thread to calculate a position/element in the output matrix P .
 - But it really depends on the size of the matrix, but this does not mean that paralleling of a matrix multiplication is always faster than the CPU implementation. It really depends on the size of the data. Small data with a few calculation fits poor on a GPU. Another reason can be memory access for slowness i.e. if we are copying a lot of memory in our algorithm. In my code i want to synchronise all the threads and if we have to wait until all threads are done, then it also can slow down the computation time.