

Домашнее задание №2 (CUDA)

Умножение матриц на GPU средствами CUDA.

Постановка задачи

$$C = A \cdot B$$

Элементы матриц A, B и C имеют тип double; матрицы имеют размер $N \times N$

Матрицы A, B и C хранятся в одномерных массивах в column-major порядке

Задания

1. Реализовать перемножение матриц с использованием глобальной памяти на CUDA.

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include "cuda_runtime.h"
4  #include "device_launch_parameters.h"
5  #define BLOCK_SIZE 2
6
7  // Filling matrix with random double numbers from 0 to 9
8  void fill_matrix(double* H, int N)
9  {
10     int i;
11     for (i = 0; i < N * N; ++i)
12     {
13         H[i] = (rand() % 9);
14     }
15 }
16
```

```

17 // Printing the matrix
18 void print_matrix(double* H, int N)
19 {
20     int i, j;
21     for (j = 0; j < N; ++j)
22     {
23         for (i = 0; i < N; ++i)
24         {
25             if (i != N - 1)
26             {
27                 printf("%f\t", H[i * N + j]);
28             }
29             else
30             {
31                 printf("%f\n", H[i * N + j]);
32             }
33         }
34     }

```

```

35     printf("\n");
36 }
37
38 // Sequential multiplication
39 void mult_pos(int N, double* A, double* B, double* C)
40 {
41     for (int i = 0; i < N; i++)
42     {
43         for (int j = 0; j < N; j++)
44         {
45             for (int k = 0; k < N; k++)
46             {
47                 C[j * N + i] += A[k * N + i] * B[j * N + k];
48             }
49         }
50     }
51 }

```

```

52
53 // Comparing matrices
54 void compare_matrices(double* C, double* hostC, int N)
55 {
56     int i, flag;
57     double eps;
58     flag = 0;
59     eps = 0.001;
60     for (i = 0; i < N * N; ++i)
61     {
62         if (abs(C[i] - hostC[i]) >= eps)
63         {
64             flag = 1;
65             break;
66         }
67     }

```

```

68     if (flag == 0)
69     {
70         printf("Results of sequential and parallel multiplications are equal \n");
71     }
72     else
73     {
74         printf("Results of sequential and parallel multiplications are different \n");
75     }
76 }
77
78 //CUDA multiplication

```

```

79 __global__ void matrixmultiplication(double* A, double* B, double* C, int N)
80 {
81     //Computing position of elemets
82     int c = blockIdx.x * blockDim.x + threadIdx.x;
83     int r = blockIdx.y * blockDim.y + threadIdx.y;
84     double sum = 0.0;
85     //Computing elements for matrix C
86     if (r < N && c < N)
87     {
88         // each thread computes one element of the block sub-matrix
89         for (int i = 0; i < N; i++)
90         {
91             sum += B[r * N + i] * A[i * N + c];
92         }
93     }
94     C[r * N + c] = sum;
95 }

```

```

96
97 int main(int argc, char* argv[])
98 {
99     int N = 4;
100     double *hostA;
101     double *hostB;
102     double *hostC;
103     double *C;
104     float time;
105
106     //Allocating memory for hosts
107     hostA = (double*) calloc(N*N, sizeof(double));
108     hostB = (double*) calloc(N*N, sizeof(double));
109     hostC = (double*) calloc(N*N, sizeof(double));
110
111     C = (double*) calloc(N*N, sizeof(double));
112
113     //Filling matrices
114     fill_matrix(hostA, N);
115     fill_matrix(hostB, N);
116
117     //Printg matrices A and B
118     printf("hostA:\n");
119     print_matrix(hostA, N);
120     printf("hostB:\n");
121     print_matrix(hostB, N);
122
123     //Allocating memory for devices
124     double *deviceA;
125     cudaMalloc((void **)&deviceA, N*N * sizeof(double));
126     double *deviceB;
127     cudaMalloc((void **)&deviceB, N*N * sizeof(double));
128     double * deviceC;
129     cudaMalloc((void **)&deviceC, N*N * sizeof(double));
130

```

```

131     //Defining dimensions of the block and the grid
132     dim3 threadsPerBlock = dim3(BLOCK_SIZE, BLOCK_SIZE);
133     int N_new = 0;
134     if (N % BLOCK_SIZE == 0)
135     |     N_new = int (N / BLOCK_SIZE);
136     else
137     |     N_new = int (N / BLOCK_SIZE) + 1;
138     dim3 blocksPerGrid = dim3(N_new,N_new);
139
140     //Defining and initializing variables for computing time
141     cudaEvent_t start, stop;
142     cudaEventCreate(&start);
143     cudaEventCreate(&stop);
144

```

```

145     //Recording the computing time
146     cudaEventRecord( start, 0);
147
148     //Copying matrices from host to device
149     cudaMemcpy(deviceA, hostA, N*N * sizeof(double), cudaMemcpyHostToDevice);
150     cudaMemcpy(deviceB, hostB, N*N * sizeof(double), cudaMemcpyHostToDevice);
151
152     //Calling multiolication function
153     matrixmultiplication<<<blocksPerGrid, threadsPerBlock>>>(deviceA, deviceB, deviceC, N);
154
155     //Copying matrices from device to host
156     cudaMemcpy(hostC, deviceC, N*N * sizeof(double), cudaMemcpyDeviceToHost);
157     cudaEventRecord( stop, 0);
158     cudaEventSynchronize(stop);
159

```

```

160     //Computing time
161     cudaEventElapsedTime( &time, start, stop);
162     printf("Time: %.2f \n",time);
163
164     //Sequantial multiplication
165     mult_pos(N, hostA, hostB, C);
166
167     //Printg result of CUDA multiplication
168     printf("hostC:\n");
169     print_matrix(hostC, N);
170
171     //Printg result of sequential multiplication
172     printf("C:\n");
173     print_matrix(C, N);
174
175     //Comparing results of CUDA multiplication and sequential multiplication
176     compare_matrices(C, hostC, N);
177

```

```

178 //Deallocationg the memory
179     cudaFree(deviceA);
180     cudaFree(deviceB);
181     cudaFree(deviceC);
182     free(hostA);
183     free(hostB);
184     free(hostC);
185     cudaEventDestroy( start );
186     cudaEventDestroy( stop );
187     return 0;
188 }

```

hostA:

1.000000	5.000000	6.000000	5.000000
7.000000	7.000000	1.000000	7.000000
0.000000	1.000000	5.000000	5.000000
7.000000	3.000000	4.000000	4.000000

hostB:

6.000000	8.000000	8.000000	1.000000
0.000000	8.000000	8.000000	1.000000
7.000000	6.000000	8.000000	5.000000
1.000000	6.000000	4.000000	0.000000

Размер матрицы 128 x 128

Размер блока 2 x 2

Time: 0.21

Размер блока 4 x 4

Time: 0.21

Размер блока 8 x 8

Time: 0.22

Размер блока 16 x 16

Time: 0.19

Размер блока 32 x 32

Time: 0.22

Размер матрицы 256 x 256

Размер блока 2 x 2

Time: 0.54

Размер блока 4 x 4

Time: 0.49

Размер блока 8 x 8

Time: 0.42

Размер блока 16 x 16

Time: 0.45

Размер блока 32 x 32

Time: 0.49

Размер матрицы 512 x 512

Размер блока 2 x 2

Time: 2.35

Размер блока 4 x 4

Time: 1.63

Размер блока 8 x 8

Time: 1.49

Размер блока 16 x 16

Time: 1.37

Размер блока 32 x 32

Time: 1.35

Размер матрицы 1024 x 1024

Размер блока 2 x 2

Time: 13.40

Размер блока 4 x 4

Time: 7.87

Размер блока 8 x 8

Time: 6.97

Размер блока 16 x 16

Time: 5.94

Размер блока 32 x 32

Time: 6.46

Размер матрицы 2048 x 2048

Размер блока 2 x 2

Time: 131.52

Размер блока 4 x 4

Time: 47.33

Размер блока 8 x 8

Time: 29.98

Размер блока 16 x 16

Time: 26.78

Размер блока 32 x 32

Time: 38.97

2. Ускорить передачу матрицы из CPU в GPU за счет использования Pinned памяти.

```
97  int main(int argc, char* argv[])
98  {
99      int N = 2048;
100     float time;
101     double *hostA;
102     double *hostB;
103     double *hostC;
104     double *C;
105
106     //Allocating memory for hosts
107     cudaMallocHost ((void**) &hostA, N*N * sizeof(double), cudaHostAllocDefault);
108     cudaMallocHost ((void**) &hostB, N*N * sizeof(double), cudaHostAllocDefault);
109     cudaMallocHost ((void**) &hostC, N*N * sizeof(double), cudaHostAllocDefault);
110
111     C = (double*) calloc(N*N, sizeof(double));
112
```

```
113     //Filling matrices
114     fill_matrix(hostA, N);
115     fill_matrix(hostB, N);
116
117     //Printg matrices A and B
118     //printf("hostA:\n");
119     //print_matrix(hostA, N);
120     //printf("hostB:\n");
121     //print_matrix(hostB, N);
122
```

```
123     //Allocating memory for devices
124     double *deviceA;
125     double *deviceB;
126     double * deviceC;
127     cudaMalloc((void **)&deviceA, N*N * sizeof(double));
128     cudaMalloc((void **)&deviceB, N*N * sizeof(double));
129     cudaMalloc((void **)&deviceC, N*N * sizeof(double));
130
131     //Defining dimensions of the block and the grid
132     dim3 threadsPerBlock = dim3(BLOCK_SIZE, BLOCK_SIZE);
133     int N_new = 0;
134     if (N % BLOCK_SIZE == 0)
135     |     N_new = int (N / BLOCK_SIZE);
136     else
137     |     N_new = int (N / BLOCK_SIZE) + 1;
138     dim3 blocksPerGrid = dim3(N_new,N_new);
139
```

```

140 //Defining and initializing variables for computing time
141 cudaEvent_t start, stop;
142 cudaEventCreate(&start);
143 cudaEventCreate(&stop);
144
145 //Recording the computing time
146 cudaEventRecord( start, 0);
147
148 //Copying matrices from host to device
149 cudaMemcpy(deviceA, hostA, N*N * sizeof(double), cudaMemcpyHostToDevice);
150 cudaMemcpy(deviceB, hostB, N*N * sizeof(double), cudaMemcpyHostToDevice);
151
152 //Calling multiplication function
153 matrixmultiplication<<<blocksPerGrid, threadsPerBlock>>>(deviceA, deviceB, deviceC, N);
154
155 //Copying matrices from device to host
156 cudaMemcpy(hostC, deviceC, N*N * sizeof(double), cudaMemcpyDeviceToHost);
157 cudaEventRecord( stop, 0);
158 cudaEventSynchronize( stop );
159
160 //Computing time
161 cudaEventElapsedTime( &time, start, stop);
162 printf("Time: %.2f \n",time);
163
164 //Sequantial multiplication
165 mult_pos(N, hostA, hostB, C);
166
167 //Printg result of CUDA multiplication
168 //printf("hostC:\n");
169 //print_matrix(hostC, N);
170
171 //Printg result of sequential multiplication
172 //printf("C:\n");
173 //print_matrix(C, N);
174
175 ///Comparing results of CUDA multiplication and sequential multiplication
176 compare_matrices(C, hostC, N);
177
178 //Deallocationg the memory
179 cudaFree(deviceA);
180 cudaFree(deviceB);
181 cudaFree(deviceC);
182 cudaFreeHost(hostA);
183 cudaFreeHost(hostB);
184 cudaFreeHost(hostC);
185 cudaEventDestroy( start );
186
185 cudaEventDestroy( start );
186 cudaEventDestroy( stop );
187 return 0;
188 }

```

Размер матрицы 128 x 128

Размер блока 2 x 2

Time: 0.18

Размер блока 4 x 4

Time: 0.19

Размер блока 8 x 8

Time: 0.18

Размер блока 16 x 16

Time: 0.16

Размер блока 32 x 32

Time: 0.16

Размер матрицы 256 x 256

Размер блока 2 x 2

Time: 0.46

Размер блока 4 x 4

Time: 0.32

Размер блока 8 x 8

Time: 0.36

Размер блока 16 x 16

Time: 0.33

Размер блока 32 x 32

Time: 0.30

Размер матрицы 512 x 512

Размер блока 2 x 2

Time: 2.12

Размер блока 4 x 4

Time: 1.25

Размер блока 8 x 8

Time: 1.10

Размер блока 16 x 16

Time: 0.93

Размер блока 32 x 32

Time: 0.98

Размер матрицы 1024 x 1024

Размер блока 2 x 2

Time: 16.62

Размер блока 4 x 4

Time: 6.32

Размер блока 8 x 8

Time: 4.15

Размер блока 16 x 16

Time: 3.53

Размер блока 32 x 32

Time: 3.59

Размер матрицы 2048 x 2048

Размер блока 2 x 2

Time: 126.24

Размер блока 4 x 4

Time: 44.93

Размер блока 8 x 8

Time: 24.87

Размер блока 16 x 16

Time: 19.64

Размер блока 32 x 32

Time: 19.76