

Міністерство освіти і науки України
Львівський національний університет Ім. І. Франка
Факультет прикладної математики та інформатики
Кафедра дискретного аналізу та інтелектуальних систем

Паралельні та розподілені обчислення

Лабораторна робота №8

Програмно-апаратна архітектура паралельних обчислень CUDA

Роботу виконала:

Студентка ПМІ-33

Багінська Маргарита

Прийняв:

доц. Пасічник Т.В.

Львів 2023

Тема: Програмно-апаратна архітектура паралельних обчислень CUDA

Мета: Реалізувати алгоритм множення матриць використовуючи програмно-апаратну архітектуру паралельних обчислень CUDA.

Виконання

- Файл **kernel.h**

```
1  #ifndef KERNEL_H_
2  #define KERNEL_H_
3
4  void matrixMultiplication(float* A, float* B, float* C, int N);
5
6  #endif
```

Файл "kernel.h" містить оголошення функції `matrixMultiplication`, яка призначена для виконання множення двох матриць, використовуючи технологію CUDA.

- Файл **dev_array.h**

```
1  #ifndef _DEV_ARRAY_H_
2  #define _DEV_ARRAY_H_
3
4  #include <stdexcept>
5  #include <algorithm>
6  #include <cuda_runtime.h>
7
8  template <class T>
9  class dev_array
10 {
11     // public functions
12 public:
13     explicit dev_array()
14         : start_(0),
15           end_(0)
16     {}
17
18     // constructor
19     explicit dev_array(size_t size)
20     {
21         allocate(size);
22     }
23
24     // destructor
25     ~dev_array()
26     {
27         free();
28     }
29
30     // resize the vector
31     void resize(size_t size)
32     {
33         free();
34         allocate(size);
35     }
36
37     // get the size of the array
38     size_t getSize() const
39     {
40         return end_ - start_;
41     }
42
43     // get data
44     const T* getData() const
45     {
46         return start_;
47     }
48
49     T* getData()
50     {
51         return start_;
52     }
53
54     // set
55     void set(const T* src, size_t size)
56     {
57         size_t min = std::min(size, getSize());
58         cudaError_t result = cudaMemcpy(start_, src, min * sizeof(T), cudaMemcpyHostToDevice);
59         if (result != cudaSuccess)
60         {
61             throw std::runtime_error("failed to copy to device memory");
62         }
63     }
64
65     // get
66     void get(T* dest, size_t size)
67     {
68         size_t min = std::min(size, getSize());
69         cudaError_t result = cudaMemcpy(dest, start_, min * sizeof(T), cudaMemcpyDeviceToHost);
70         if (result != cudaSuccess)
71         {
72             throw std::runtime_error("failed to copy to host memory");
73         }
74     }
75
76     // private functions
77 private:
78     // allocate memory on the device
79     void allocate(size_t size)
80     {
81         cudaError_t result = cudaMalloc((void**)&start_, size * sizeof(T));
82         if (result != cudaSuccess)
83         {
84             start_ = end_ = 0;
85             throw std::runtime_error("failed to allocate device memory");
86         }
87         end_ = start_ + size;
88     }
89
90     // free memory on the device
91     void free()
92     {
93         if (start_ != 0)
94         {
95             cudaFree(start_);
96             start_ = end_ = 0;
97         }
98     }
99
100     T* start_;
101     T* end_;
102 };
103 #endif
```

Файл "dev_array.h" містить клас шаблонів `dev_array`, який є обгорткою над пам'яттю пристрою (пам'ять GPU) для простої алокації пам'яті, зміни розміру, зчитування та запису. Цей клас допомагає забезпечити легкий та безпечний спосіб керування пам'яттю пристроїв.

- Файл kernel.cu

```

1  #include <math.h>
2  #include <iostream>
3  #include "cuda_runtime.h"
4  #include "kernel.h"
5  #include <stdlib.h>
6  #include "device_launch_parameters.h"
7
8  using namespace std;
9
10 __global__ void matrixMultiplicationKernel(float* A, float* B, float* C, int N) {
11
12     int ROW = blockIdx.y * blockDim.y + threadIdx.y;
13     int COL = blockIdx.x * blockDim.x + threadIdx.x;
14
15     float tmpSum = 0;
16
17     if (ROW < N && COL < N) {
18         // each thread computes one element of the block sub-matrix
19         for (int i = 0; i < N; i++) {
20             tmpSum += A[ROW * N + i] * B[i * N + COL];
21         }
22         C[ROW * N + COL] = tmpSum;
23     }
24 }
25
26
27 void matrixMultiplication(float* A, float* B, float* C, int N) {
28
29     // declare the number of blocks per grid and the number of threads per block
30     // use 1 to 512 threads per block
31     dim3 threadsPerBlock(N, N);
32     dim3 blocksPerGrid(1, 1);
33     if (N * N > 512) {
34         threadsPerBlock.x = 512;
35         threadsPerBlock.y = 512;
36         blocksPerGrid.x = ceil(double(N) / double(threadsPerBlock.x));
37         blocksPerGrid.y = ceil(double(N) / double(threadsPerBlock.y));
38     }
39
40     matrixMultiplicationKernel << blocksPerGrid, threadsPerBlock >> (A, B, C, N);
41 }
42

```

Файл "kernel.cu" містить визначення функції **matrixMultiplication**, яка використовує ядро CUDA **matrixMultiplicationKernel** для перемноження двох матриць на GPU. Основні частини коду:

__global__ void matrixMultiplicationKernel(float* A, float* B, float* C, int N):

Це ядро CUDA, яке виконує матричне множення на рівні елементів. Кожен потік відповідає за обчислення одного елемента результуючої матриці C. Ітерація через рядки і стовпці вхідних матриць A та B відбувається всередині цього ядра. Ядро виконується у паралель для кожної комбінації `blockIdx` та `threadIdx`, середнє значення цих значень перемножується з розмірами блоку, щоб отримати індекс рядка (ROW) та стовпця (COL).

Функція **void matrixMultiplication(float* A, float* B, float* C, int N):**

керівництво CUDA, яке приймає вказівники на дані матриць A, B та C та розмір матриці N. Функція налаштовує рядки, стовпці блоків на сітці та потоки на блок для керування виконанням ядра. Він використовує блоки та потоки для паралельного обчислення елементів результуючої матриці C.

- Файл matrixmul.cu

```

1 #include <iostream>
2 #include <vector>
3 #include <stdlib.h>
4 #include <time.h>
5 #include <cuda_runtime.h>
6 #include <kernel.h>
7 #include <dev_array.h>
8 #include <math.h>
9 #include <chrono>
10
11 using namespace std;
12 using namespace std::chrono;
13
14 int main()
15 {
16     // Perform matrix multiplication C = A*B
17     // where A, B and C are NxM matrices
18     int N = 30000;
19     int SIZE = N * N;
20
21     // Allocate memory on the host
22     vector<float> h_A(SIZE);
23     vector<float> h_B(SIZE);
24     vector<float> h_C(SIZE);
25
26     // Initialize matrices on the host
27     for (int i = 0; i < N; i++) {
28         for (int j = 0; j < N; j++) {
29             h_A[i * N + j] = sin(i);
30             h_B[i * N + j] = cos(j);
31         }
32     }
33
34     // Measure time taken by CUDA
35     auto start_cuda = high_resolution_clock::now();
36
37     // Allocate memory on the device
38     dev_array<float> d_A(SIZE);
39     dev_array<float> d_B(SIZE);
40     dev_array<float> d_C(SIZE);
41
42     d_A.set(&h_A[0], SIZE);
43     d_B.set(&h_B[0], SIZE);
44
45     matrixMultiplication(d_A.getData(), d_B.getData(), d_C.getData(), N);
46     cudaDeviceSynchronize();
47
48     d_C.get(&h_C[0], SIZE);
49     cudaDeviceSynchronize();
50
51     auto stop_cuda = high_resolution_clock::now();
52     auto duration_cuda = duration_cast<milliseconds>(stop_cuda - start_cuda);
53
54     float* cpu_C;
55     cpu_C = new float[SIZE];
56
57     // Measure time taken by CPU
58     auto start_cpu = high_resolution_clock::now();
59
60     // Now do the matrix multiplication on the CPU
61     float sum;
62     for (int row = 0; row < N; row++) {
63         for (int col = 0; col < N; col++) {
64             sum = 0.f;
65             for (int n = 0; n < N; n++) {
66                 sum += h_A[row * N + n] * h_B[n * N + col];
67             }
68             cpu_C[row * N + col] = sum;
69         }
70     }
71
72     auto stop_cpu = high_resolution_clock::now();
73     auto duration_cpu = duration_cast<milliseconds>(stop_cpu - start_cpu);
74
75     cout << "[" << N << "]" << "[" << N << "]" << endl;
76     cout << "CUDA execution time: " << duration_cuda.count() << " ms = " << duration_cuda.count() * 0.001 << " s" << endl;
77     cout << "CPU execution time: " << duration_cpu.count() << " ms = " << duration_cpu.count() * 0.001 << " s" << endl;
78
79     // Release device memory
80     cudaFree(d_A.getData());
81     cudaFree(d_B.getData());
82     cudaFree(d_C.getData());
83
84     // Release host memory
85     delete[] cpu_C;
86
87     return 0;
88 }

```

В цьому файлі реалізовано порівняння часу виконання матричного множення на CPU та GPU (CUDA). Використовуючи структуру chrono, код обчислює час виконання операцій на CPU та GPU.

Результати

Згадаємо результати третьої лабораторної роботи, де ми множили матриці синхронно на паралельно на CPU.

```

[2000][2000] * [2000][2000]
Sync time ~ 00:00:54.6495798
Async time ~ 00:00:09.2754881 with 8 threads
Acceleration ~ 5,891827924397854
Efficiency ~ 0,7364784905497318

```

Тепер результати CUDA:

```

[2000][2000]
CUDA execution time: 153 ms = 0.153 s
CPU execution time: 56636 ms = 56.636 s

```

Під CPU execution time мається на увазі час роботи послідовного алгоритму, але навіть порівняно з часом паралельного алгоритму на результати CUDA вражають.

Збільшимо розмірність.

```
[4000][4000]  
CUDA execution time: 1907 ms = 1.907 s  
CPU execution time: 545581 ms = 545.581 s
```

Різниця насправді колосальна. Те, що синхронно рахувалось CPU 9 хвилин, на CUDA ядрах виконалось за менш ніж 2 секунди.

Висновки: У результаті виконання лабораторної роботи було реалізовано алгоритм множення матриць використовуючи програмно-апаратну архітектуру паралельних обчислень CUDA.