



## 1 Постановка задачи

Условие задачи (№5) звучит следующим образом: "Задан вещественный массив. Найти сумму элементов, размещенных в тредах с четными номерами".

## 2 Описание тестового стенда

Вычисления проводились на выданном вычислительном узле:

- Операционная система: Ubuntu 20.04.4 LTS
- Процессор: Intel(R) Xeon(R) E-2136 CPU @ 3.30GHz, 6 ядер, 2 потока на ядро
- Доступная оперативная память: 62-63 ГБ
- GPU: Quadro P2000
- Компилятор: nvcc (NVIDIA (R) Cuda compiler driver, Cuda compilation tools, release 12.0, V12.0.140)

Время работы алгоритмов вычислялось при помощи команды **clock()** из стандартной библиотеки **time** языка **C++**. Проводилось по три серии экспериментов, затем полученные значения времен усреднялись.

Проверка правильности суммирования осуществлялась сравнением с последовательным вычислением на CPU (время не включается в время выполнения).

## 3 Описание алгоритма решения задачи

Условие задачи требует вычислить сумму элементов в тредах с четными номерами, что при последовательном размещении массива в памяти соответствует вычислению суммы элементов на четных позициях. В последовательной реализации для данной цели использовался цикл, в котором считалась сумма каждого второго элемента. Код:

```
1 // Sum elements on even positions on CPU
2 double SumEvenCPU(double* hostArray, int N){
3     double result = 0.0;
4     // iterate over half of the array size
5     for (int i = 0; i < N / 2; i++){
6         // ... end sum every second element
7         result+= hostArray[2*i];
8     }
9     return result;
10 }
```

Идея использования нитей CUDA в такой сортировке заключается в том, что можно разбить массив на блоки, и средствами нитей вычислять сумму элементов на четных местах сначала в каждом блоке (провести редукцию внутри блока), а затем на хосте просуммировать результаты поблочных вычислений (финальная редукция). Реализация кода нити и алгоритма представлена ниже:

```
1 // Thread code for sum of the elements in the even threads
2 __global__ void SumEven(double* array, int N, double* deviceRes){
3     // shared array for threads in the block
4     __shared__ double helper[THREADS_PER_BLOCK];
5
6     // Get thread ID
7     int idx = blockDim.x * blockIdx.x + threadIdx.x;
8     // Create thread variable for storing sum
9     double sum = 0.0;
10    // If thread ID is even
11    if (idx % 2 == 0){
12        helper[threadIdx.x] = array[idx];
13    }else{
14        // ... or if odd sum equals zero
15        helper[threadIdx.x] = 0.0;
16    }
17
18    __syncthreads();
19
20    // Reduction
21    int amountOfThreads = blockDim.x / 2;
22    while (amountOfThreads > 0){
23        if (threadIdx.x < amountOfThreads){
24            helper[threadIdx.x] += helper[threadIdx.x + amountOfThreads];
25        }
26
27        __syncthreads();
28
29        amountOfThreads /= 2;
30    }
31
32    // Save the results for block
33    if (threadIdx.x == 0) {
34        deviceRes[blockIdx.x] = helper[0];
35    }
36 }
37
38 // Sum elements on even positions using CUDA
39 double SumEvenCUDA(double* hostArray, int N){
40     // Make two arrays on device variables
41     double *deviceArray, *deviceRes;
42     // helper for intermediate results
43     double *hostInterRes;
44     // Variable to accumulate intermediate results
45     double res = 0.0;
46
47     // Thread blocks per Grid
48     int gridDim = (N + THREADS_PER_BLOCK - 1) / THREADS_PER_BLOCK;
```

```

49     size_t interResSize = gridDim * sizeof(double);
50
51     // Create vectors
52     hostInterRes = (double*)malloc(interResSize);
53     cudaMalloc(&deviceArray, N*sizeof(double));
54     cudaMalloc(&deviceRes, interResSize);
55     // Copy from host to device
56     cudaMemcpy(deviceArray, hostArray, N * sizeof(double), cudaMemcpyHostToDevice);
57
58     // Sum
59     SumEven<<<gridDim, THREADS_PER_BLOCK>>>(deviceArray, N, deviceRes);
60
61     // Copy from device to host
62     cudaMemcpy(hostInterRes, deviceRes, interResSize, cudaMemcpyDeviceToHost);
63
64     // Free memory
65     cudaFree(deviceArray);
66     cudaFree(deviceRes);
67
68     // Final reduction
69     for (int i = 0; i < gridDim; i++){
70         res += hostInterRes[i];
71     }
72     free(hostInterRes);
73     return res;
74 }

```

---

## 4 Программный код

Полный код программы представлен в Приложении 1 в разделе 8.

## 5 Результаты измерений

Вывод программы представлен в Приложении 2 в разделе 9.

## 6 Анализ и обсуждение результатов

Рассмотрим затраты времени при проверке алгоритма на векторах разного размера. На рисунке 1 представлена зависимость времени, затраченного на вычисления, от размера массива. При вычислениях на CPU зависимость является линейной ( $O(N)$ ), что наблюдается на всём интервале размеров массива (с учетом погрешностей усреднения). При вычислениях на GPU же, до размера массива  $10^5$  элементов, затраты времени практически постоянные и много больше времени для CPU, что связано с затратами на создание нитей, пересылку данных. При увеличении размера массива, зависимость становится так же линейной, причем угол наклона меньше, чем для CPU, что говорит о некотором ускорении. И действительно, начиная с размера массива в  $10^8$  элементов, вычисления на GPU начинают выигрывать во времени у CPU.

Аналогичные тренды можно наблюдать и на графике ускорения GPU 2. Рост ускорения на первом участке до  $10^5$  элементов обусловлен тем, что вычисления на CPU замедлялись,

а на GPU были бы постоянными по времени. На оставшемся участке наблюдается истинное увеличение ускорения, обусловленное эффективностью использования GPU.

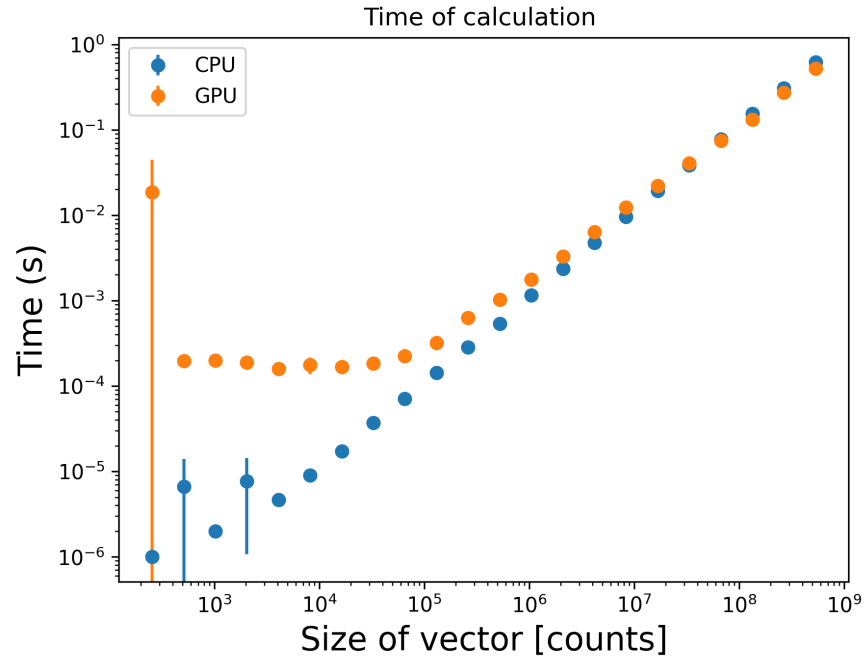


Рис. 1: Зависимость времени вычисления от размера массива.

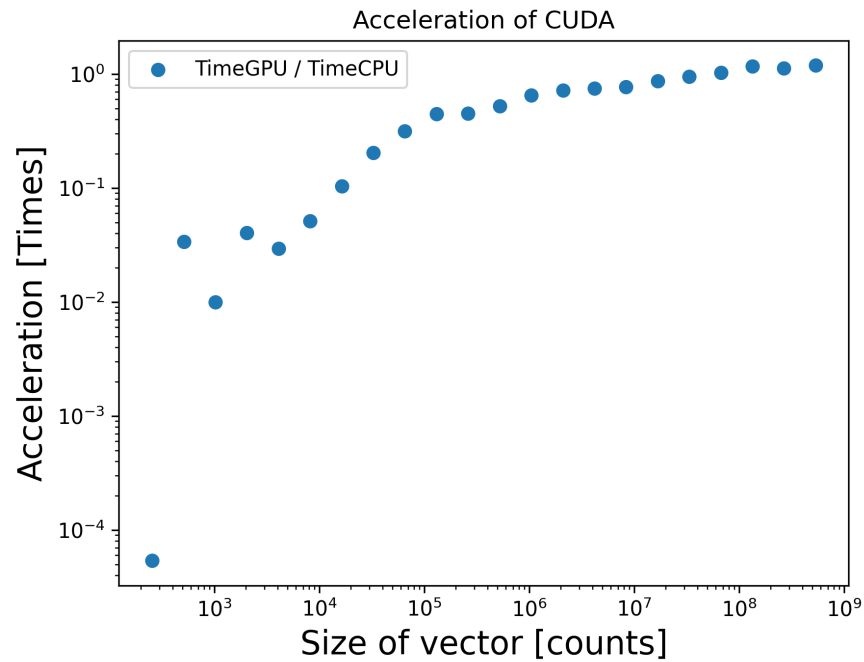


Рис. 2: Зависимость ускорения вычисления от размера массива.

## 7 Выводы

Поставленные задачи были выполнены. Было проведено сравнение работы алгоритмов на CPU и GPU и наблюдалось ускорение работы, которое может обеспечить GPU, и был определен размер массива, для которого использование GPU становится оправданным.

## 8 Приложение 1. Код программы

---

```
1  #include <iostream>
2  #include <stdio.h>
3  #include <time.h>
4  #include <iomanip>
5  #define THREADS_PER_BLOCK 1024
6
7  using namespace std;
8
9
10 __global__ void SumEven(double* array, int N, double* deviceRes);
11
12
13 double SumEvenCUDA(double* hostArray, int N);
14
15 double SumEvenCPU(double* hostArray, int N);
16
17 void PrintArray(const double* array, const int n);
18
19 void measureNoCUDA(int N);
20
21 void measureCUDA(int N);
22
23
24 int main (int argc, char * argv []){
25     srand(static_cast <unsigned> (time(0)));
26
27     cout << "CPU:" << endl;
28     cout << "N          Time1      Time2      Time3" << endl;
29     for (int n = 256; n < 1000000000; n*=2 ){
30         cout<< setw(9) << n << " ";
31         for (int i = 0; i < 3; i++){
32             measureNoCUDA(n);
33             cout << " ";
34         }
35         cout << endl;
36     }
37
38     cout << "GPU:" << endl;
39     cout << "N          Time1      Time2      Time3" << endl;
40     //cout << "N          Time" << endl;
41     for (int n = 256; n < 1000000000; n*=2 ){
42         cout<< setw(9) << n << " ";
43         for (int i = 0; i < 3; i++){
44             measureCUDA(n);
```

```

45         cout << " ";
46     }
47     cout << endl;
48 }
49 cout << endl;
50 return 0;
51 }
52
53 // Thread code for sum of the elements in the even threads
54 __global__ void SumEven(double* array, int N, double* deviceRes){
55     // shared array for threads in the block
56     __shared__ double helper[THREADS_PER_BLOCK];
57
58     // Get thread ID
59     int idx = blockDim.x * blockIdx.x + threadIdx.x;
60     // Create thread variable for storing sum
61     double sum = 0.0;
62     // If thread ID is even
63     if (idx % 2 == 0){
64         helper[threadIdx.x] = array[idx];
65     }else{
66         // ... or if odd sum equals zero
67         helper[threadIdx.x] = 0.0;
68     }
69
70     __syncthreads();
71
72     // Reduction
73     int amountOfThreads = blockDim.x / 2;
74     while (amountOfThreads > 0){
75         if (threadIdx.x < amountOfThreads){
76             helper[threadIdx.x] += helper[threadIdx.x + amountOfThreads];
77         }
78
79         __syncthreads();
80
81         amountOfThreads /= 2;
82     }
83
84     // Save the results for block
85     if (threadIdx.x == 0) {
86         deviceRes[blockIdx.x] = helper[0];
87     }
88 }
89
90 // Sum elements on even positions using CUDA
91 double SumEvenCUDA(double* hostArray, int N){
92     // Make two arrays on device variables
93     double *deviceArray, *deviceRes;
94     // helper for intermediate results
95     double *hostInterRes;
96     // Variable to accumulate intermediate results
97     double res = 0.0;
98

```

```

99 // Thread blocks per Grid
100 int gridDim = (N + THREADS_PER_BLOCK - 1) / THREADS_PER_BLOCK;
101 size_t interResSize = gridDim * sizeof(double);
102
103 // Create vectors
104 hostInterRes = (double*)malloc(interResSize);
105 cudaMalloc(&deviceArray, N*sizeof(double));
106 cudaMalloc(&deviceRes, interResSize);
107 // Copy from host to device
108 cudaMemcpy(deviceArray, hostArray, N * sizeof(double), cudaMemcpyHostToDevice);
109
110 // Sum
111 SumEven<<<gridDim, THREADS_PER_BLOCK>>>(deviceArray, N, deviceRes);
112
113 // Copy from device to host
114 cudaMemcpy(hostInterRes, deviceRes, interResSize, cudaMemcpyDeviceToHost);
115
116 // Free memory
117 cudaFree(deviceArray);
118 cudaFree(deviceRes);
119
120 // Final reduction
121 for (int i = 0; i < gridDim; i++){
122     res += hostInterRes[i];
123 }
124 free(hostInterRes);
125 return res;
126 }
127
128 // Sum elements on even positions on CPU
129 double SumEvenCPU(double* hostArray, int N){
130     double result = 0.0;
131     // iterate over half of the array size
132     for (int i = 0; i < N / 2; i++){
133         // ... end sum every second element
134         result+= hostArray[2*i];
135     }
136     return result;
137 }
138
139 // measure time taken for all steps on CPU CUDA
140 void measureNoCUDA(int N){
141     // Create array
142     double *array = (double *)malloc(N * sizeof(double));
143     // Create variables for saving time results
144     clock_t start_time, end_time;
145     float timeResult;
146     double resCPU = 0.0;
147     // Fill array with random numbers
148     for (int i = 0; i < N; i++) {
149         array[i] = static_cast <double> (rand()) / static_cast <double> (RAND_MAX);
150     }
151
152     // Calculation for the array

```

```

153     start_time = clock();
154     resCPU = SumEvenCPU(array, N);
155     end_time = clock();
156     // Time taken
157     timeResult = (float)(end_time - start_time) / CLOCKS_PER_SEC;
158
159     // Free memory
160     free(array);
161     // Print time info
162     cout << setw(10) << timeResult;
163 }
164
165 // measure time taken for algorithm on GPU CUDA
166 void measureCUDA(int N){
167     // Create array
168     double *array = (double *)malloc(N * sizeof(double));
169     // Create variables for saving time results
170     clock_t start_time, end_time;
171     float timeResult;
172     double resGPU = 0.0, resCPU = 0.0;
173     // Fill array with random numbers
174     for (int i = 0; i < N; i++) {
175         array[i] = static_cast <double> (rand()) / static_cast <double> (RAND_MAX);
176     }
177
178     // Calculation for the array
179     start_time = clock();
180     resGPU = SumEvenCUDA(array, N);
181     end_time = clock();
182     // Time taken
183     timeResult = (float)(end_time - start_time) / CLOCKS_PER_SEC;
184
185     // compare to the CPU result
186     resCPU = SumEvenCPU(array, N);
187     //cout << "CPU : " << resCPU << endl;
188     //cout << "GPU : " << resGPU << endl;
189     if (abs(resCPU - resGPU) < 1e-2){
190         //cout << timeDotProduct << " ";
191     }
192     else {
193         free(array);
194         cout << "ERROR, RESULT DOESN'T MATCH CPU" << endl;
195         return;
196     }
197
198     // Free memory
199     free(array);
200     // Print time info
201     cout << setw(10) << timeResult;
202 }
203
204 // Print double array into console
205 void PrintArray(const double* array, const int n){
206     for(int i = 0; i < n; i++){

```



```

207     cout << array[i] << ' ';
208 }
209 cout << endl;
210 }

```

---

## 9 Приложение 2. Результаты

---

1	CPU:			
2	N	Time1	Time2	Time3
3	256	1e-06	1e-06	1e-06
4	512	1.7e-05	1e-06	2e-06
5	1024	2e-06	2e-06	2e-06
6	2048	3e-06	3e-06	1.7e-05
7	4096	5e-06	4e-06	5e-06
8	8192	9e-06	9e-06	9e-06
9	16384	1.8e-05	1.7e-05	1.7e-05
10	32768	3.3e-05	4.6e-05	3.3e-05
11	65536	7.1e-05	7.1e-05	7.1e-05
12	131072	0.000142	0.000143	0.000142
13	262144	0.000286	0.000284	0.000283
14	524288	0.000559	0.000532	0.000527
15	1048576	0.001163	0.001142	0.001151
16	2097152	0.002417	0.002354	0.002350
17	4194304	0.004758	0.004843	0.004733
18	8388608	0.009491	0.009605	0.009665
19	16777216	0.019017	0.019326	0.019348
20	33554432	0.038119	0.038327	0.038573
21	67108864	0.077164	0.076971	0.077040
22	134217728	0.153913	0.153270	0.153367
23	268435456	0.306564	0.309114	0.306612
24	536870912	0.614487	0.615426	0.614708
25				
26	GPU:			
27	N	Time1	Time2	Time3
28	256	0.055382	0.000203	0.000208
29	512	0.000184	0.000203	0.000204
30	1024	0.000200	0.000200	0.000199
31	2048	0.000202	0.000184	0.000182
32	4096	0.000166	0.000155	0.000156
33	8192	0.000231	0.000149	0.000147
34	16384	0.000161	0.000160	0.000179
35	32768	0.000196	0.000180	0.000174
36	65536	0.000231	0.000221	0.000220
37	131072	0.000325	0.000329	0.000305
38	262144	0.000634	0.000626	0.000629
39	524288	0.001025	0.001019	0.001038
40	1048576	0.001774	0.001744	0.001770
41	2097152	0.003340	0.003301	0.003262
42	4194304	0.006390	0.006367	0.006392
43	8388608	0.012403	0.012343	0.012417
44	16777216	0.021943	0.022332	0.022116
45	33554432	0.040394	0.040898	0.040242

46	67108864	0.076952	0.074762	0.072239
47	134217728	0.132258	0.131437	0.131514
48	268435456	0.288641	0.270241	0.258923
49	536870912	0.515918	0.517701	0.517307

---