

Processing Time in the GPU and CPU

If we were to test the processing time of the same operation over and over, say $z = 2x + y$, performed both by the CPU and GPU, there would be a noticeable difference if the size of the vectors to store this repeated operations is enormous (102,400 cells). Consider the image below:

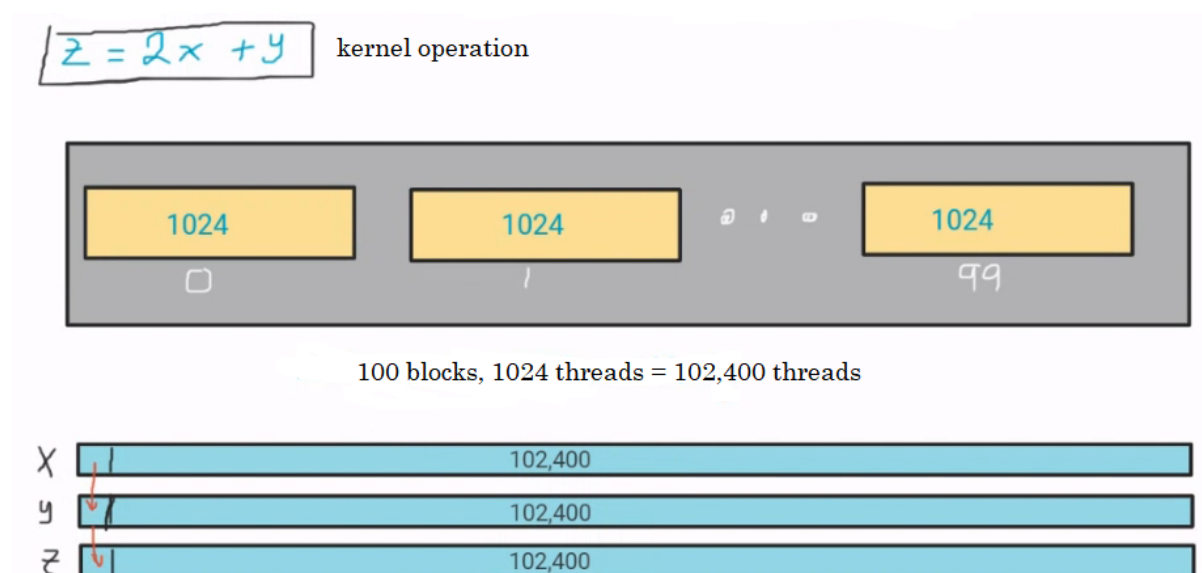


Figure 1: img

Implementation

```

1  #include "cuda_runtime.h"
2  #include "device_launch_parameters.h"
3
4  #include <stdio.h>
5  #include <iostream>
6  #include <time.h>
7
8  using namespace std;
9
10 __global__ void GPU_fn(int* x, int* y, int* z) {
11     int gId = blockIdx.x * blockDim.x + threadIdx.x;
12     z[gId] = 2 * x[gId] + y[gId];
13 }
14
15
16 __host__ void CPU_fn(int* x, int* y, int* z, int vecSize) {
17     for (int i = 0; i < vecSize; i++) {
18         z[i] = 2 * x[i] + y[i];
19     }
20 }

```

```

21
22 __host__ void checkCUDAError(const char* msg) {
23     cudaError_t error;
24     cudaDeviceSynchronize();
25     error = cudaGetLastError();
26     if (error != cudaSuccess) {
27         printf("ERROR %d: %s (%s)\n", error,
28             cudaGetErrorString(error), msg);
29     }
30 }
31 __host__ void validate(int* result_CPU, int* result_GPU, int N
32 ) {
33     for (int i = 0; i < N; i++) {
34         if (result_CPU[i] != result_GPU[i]) {
35             printf("The vectors are not equal\n");
36             return;
37         }
38     }
39     printf("Kernel validated successfully\n");
40     return;
41 }
42 int main()
43 {
44     cudaDeviceProp prop;
45     cudaGetDeviceProperties(&prop, 0);
46     printf("maxThreadsPerBlock: %d\n", prop.maxThreadsPerBlock
47 ); // 1024 in all its block dimension
48     printf("maxThreadsDim[0]: %d\n", prop.maxThreadsDim[0]);
49     // 1024 in a block's x dim
50     // 100 blocks x 1024 threads = 102 400 threads
51     // x,y,z vectors of size 1024
52
53     int numBlocks = 100000; // add one zero and you get ERROR
54     // 2: run out of global memory
55     int numThreadsPerBlock = 1024;
56     int vecSize = numBlocks * numThreadsPerBlock;
57
58     int* hostx = (int*)malloc(vecSize * sizeof(int));
59     int* hosty = (int*)malloc(vecSize * sizeof(int));
60     int* hostzCPU = (int*)malloc(vecSize * sizeof(int));
61     int* hostzGPU = (int*)malloc(vecSize * sizeof(int));
62
63     int* devx, * devy, * devz;
64     cudaMalloc((void**)&devx, vecSize * sizeof(int));
65     checkCUDAError("Error at cudaMalloc: devx");
66     cudaMalloc((void**)&devy, vecSize * sizeof(int));
67     checkCUDAError("Error at cudaMalloc: devy");
68     cudaMalloc((void**)&devz, vecSize * sizeof(int));

```

```
67     checkCUDAError("Error at cudaMalloc: devz");
68
69     for (int i = 0; i < vecSize; i++) {
70         hostx[i] = 1;
71         hosty[i] = 2;
72     }
73
74     cudaMemcpy(devx, hostx, vecSize * sizeof(int),
75               cudaMemcpyHostToDevice);
76     cudaMemcpy(devy, hosty, vecSize * sizeof(int),
77               cudaMemcpyHostToDevice);
78
79     dim3 block(numThreadsPerBlock);
80     dim3 grid(numBlocks);
81
82     cudaEvent_t startGPU;
83     cudaEvent_t endGPU;
84     cudaEventCreate(&startGPU);
85     cudaEventCreate(&endGPU); // be able to mark the time
86     cudaEventRecord(startGPU); // save current time
87     GPU_fn << <grid, block >> > (devx, devy, devz);
88     cudaEventRecord(endGPU);
89     cudaEventSynchronize(endGPU); // so that cudaEventRecord(
90                                     startGPU) and cudaEventSynchronize(endGPU) are not done
91                                     at the same time
92
93     float elapsedTimeGPU;
94     cudaEventElapsedTime(&elapsedTimeGPU, startGPU, endGPU);
95     cudaMemcpy(hostzGPU, devz, vecSize * sizeof(int),
96               cudaMemcpyDeviceToHost);
97
98     clock_t startCPU = clock(); // save current time
99     CPU_fn(hostx, hosty, hostzCPU, vecSize);
100    clock_t endCPU = clock();
101    float elapsedTimeCPU = endCPU - startCPU;
102    printf("Time elapsed CPU: %f milliseconds\n",
103           elapsedTimeCPU);
104    printf("Time elapsed GPU: %f milliseconds\n",
105           elapsedTimeGPU);
106
107    validate(hostzCPU, hostzGPU, vecSize);
108
109    free(hostx);
110    free(hosty);
111    free(hostzCPU);
112    free(hostzGPU);
113    cudaFree(devx);
114    cudaFree(devy);
115    cudaFree(devz);
116
117    cudaEventDestroy(startGPU);
118    cudaEventDestroy(endGPU);
```

```
111     return 0;  
112 }
```

Output

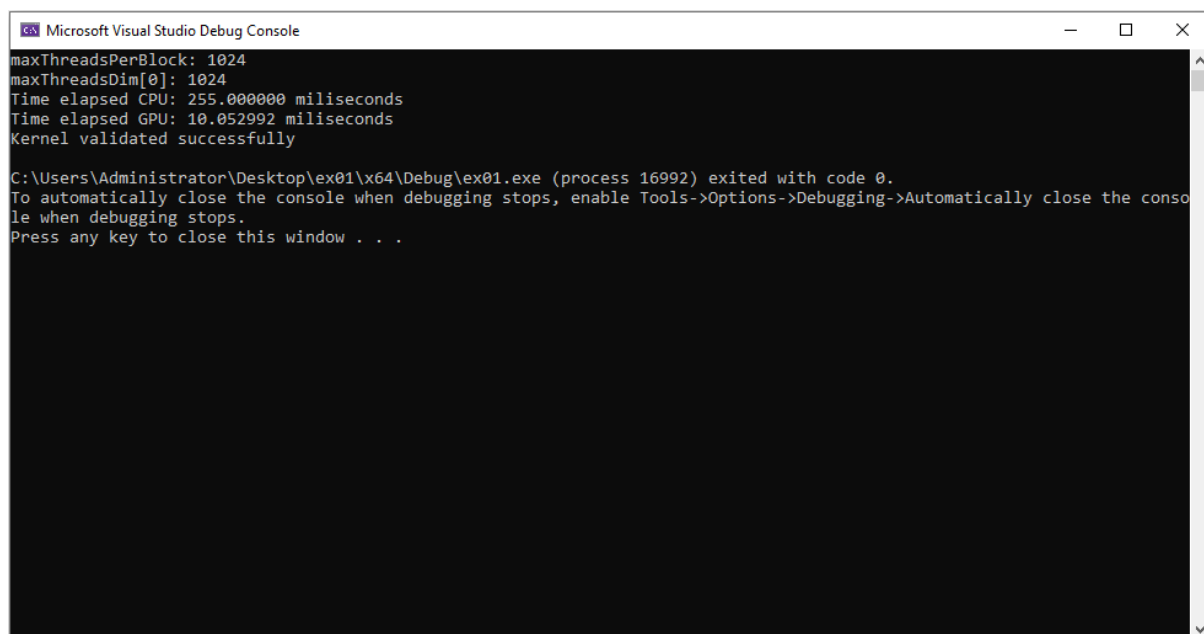


Figure 2: img