

# Metody Programowania Równoległego

## Cuda - VectorAdd

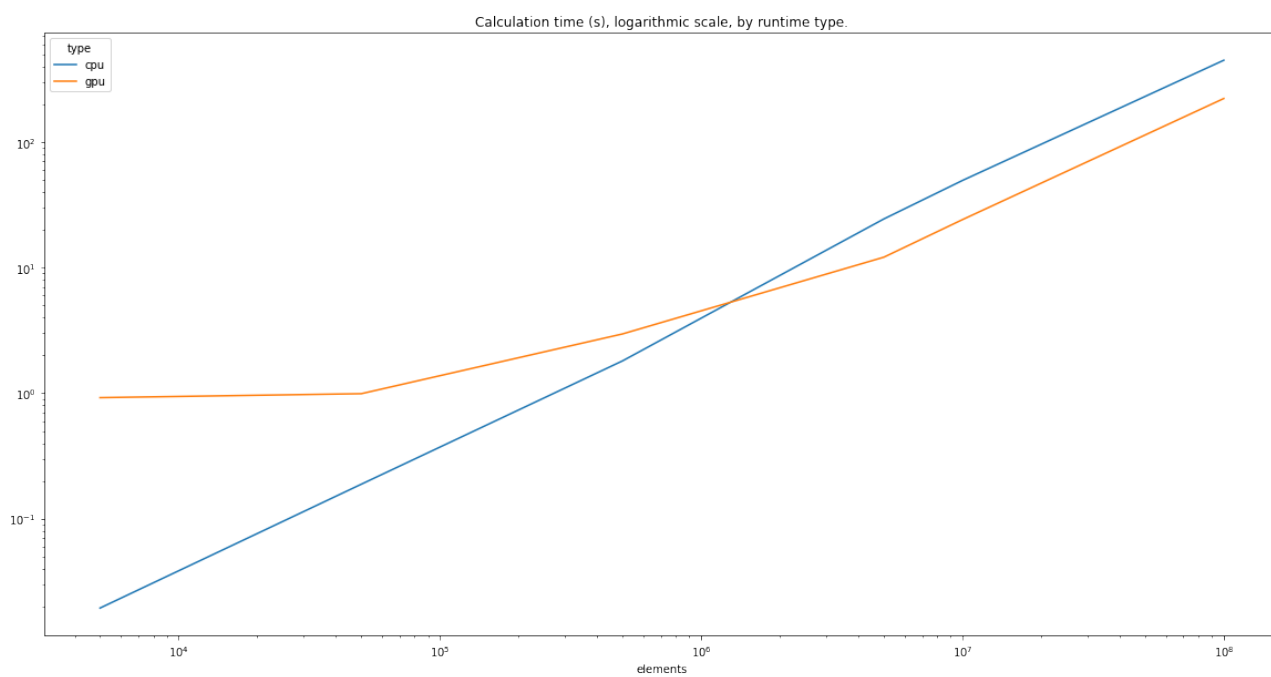
Maciej Trątnowiecki

AGH, Semestr Letni, 2022

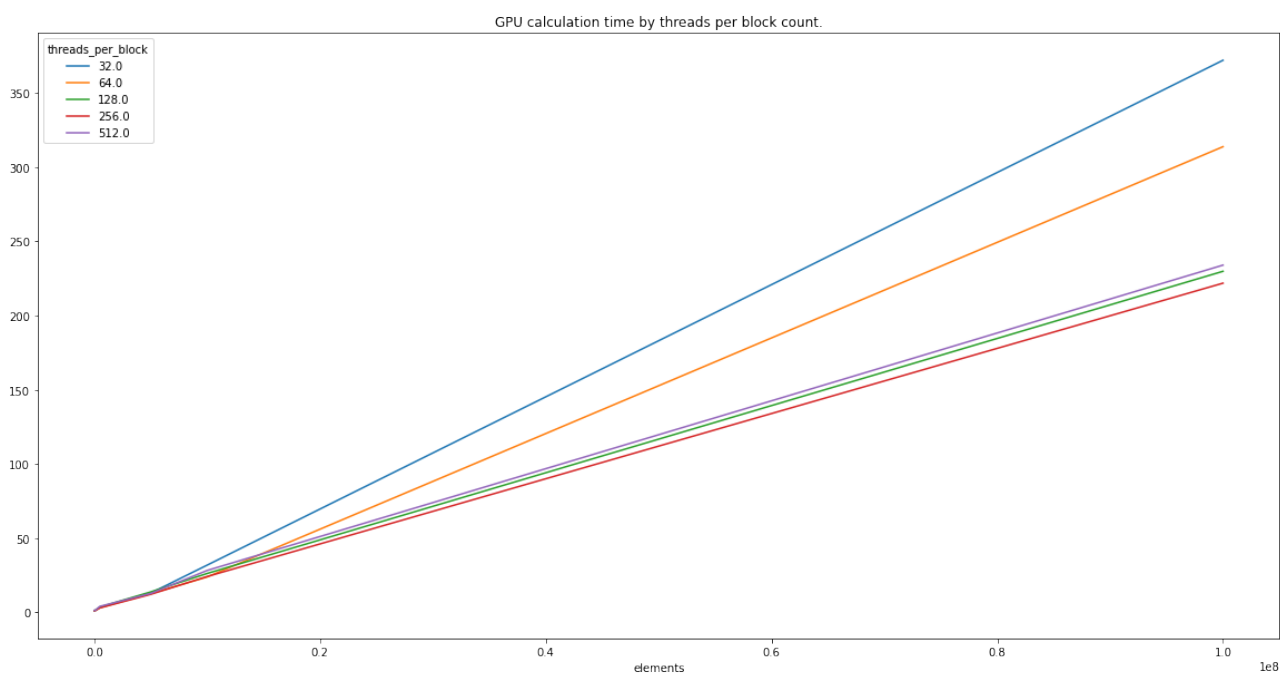
### 1 Pomiary i wyniki

W ramach laboratorium przygotowałem odpowiednio zmodyfikowany kod programu implementującego dodawanie wektorów (opierając się na przygotowanym przez nvidię programie przykładowym) i dokonałem pomiarów czasu wykonania. Wyniki zebrałem w poniższej tabeli oraz zilustrowałem za pomocą poniższych wykresów. Analizując wykresy widzimy, że dla odpowiednio małych rozmiarów wektorów czas wykonania programu może być mniejszy bez użycia akceleracji za pomocą karty graficznej. Wynika to z narzutu czasowego generowanego przez obsługę pamięci karty graficznej i kopiowanie danych. Obserwując wykres wpływu zmiany ilości wątków na blok dla obliczeń wykonywanych z wykorzystaniem karty graficznej zauważyć możemy, że ma on ten sam kształt (w skali logarytmicznej) dla każdego rozmiaru dodawanych wektorów.

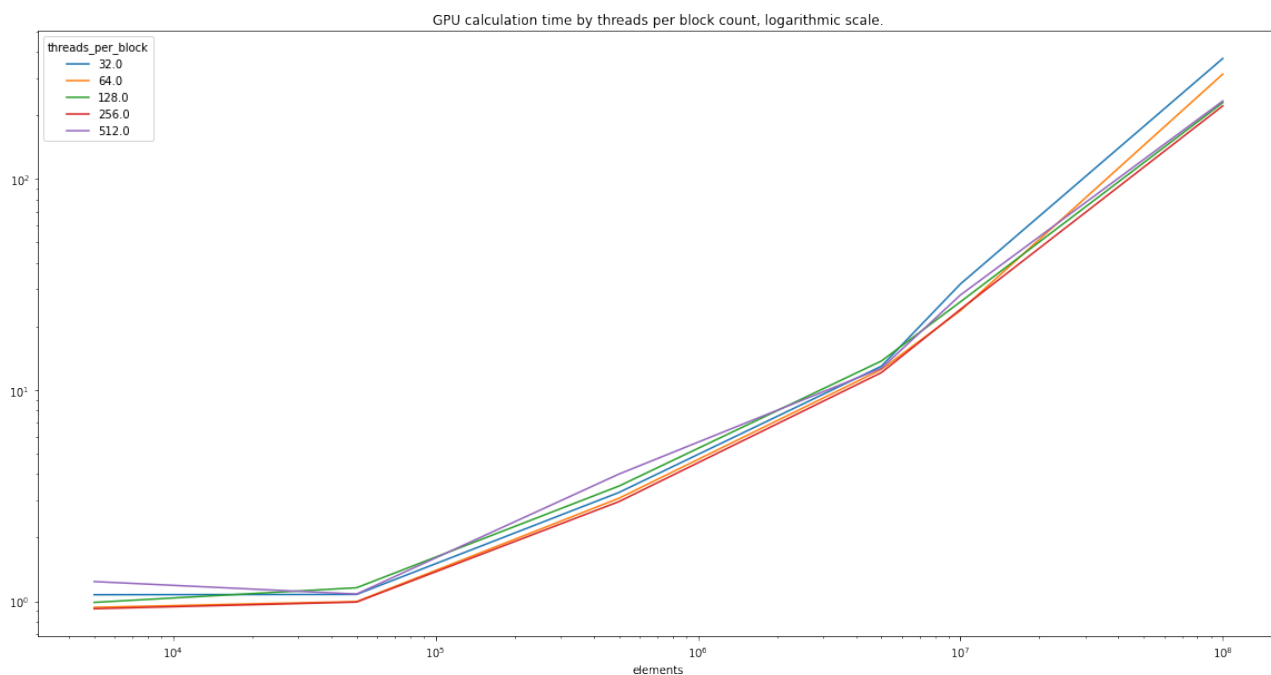
	elements	time	type	threads_per_block	blocks_per_grid
0	5000	0.019456	cpu	nan	nan
1	5000	1.073056	gpu	32.000000	157.000000
2	5000	0.935136	gpu	64.000000	79.000000
3	5000	0.986880	gpu	128.000000	40.000000
4	5000	0.921024	gpu	256.000000	20.000000
5	5000	1.238368	gpu	512.000000	10.000000
6	50000	0.188416	cpu	nan	nan
7	50000	1.054368	gpu	32.000000	1563.000000
8	50000	1.076416	gpu	32.000000	1563.000000
9	50000	0.995840	gpu	64.000000	782.000000
10	50000	1.157600	gpu	128.000000	391.000000
11	50000	0.990336	gpu	256.000000	196.000000
12	50000	1.080864	gpu	512.000000	98.000000
13	500000	1.808384	cpu	nan	nan
14	500000	3.273312	gpu	32.000000	15625.000000
15	500000	3.073728	gpu	64.000000	7813.000000
16	500000	3.510112	gpu	128.000000	3907.000000
17	500000	2.962688	gpu	256.000000	1954.000000
18	500000	4.000768	gpu	512.000000	977.000000
19	5000000	24.389631	cpu	nan	nan
20	5000000	12.969472	gpu	32.000000	156250.000000
21	5000000	12.517504	gpu	64.000000	78125.000000
22	5000000	13.730528	gpu	128.000000	39063.000000
23	5000000	12.097664	gpu	256.000000	19532.000000
24	5000000	12.685152	gpu	512.000000	9766.000000
25	10000000	49.402882	cpu	nan	nan
26	10000000	31.789600	gpu	32.000000	312500.000000
27	10000000	23.785185	gpu	64.000000	156250.000000
28	10000000	26.165279	gpu	128.000000	78125.000000
29	10000000	24.098751	gpu	256.000000	39063.000000
30	10000000	28.189184	gpu	512.000000	19532.000000
31	100000000	447.545349	cpu	nan	nan
32	100000000	372.130005	gpu	32.000000	3125000.000000
33	100000000	313.858063	gpu	64.000000	1562500.000000
34	100000000	229.868637	gpu	128.000000	781250.000000
35	100000000	221.865829	gpu	256.000000	390625.000000
36	100000000	234.037064	gpu	512.000000	195313.000000



Rysunek 1: Czas wykonania programu (w sekundach) w zależności od jednostki wykonującej obliczenia, skala logarytmiczna.



Rysunek 2: Czas wykonania programu (w sekundach) w zależności od jednostki wykonującej obliczenia, skala logarytmiczna.



Rysunek 3: Czas wykonania programu (w sekundach) w zależności od jednostki wykonującej obliczenia, skala logarytmiczna.

## 2 Kod programu

```
1000 /**
1001  * Copyright 1993–2015 NVIDIA Corporation. All rights reserved.
1002  *
1003  * Please refer to the NVIDIA end user license agreement (EULA) associated
1004  * with this source code for terms and conditions that govern your use of
1005  * this software. Any use, reproduction, disclosure, or distribution of
1006  * this software and related documentation outside the terms of the EULA
1007  * is strictly prohibited.
1008  *
1009  */
1010
1011 /**
1012  * Vector addition: C = A + B.
1013  *
1014  * This sample is a very basic sample that implements element by element
1015  * vector addition. It is the same as the sample illustrating Chapter 2
1016  * of the programming guide with some additions like error checking.
1017  */
1018 #include <stdio.h>
1019
1020 // For the CUDA runtime routines (prefixed with "cuda_")
1021 #include <cuda_runtime.h>
1022 #include <helper_cuda.h>
1023 #include "gputimer.h"
1024 #define VERBOSE false
1025
1026 /**
1027  * CUDA Kernel Device code
1028  *
1029  * Computes the vector addition of A and B into C. The 3 vectors have the same
1030  * number of elements numElements.
1031  */
1032 __global__ void
1033 vectorAdd(const float *A, const float *B, float *C, int numElements)
1034 {
1035     int i = blockDim.x * blockIdx.x + threadIdx.x;
1036
1037     if (i < numElements)
1038     {
1039         C[i] = A[i] + B[i];
1040     }
1041 }
1042
1043 void run_on_gpu(size_t size, int numElements, float *h_A, float *h_B, float *h_C, int
1044 threadsPerBlock, int blocksPerGrid) {
1045     // Error code to check return values for CUDA calls
1046     cudaError_t err = cudaSuccess;
1047
1048     // Allocate the device input vector A
1049     float *d_A = NULL;
1050     err = cudaMalloc((void **)&d_A, size);
1051
1052     if (err != cudaSuccess)
1053     {
1054         fprintf(stderr, "Failed to allocate device vector A (error code %s)!\n",
1055             cudaGetErrorString(err));
1056         exit(EXIT_FAILURE);
1057     }
1058
1059     // Allocate the device input vector B
1060     float *d_B = NULL;
1061     err = cudaMalloc((void **)&d_B, size);
1062
1063     if (err != cudaSuccess)
1064     {
1065         fprintf(stderr, "Failed to allocate device vector B (error code %s)!\n",
1066             cudaGetErrorString(err));
1067         exit(EXIT_FAILURE);
1068     }
1069
1070     // Allocate the device output vector C
1071     float *d_C = NULL;
1072     err = cudaMalloc((void **)&d_C, size);
```

```

1072     if (err != cudaSuccess)
1073     {
1074         fprintf(stderr, "Failed to allocate device vector C (error code %s)!\n",
cudaGetErrorString(err));
        exit(EXIT_FAILURE);
1076     }

1078     // Copy the host input vectors A and B in host memory to the device input vectors in
// device memory
1080     #if (VERBOSE == true)
        printf("Copy input data from the host memory to the CUDA device\n");
1082     #endif
    err = cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);

1084     if (err != cudaSuccess)
1085     {
1086         fprintf(stderr, "Failed to copy vector A from host to device (error code %s)!\n",
cudaGetErrorString(err));
1088         exit(EXIT_FAILURE);
1089     }

    err = cudaMemcpy(d_B, h_B, size, cudaMemcpyHostToDevice);

1092     if (err != cudaSuccess)
1093     {
1094         fprintf(stderr, "Failed to copy vector B from host to device (error code %s)!\n",
cudaGetErrorString(err));
1096         exit(EXIT_FAILURE);
1097     }

1098     // Launch the Vector Add CUDA Kernel
1100     #if (VERBOSE == true)
        printf("CUDA kernel launch with %d blocks of %d threads\n", blocksPerGrid,
threadsPerBlock);
1102     #endif
    vectorAdd<<<blocksPerGrid, threadsPerBlock>>>(d_A, d_B, d_C, numElements);
1104     err = cudaGetLastError();

1106     if (err != cudaSuccess)
1107     {
1108         fprintf(stderr, "Failed to launch vectorAdd kernel (error code %s)!\n",
cudaGetErrorString(err));
1110         exit(EXIT_FAILURE);
1111     }

1112     // Copy the device result vector in device memory to the host result vector
// in host memory.
1114     #if (VERBOSE == true)
        printf("Copy output data from the CUDA device to the host memory\n");
1116     #endif
    err = cudaMemcpy(h_C, d_C, size, cudaMemcpyDeviceToHost);

1118     if (err != cudaSuccess)
1119     {
1120         fprintf(stderr, "Failed to copy vector C from device to host (error code %s)!\n",
cudaGetErrorString(err));
1122         exit(EXIT_FAILURE);
1123     }

1124     // Free device global memory
1126     err = cudaFree(d_A);

1128     if (err != cudaSuccess)
1129     {
1130         fprintf(stderr, "Failed to free device vector A (error code %s)!\n",
cudaGetErrorString(err));
1132         exit(EXIT_FAILURE);
1133     }

    err = cudaFree(d_B);

1136     if (err != cudaSuccess)
1137     {
1138         fprintf(stderr, "Failed to free device vector B (error code %s)!\n",
cudaGetErrorString(err));

```

```

1140     exit(EXIT_FAILURE);
1141 }
1142 err = cudaFree(d_C);
1143
1144 if (err != cudaSuccess)
1145 {
1146     fprintf(stderr, "Failed to free device vector C (error code %s)!\n",
1147         cudaGetErrorString(err));
1148     exit(EXIT_FAILURE);
1149 }
1150
1151 void run_on_cpu(size_t size, int numElements, float *h_A, float *h_B, float *h_C) {
1152     for (int i = 0; i < numElements; ++i)
1153         h_C[i] = h_A[i] + h_B[i];
1154 }
1155
1156 /**
1157  * Host main routine
1158  */
1159 int
1160 main(int argc, char *argv[]) {
1161     // Print the vector length to be used, and compute its size
1162     int numElements = 0;
1163     if (argc > 1) {
1164         numElements = atoi(argv[1]);
1165         #if (VERBOSE == true)
1166             printf("Number of elements: %d\n", numElements);
1167         #endif
1168     } else {
1169         printf("Please pass number of elements as command line argument.\n");
1170         exit(EXIT_FAILURE);
1171     }
1172
1173     bool cpu_only = false;
1174
1175     if (argc > 2 && strcmp(argv[2], "true") == 0) {
1176         cpu_only = true;
1177         #if (VERBOSE == true)
1178             printf("Running on CPU only!\n");
1179         #endif
1180     }
1181
1182     size_t size = numElements * sizeof(float);
1183     // printf("[Vector addition of %d elements]\n", numElements);
1184
1185     // Allocate the host input vector A
1186     float *h_A = (float *)malloc(size);
1187
1188     // Allocate the host input vector B
1189     float *h_B = (float *)malloc(size);
1190
1191     // Allocate the host output vector C
1192     float *h_C = (float *)malloc(size);
1193
1194     // Verify that allocations succeeded
1195     if (h_A == NULL || h_B == NULL || h_C == NULL)
1196     {
1197         fprintf(stderr, "Failed to allocate host vectors!\n");
1198         exit(EXIT_FAILURE);
1199     }
1200
1201     // Initialize the host input vectors
1202     for (int i = 0; i < numElements; ++i)
1203     {
1204         h_A[i] = rand() / (float)RAND_MAX;
1205         h_B[i] = rand() / (float)RAND_MAX;
1206     }
1207
1208     int threadsPerBlock = 256;
1209     if (argc > 3)
1210         threadsPerBlock = atoi(argv[3]);
1211
1212     int blocksPerGrid = (numElements + threadsPerBlock - 1) / threadsPerBlock;

```

```

1214     if(argc > 4)
1215         blocksPerGrid = atoi(argv[4]);
1216
1217     GpuTimer timer;
1218     timer.Start();
1219
1220     if(cpu_only)
1221         run_on_cpu(size, numElements, h_A, h_B, h_C);
1222     else
1223         run_on_gpu(size, numElements, h_A, h_B, h_C, threadsPerBlock, blocksPerGrid);
1224
1225     timer.Stop();
1226     #if (VERBOSE == true)
1227         printf("Elapsed: %f\n", timer.Elapsed());
1228     #endif
1229
1230     // Verify that the result vector is correct
1231     for (int i = 0; i < numElements; ++i)
1232     {
1233         if (fabs(h_A[i] + h_B[i] - h_C[i]) > 1e-5) {
1234             fprintf(stderr, "Result verification failed at element %d!\n", i);
1235             exit(EXIT_FAILURE);
1236         }
1237     }
1238
1239     #if (VERBOSE == true)
1240         printf("Test PASSED\n");
1241     #endif
1242
1243     // Free host memory
1244     free(h_A);
1245     free(h_B);
1246     free(h_C);
1247
1248     #if (VERBOSE == true)
1249         printf("Done\n");
1250     #endif
1251
1252     // printf("elements, time, type, threads_per_block, blocks_per_grid\n");
1253     printf("%d, %f, %s, %d, %d\n", numElements, timer.Elapsed(), (cpu_only ? "cpu" : "gpu"),
1254           threadsPerBlock, blocksPerGrid);
1255     return 0;
1256 }

```

../0.Simple/vectorAdd/vectorAdd.cu