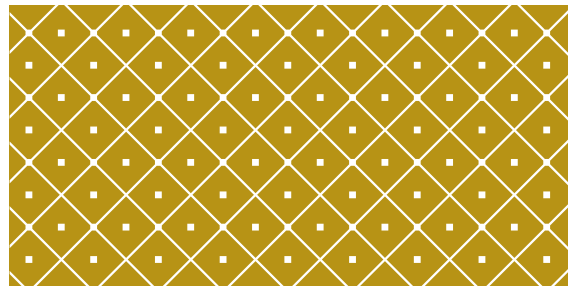


PROGRAMOWANIE WEKTOROWE I RÓWNOLEGŁE

Wykład dla kierunku: Matematyka stosowana i technologie informatyczne



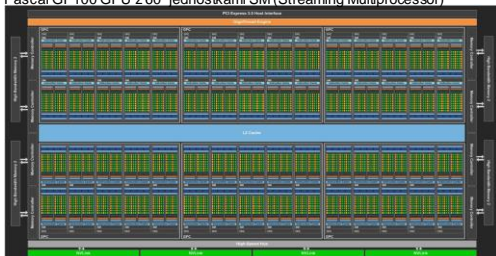
RODZAJE PAMIĘCI I OPTYZMALIZACJA JEJ WYKORZYSTANIA.

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGŁE 2

ARCHITEKTURA PASCAL – GTX10X0

Pascal GP100 GPU z 60 jednostkami SM (Streaming Multiprocessor)



(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGŁE 3

ARCHITEKTURA PASCAL – GTX10X0

- Układ zawiera:
- 6 jednostek Graphics Processing Cluster (GPC) po 10 SM
 - 4 MB cache L2
- Moduł SM zawiera:
- 64 rdzenie pojedynczej precyzji
 - 32 rdzenie podwójnej precyzji
 - 16 jednostek specjalne,
 - 16 jednostek odczytująco-zapisujących
 - 64 KB pamięci współdzielonej.



(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGŁE 4

ARCHITEKTURA PASCAL – GTX10X0

Porównanie Kepler, Maxwell i Pascal GPU

GPU	Kepler GK110	Maxwell GM200	Pascal GP100
Compute Capability	3.5	5.2	6.0
Threads / Warp	32	32	32
Max Warps / Multiprocessor	64	64	64
Max Threads / Multiprocessor	2048	2048	2048
Max Thread Blocks / Multiprocessor	16	32	32
Max 32-bit Registers / SM	65536	65536	65536
Max Registers / Block	65536	32768	65536
Max Registers / Thread	255	255	255
Max Thread Block Size	1024	1024	1024
Shared Memory Size / SM	16 KB/32 KB/48 KB	96 KB	64 KB

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGŁE 5

ARCHITEKTURA AMPERE – RTX30X0

- Układ zawiera:
- 7 jednostek Graphics Processing Cluster (GPC) po 12 SM
 - 6 MB cache L2
- Moduł SM zawiera:
- 2 rdzenie podwójnej precyzji
 - 128 rdzeni pojedynczej precyzji
 - w tym 64 rdzenie całkowite
 - 4 jednostki specjalne,
 - 16 jednostek odczytująco-zapisujących
 - 4 rdzenie Tensor
 - 1 jednostkę śledzenia promieni
 - 128 KB pamięci współdzielonej/L1.



(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGŁE 6

STRUKTURA CUDADEVICEPROP - WYCIĄG

```

struct cudaDeviceProp {
    char name[256];
    size_t totalGlobalMem;
    size_t totalConstMem;
    size_t sharedMemPerBlock;
    size_t
sharedMemPerMultiProcessor;
    int L2CacheSize;
    int regsPerBlock;
    int regsPerMultiProcessor;
    int warpSize;
    size_t memPitch;
    int maxThreadsPerBlock;
    int maxThreadsPerMultiProcessor;
    int maxBlocksPerMultiProcessor;
    int maxThreadsDim[3];
    int maxGridSize[3];
    int major;
    int minor;

    int clockRate;
    size_t textureAlignment;
    int deviceOverlap;
    int multiProcessorCount;
    int kernelExecTimeoutEnabled;
    int integrated;
    int canMapHostMemory;
    int computeMode;
    int maxTexture1D;
    int maxTexture2D[2];
    int maxTexture3D[3];
    int maxTexture2DArray[3];
    int concurrentKernels;
}

```

(C) KISI d.KIK PCz 2023

PROGRAMOWANIEWEKTOROWEIRÓWNOLEGLE 7

WYDOBYCIE INFORMACJIO URZĄDZENIU

```

25 HANDLE_ERROR( cudaGetDeviceProperties( &prop, i ) );
26 printf( " --- General information for device %d ---\n", i );
27 printf( "Name: %s\n", prop.name );
28 printf( "Potencjal obliczeniowy: %d.%d\n", prop.major, prop.minor );
29 printf( "Zegar: %d\n", prop.clockRate );
30 printf( "Wartość deviceOverlap: " );
31 if (prop.deviceOverlap)
32     printf( "Włączona" );
33 else
34     printf( "Wyłączona" );
35 printf( "Limit czasu działania jadra: " );
36 if (prop.kernelExecTimeoutEnabled)
37     printf( "Włączony" );
38 else
39     printf( "Wyłączony" );

```

(C) KISI d.KIK PCz 2023

PROGRAMOWANIEWEKTOROWEIRÓWNOLEGLE 8

WYDOBYCIE INFORMACJIO URZĄDZENIU

```

41 printf( " --- Informacje o pamięci urządzenia %d ---\n", i );
42 printf( "Ilość pamięci globalnej: %d\n", prop.totalGlobalMem );
43 printf( "Ilość pamięci stałej: %d\n", prop.totalConstMem );
44 printf( "Maks. szerokość pamięci: %d\n", prop.memPitch );
45 printf( "Wyrównanie tekstur: %d\n", prop.textureAlignment );
46 printf( " --- Informacje na temat wieloprocesorów urządzenia %d ---\n", i );
47 printf( "Liczba wieloprocesorów: %d\n",
48         prop.multiProcessorCount );
49 printf( "Pamięć wspólna na wieloprocesor: %d\n", prop.sharedMemPerBlock );
50 printf( "Rejestry na wieloprocesor: %d\n", prop.regsPerBlock );
51 printf( "Liczba wątków w osnowie: %d\n", prop.warpSize );
52 printf( "Maks. liczba wątków na blok: %d\n",
53         prop.maxThreadsPerBlock );
54 printf( "Maks. liczba wymiarów wątków: (%d, %d, %d)\n",
55         prop.maxThreadsDim[0], prop.maxThreadsDim[1],
56         prop.maxThreadsDim[2] );
57 printf( "Maks. liczba wymiarów siatki: (%d, %d, %d)\n",
58         prop.maxGridSize[0], prop.maxGridSize[1],
59         prop.maxGridSize[2] );

```

(C) KISI d.KIK PCz 2023

PROGRAMOWANIEWEKTOROWEIRÓWNOLEGLE 9

Wybrane właściwości

char name[256] - łańcuch znaków ASCII identyfikujący urządzenie;

size_t totalGlobalMem - globalna pamięć dostępna na urządzeniu w bajtach;

size_t totalConstMem - pamięć stała dostępna na urządzeniu w bajtach;

size_t sharedMemPerBlock - maksymalny rozmiar pamięci wspólnej na blok w bajtach;

size_t sharedMemPerMultiProcessor - maksymalny rozmiar pamięci wspólnej na multiprocessor w bajtach;

int L2CacheSize - rozmiar w bajtach pamięci podręcznej L2;

(C) KISI d.KIK PCz 2023

PROGRAMOWANIEWEKTOROWEIRÓWNOLEGLE 10

Wybrane właściwości

int multiProcessorCount - liczba multiprocessorów w urządzeniu;

int regsPerBlock - maksymalna liczba rejestrów na blok;

int regsPerMultiProcessor - maksymalna liczba rejestrów na multiprocessor;

int maxThreadsPerBlock - maksymalna liczba wątków na blok;

int maxThreadsPerMultiProcessor - maksymalna liczba wątków na multiprocessor;

int maxBlocksPerMultiProcessor - maksymalna liczba bloków na multiprocessor;

int maxThreadsDim[3] - maksymalny rozmiar każdego wymiaru bloku;

int maxGridSize[3] - maksymalny rozmiar każdego wymiaru siatki;

int warpSize - liczba wątków w osnowie;

(C) KISI d.KIK PCz 2023

PROGRAMOWANIEWEKTOROWEIRÓWNOLEGLE 11

Wybrane właściwości

size_t textureAlignment - wymagania dotyczące wyrównania tekstur;

int maxTexture1D - maksymalny rozmiar tekstur 1D;

int maxTexture2D[2] - maksymalny rozmiar tekstur 2D;

int maxTexture3D[3] - maksymalny rozmiar tekstur 3D;

int maxTexture2DArray[3] - maksymalny rozmiar tablic tekstur 2D;

(C) KISI d.KIK PCz 2023

PROGRAMOWANIEWEKTOROWEIRÓWNOLEGLE 12

Wybrane właściwości

```

int concurrentKernels - możliwość uruchomienia wielu
jąderek;
int major - zdolność obliczeniowa wartość główna;
int minor - zdolność obliczeniowa wartość po przecinku;
int clockRate - częstotliwość zegara w kHz;
int deviceOverlap - możliwość jednoczesnego wykonywania
jądra i kopiowania;
int kernelExecTimeoutEnabled - ograniczenie czasowe
wykonywania jądra;
int integrated - czy układ jest zintegrowany;
int canMapHostMemory - czy można pamięć hosta na
urządzenie;

```

(C) KISI d.KIK PCz 2023 PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 13

Dodawanie wektorów – 1 – bloki

```

21 | global void add( int *a, int *b, int *c ) {
22 |     int tid = blockIdx.x; // Ten wątek przetwarza dane pod określonym indeksem
23 |     if (tid < N)
24 |         c[tid] = a[tid] + b[tid];
25 | }

```

48 | add<<<N,1>>>(dev_a, dev_b, dev_c);

(C) KISI d.KIK PCz 2023 PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 14

Dodawanie wektorów – 2 – bloki

```

21 | global void add( int *a, int *b, int *c ) {
22 |     int tid = blockIdx.x;
23 |     while (tid < N) {
24 |         c[tid] = a[tid] + b[tid];
25 |         tid += blockDim.x;
26 |     }
27 | }

```

55 | add<<<128,1>>>(dev_a, dev_b, dev_c);

(C) KISI d.KIK PCz 2023 PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 15

Dodawanie wektorów – 3 – wątki

```

21 | global void add( int *a, int *b, int *c ) {
22 |     int tid = threadIdx.x;
23 |     if (tid < N)
24 |         c[tid] = a[tid] + b[tid];
25 | }

```

48 | add<<<1,N>>>(dev_a, dev_b, dev_c);

(C) KISI d.KIK PCz 2023 PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 16

Dodawanie wektorów – 4 – wątki + bloki

Block 0	Thread 0	Thread 1	Thread 2	Thread 3
Block 1	Thread 0	Thread 1	Thread 2	Thread 3
Block 2	Thread 0	Thread 1	Thread 2	Thread 3
Block 3	Thread 0	Thread 1	Thread 2	Thread 3

(C) KISI d.KIK PCz 2023 PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 17

Dodawanie wektorów – 4 – wątki + bloki

```

21 | global void add( int *a, int *b, int *c ) {
22 |     int tid = threadIdx.x + blockIdx.x * blockDim.x;
23 |     if (tid < N)
24 |         c[tid] = a[tid] + b[tid];
25 | }

```

```

55 | add<<<(N+127)/128,128>>>( dev_a, dev_b, dev_c );
56 | //lub
57 | add<<<128,128>>>( dev_a, dev_b, dev_c );

```

(C) KISI d.KIK PCz 2023 PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 18

Dodawanie wektorów – 5 – wątki + bloki

```

21 | global __void add( int *a, int *b, int *c ) {
22 |     int tid = threadIdx.x + blockIdx.x * blockDim.x;
23 |     while (tid < N) {
24 |         c[tid] = a[tid] + b[tid];
25 |         tid += blockDim.x * gridDim.x;
26 |     }
27 | }

```

```

55 | add<<<128,128>>>( dev_a, dev_b, dev_c );

```

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 19

Klasyfikatory zmiennych w przestrzeni pamięci

Klasyfikatory zmiennych w przestrzeni pamięci określają lokalizację pamięci na urządzeniu, w której znajduje się zmienna.

Zmienna automatycznie zadeklarowana w kodzie urządzenia bez żadnego klasyfikatora przestrzeni pamięci `__device__`, `__shared__` lub `__constant__` zazwyczaj rezyduje w rejestrze.

Jednak w niektórych przypadkach kompilator może zdecydować się na umieszczenie go w pamięci lokalnej, co może mieć niekorzystny wpływ na wydajność.

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 20

Klasyfikator `__device__`

Klasyfikator przestrzeni pamięci `__device__` deklaruje zmienną rezydującą na urządzeniu. Co najwyżej jeden z innych klasyfikatorów przestrzeni pamięci może być użyty razem z `__device__` w celu dalszego określenia, do której przestrzeni pamięci należy zmienna. Jeśli żaden z nich nie jest obecny, zmienna:

- Rezyduje w globalnej przestrzeni pamięci,
- Posiada czas życia kontekstu CUDA, w którym jest tworzona,
- Ma odrębny obiekt na urządzenie,
- Jest dostępna ze wszystkich wątków w siatce oraz z hosta poprzez bibliotekę wykonawczą (`cudaGetSymbolAddress()` / `cudaGetSymbolSize()` / `cudaMemcpyToSymbol()` / `cudaMemcpyFromSymbol()`).

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 21

Klasyfikator `__shared__`

Klasyfikator przestrzeni pamięci `__shared__`, opcjonalnie używany razem z `__device__`, deklaruje zmienną, która:

- Rezyduje we współdzielonej przestrzeni pamięci bloku wątków,
- Ma żywotność bloku,
- Ma odrębny obiekt na blok,
- Jest dostępny tylko ze wszystkich wątków w bloku,
- Nie ma stałego adresu.

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 22

Klasyfikator `__constant__`

Klasyfikator przestrzeni pamięci `__constant__`, opcjonalnie używany razem z `__device__`, deklaruje zmienną, która:

- Rezyduje w przestrzeni pamięci stałych,
- Posiada czas życia kontekstu CUDA, w którym jest tworzona,
- Ma odrębny obiekt na urządzenie,
- Jest dostępna ze wszystkich wątków w siatce oraz z hosta poprzez bibliotekę wykonawczą (`cudaGetSymbolAddress()` / `cudaGetSymbolSize()` / `cudaMemcpyToSymbol()` / `cudaMemcpyFromSymbol()`).

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 23

Klasyfikator `__grid_constant__`

Klasyfikator `__grid_constant__` dla architektury obliczeniowych wejściowych lub równych 7.0 zawiera adnotację do parametru funkcji `__global__` z kwalifikacją `const` typu innego niż referencyjny, który:

- Ma żywotność sieci,
 - Jest prywatny dla siatki, tj. obiekt nie jest dostępny dla wątków hosta i wątków z innych siatek, w tym podsiatek,
 - Ma odrębny obiekt na siatkę, tj. Wszystkie wątki w siatce widzą ten sam adres,
 - Jest tylko do odczytu, tj. modyfikacja obiektu `__grid_constant__` lub któregośkolwiek z jego obiektów podrzędnych jest niezdefiniowanym zachowaniem.
- Wymagania:
- Parametry jądra z adnotacją `__grid_constant__` muszą mieć kwalifikowane przez `const` typy niereferencyjne.
 - Wszystkie deklaracje funkcji muszą być zgodne w odniesieniu do wszelkich parametrów `__grid_constant__`.
 - Specjalizacja szablonu funkcji musi być zgodna z deklaracją szablonu podstawowego w odniesieniu do wszelkich parametrów `__grid_constant__`.
 - Dyrektywa tworzenia instancji szablonu funkcji musi być zgodna z deklaracją szablonu podstawowego w odniesieniu do wszelkich parametrów `__grid_constant__`.

Jeśli zostanie pobrany adres parametru funkcji `__global__`, kompilator zwykle utworzy kopie parametru jądra w pamięci lokalnej wątku i użyje adresu kopii, aby częściowo obsługiwać semantykę C++, która pozwala każdemu wątkowi modyfikować własną kopię lokalnej parametru funkcji. Adnotacja parametru funkcji `__global__` za pomocą `__grid_constant__` gwarantuje, że kompilator nie utworzy kopii parametru jądra w pamięci lokalnej wątku, ale zamiast tego użyje ogólnego adresu samego parametru. Unikanie kopii lokalnej może poprawić wydajność.

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 24

Klasyfikator __managed__

Klasyfikator __managed__ przestrzeni pamięci, opcjonalnie używany razem z __device__, deklaruje zmienną, która:

- Można się do niej odwoływać zarówno z kodu urządzenia, jak i hosta, na przykład jej adres można pobrać lub można go odczytać lub zapisać bezpośrednio z urządzenia lub funkcji hosta.
- Ma czas życia aplikacji.

Iloczyn skalarny

Iloczyn skalarny:

$$c = \sum_{i=1}^N a_i b_i$$

Iloczyn skalarny

```

17 #include "../common/book.h"
18
19 #define imin(a,b) (a<b?a:b)
20
21 const int N = 33 * 1024;
22 const int threadsPerBlock = 256;
23 const int blocksPerGrid =
24     imin( 32, (N+threadsPerBlock-1) / threadsPerBlock );
25
26
27 global void dot( float *a, float *b, float *c ) {
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59 int main( void ) {

```

Iloczyn skalarny

```

39 int main( void ) {
40     float *a, *b, *c, *partial_c;
41     float *dev_a, *dev_b, *dev_partial_c;
42
43     // Alokacja pamięci na CPU
44     a = (float*)malloc( N*sizeof(float) );
45     b = (float*)malloc( N*sizeof(float) );
46     partial_c = (float*)malloc( blocksPerGrid*sizeof(float) );
47
48     // Alokacja pamięci na GPU
49     HANDLE_ERROR( cudaMalloc( (void**)&dev_a,
50                             N*sizeof(float) ) );
51     HANDLE_ERROR( cudaMalloc( (void**)&dev_b,
52                             N*sizeof(float) ) );
53     HANDLE_ERROR( cudaMalloc( (void**)&dev_partial_c,
54                             blocksPerGrid*sizeof(float) ) );
55
56     // Zapełnienie pamięci hosta danymi
57     for( int i=0; i<N; i++) {
58         a[i] = i;
59         b[i] = i*i;
60     }
61
62     // Skopiowanie tablic a i b do GPU
63     HANDLE_ERROR( cudaMemcpy( dev_a, a, N*sizeof(float),
64                               cudaMemcpyHostToDevice ) );
65     HANDLE_ERROR( cudaMemcpy( dev_b, b, N*sizeof(float),
66                               cudaMemcpyHostToDevice ) );
67
68     dot<<<blocksPerGrid, threadsPerBlock>>>( dev_a, dev_b,
69                                             dev_partial_c );

```

Iloczyn skalarny

```

88 dot<<<blocksPerGrid, threadsPerBlock>>>( dev_a, dev_b,
89                                           dev_partial_c );
90
91 // Skopiowanie tablicy c z powrotem z GPU do CPU
92 HANDLE_ERROR( cudaMemcpy( partial_c, dev_partial_c,
93                           blocksPerGrid*sizeof(float),
94                           cudaMemcpyDeviceToHost ) );
95
96 // Zakoczenie pracy na CPU
97 c = 0;
98 for( int i=0; i<blocksPerGrid; i++) {
99     c += partial_c[i];
100 }
101
102 #define sum_squares(x) (x*(x+1)*(2*x+1))/6
103 printf( "Dose GPU value %g = %g\n", c,
104        sum_squares( (float)(N - 1) ) );
105
106 // Zwolnienie pamięci na GPU
107 HANDLE_ERROR( cudaFree( dev_a ) );
108 HANDLE_ERROR( cudaFree( dev_b ) );
109 HANDLE_ERROR( cudaFree( dev_partial_c ) );
110
111 // Zwolnienie pamięci na CPU
112 free( a );
113 free( b );
114 free( partial_c );
115 }

```

Iloczyn skalarny

```

27 global void dot( float *a, float *b, float *c ) {
28     __shared__ float cache[threadsPerBlock];
29     int tid = threadIdx.x * blockDim.x + blockIdx.x;
30     int cacheIndex = threadIdx.x;
31
32     float temp = 0;
33     while( tid < N ) {
34         temp += a[tid] * b[tid];
35         tid += blockDim.x * gridDim.x;
36     }
37
38     // Zastąpienie wartości pamięci podręcznej
39     cache[cacheIndex] = temp;
40
41     // Synchronizacja wątków w tym bloku
42     __syncthreads();
43
44     // W przypadku redukcji threadsPerBlock musi być potęgą 2.
45     // ze względu na poniższy kod
46     int i = blockDim.x/2;
47     while( i != 1 ) {
48         if( cacheIndex < i )
49             cache[cacheIndex] += cache[cacheIndex + i];
50         __syncthreads();
51         i /= 2;
52     }
53
54     if( cacheIndex == 0 )
55         o[blockDim.x] = cache[0];
56 }

```

Funkcja __syncthreads()

void __syncthreads();

czeka, aż wszystkie wątki w bloku wątków osiągną ten punkt, a wszystkie dostępy do pamięci globalnej i współdzielonej dokonane przez te wątki przed __syncthreads() będą widoczne dla wszystkich wątków w bloku.

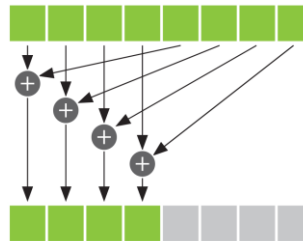
__syncthreads() służy do koordynowania komunikacji między wątkami tego samego bloku. Gdy niektóre wątki w bloku uzyskują dostęp do tych samych adresów w pamięci współdzielonej lub globalnej, w przypadku niektórych z tych dostępu do pamięci istnieją potencjalne zagrożenia związane z odczytem po zapisie, zapisem po odczycie lub zapisem po zapisie. Tych zagrożeń związanych z danymi można uniknąć, synchronizując wątki pomiędzy tymi dostęпами.

__syncthreads() jest dozwolone w kodzie warunkowym, ale tylko wtedy, gdy warunek daje identyczną wartość w całym bloku wątków, w przeciwnym razie wykonanie kodu może się zawiesić lub wywołać niezamierzone skutki uboczne.

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 31

Iloczyn skalarny - redukcja



(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 32

Iloczyn skalarny

```

58 __global__ void dot21(float *a, float *b, float *c) {
59     __shared__ float cache[threadsPerBlock];
60     int tid = threadIdx.x + blockIdx.x * blockDim.x;
61     int cacheIndex = threadIdx.x;
62
63     float temp = 0;
64     while (tid < N) {
65         temp += a[tid] * b[tid];
66         tid += blockDim.x * gridDim.x;
67     }
68
69     // Ustawienie wartości pamięci podręcznej
70     cache[cacheIndex] = temp;
71
72     // Synchronizacja wątków w tym bloku
73     __syncthreads();
74
75     // Redukcja - threadsPerBlock nie musi być potęgą 2!!!
76     int i = blockDim.x;
77     int cor = 1 + blockDim.x;
78     while (i != 0) {
79         if (cacheIndex < i)
80             cache[cacheIndex] += cache[cacheIndex + cor];
81         __syncthreads();
82         i = cor/2;
83         cor = i + cor;
84     }
85     if (cacheIndex == 0)
86         c[blockIdx.x] = cache[0];
87 }

```

(C) KISI d.KIK PCz 2023

PROGRAMOWANIE WEKTOROWE I RÓWNOLEGLE 33