

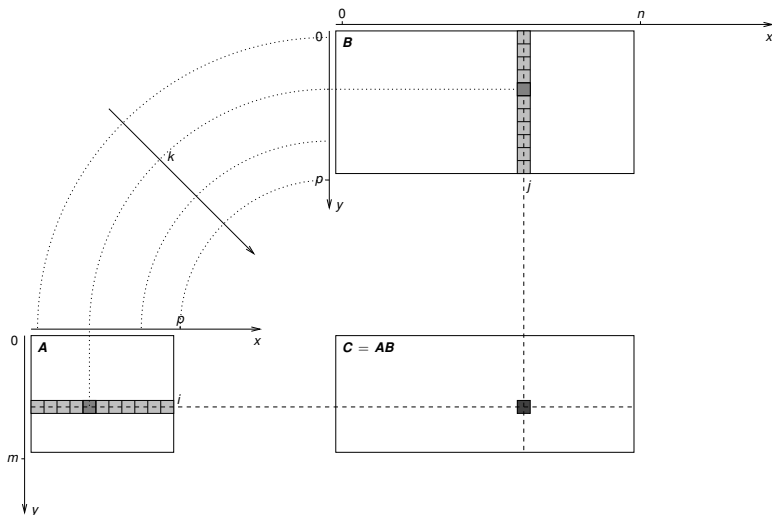
Mnożenie macierzy rzeczywistych

- Problem: $\mathbf{C}_{m \times n} = \mathbf{A}_{m \times p} \cdot \mathbf{B}_{p \times n}$, gdzie $a_{ik}, b_{kj}, c_{ij} \in \mathbb{R}$.
- Każdy element macierzy wynikowej \mathbf{C} jest **iloczynem skalarnym** wiersza \mathbf{A} przez kolumnę \mathbf{B} (indeksy od 0):

$$c_{ij} = \sum_{k=0}^{p-1} a_{ik} b_{kj}, \quad i = 0, \dots, m-1, \quad j = 0, \dots, n-1$$

- Czy warto się tym problemem zajmować? Czy potrafimy dostarczyć procesorowi GPU dostatecznie dużo pracy?
 Złożoność obliczeniowa (\cdot i $+$) $mnp + mn(p-1)$ [FLOP]
 Liczba dostępow do pamięci $mn + mp + pn$ [T]
- Współczynnik pamięciowy złożoności obliczeniowej CGMA (*Compute to Global Memory Access ratio*) dla $m = p = n$ wynosi $\frac{2}{3}n - \frac{1}{3}$ i dla $n \rightarrow \infty$ może być dostatecznie duży.
- Uwaga! Dla algorytmu iloczynu skalarnego współczynnik CGMA wynosi już tylko $\frac{p+(p-1)}{1+p+p} = \frac{2p-1}{2p+1} \approx 1$.

Podstawowy schemat mnożenia macierzy



Generacja macierzy w MATLAB-ie – `matgen.m`

```
1  function matgen(m, p, n)
2
3  %% utworzenie losowych macierzy czynników i iloczynu
4  A = randn(m, p);
5  B = randn(p, n);
6  C = A * B;
7
8  %% zapamiętanie rozmiarów i macierzy w pliku binarnym
9  f = fopen('matmul.dat', 'wb');
10 fwrite(f, [m, p, n], 'int');
11 % zapamiętujemy transponowane macierze, bo język C
12 % przechowuje macierze wierszami, a nie kolumnami
13 fwrite(f, A.', 'single');
14 fwrite(f, B.', 'single');
15 fwrite(f, C.', 'single');
16 fclose(f);
```

Szkielet mnożenia macierzy w języku C – matmul.c

```
1  #include <stdlib.h>
2  #include <stdio.h>
3
4  #include "matmul.h"
5
6  int main(int argc, char *argv[])
7  {
8      app_timer_t start, stop;
9      int         m,  p,  n;
10     float        *A, *B, *C, *D;
11     read_data(&m, &p, &n, &A, &B, &C, &D);
12     timer(&start);
13     matmul(C, A, B, m, p, n, argc, argv);
14     timer(&stop);
15     elapsed_time(start, stop, 2 * m * p * n);
16     matcmp(C, D, m, n);
17     free_mem(A, B, C, D);
18     if (IsDebuggerPresent()) getchar();
19     return 0;
20 }
```

Szkielet mnożenia macierzy w języku C – matmul.c

```
1 void read_data(int *m_ptr, int *p_ptr, int *n_ptr,
2               float **A_ptr, float **B_ptr,
3               float **C_ptr, float **D_ptr)
4 {
5     FILE *f = fopen("matmul.dat", "rb");
6
7     fread(m_ptr, sizeof(int), 1, f);
8     fread(p_ptr, sizeof(int), 1, f);
9     fread(n_ptr, sizeof(int), 1, f);
10
11     alloc_mem(*m_ptr, *p_ptr, *n_ptr,
12             A_ptr, B_ptr, C_ptr, D_ptr);
13
14     fread(*A_ptr, sizeof(float), *m_ptr * *p_ptr, f);
15     fread(*B_ptr, sizeof(float), *p_ptr * *n_ptr, f);
16     fread(*D_ptr, sizeof(float), *m_ptr * *n_ptr, f);
17
18     fclose(f);
19 }
```

Szkielet mnożenia macierzy w języku C – `matmul.c`

```
1 void alloc_mem(int m, int p, int n,  
2               float **A_ptr, float **B_ptr,  
3               float **C_ptr, float **D_ptr)  
4 {  
5     *A_ptr = (float *) malloc(m * p * sizeof(float));  
6     *B_ptr = (float *) malloc(p * n * sizeof(float));  
7     *C_ptr = (float *) malloc(m * n * sizeof(float));  
8     *D_ptr = (float *) malloc(m * n * sizeof(float));  
9 }  
10  
11 void free_mem(float *A, float *B, float *C, float *D)  
12 {  
13     free(A);  
14     free(B);  
15     free(C);  
16     free(D);  
17 }
```

Szkielet mnożenia macierzy w języku C – matmul.c

```
1  #define WINDOWS_LEAN_AND_MEAN
2  #include <windows.h>
3
4  typedef LARGE_INTEGER app_timer_t;
5
6  #define timer(t_ptr) QueryPerformanceCounter(t_ptr)
7
8  void elapsed_time(app_timer_t start, app_timer_t stop,
9                  unsigned long flop)
10 {
11     double etime;
12     LARGE_INTEGER clk_freq;
13     QueryPerformanceFrequency(&clk_freq);
14     etime = (stop.QuadPart - start.QuadPart) /
15             (double) clk_freq.QuadPart;
16     printf("CPU_(total!)_time_=%_.3f_ms_(%6.3f_GFLOP/s)\n",
17           etime * 1e3, 1e-9 * flop / etime);
18 }
```

Szkielet mnożenia macierzy w języku C – `matmul.c`

```
1  #include <math.h>
2
3  void matcmp(float *C, float *D, int m, int n)
4  {
5      int k;
6      float d, e = -1.0f;
7      for (k = 0; k < m * n; k++)
8          if ((d = fabsf(C[k] - D[k])) > e)
9              e = d;
10     printf("max. abs. err. = %.1e\n", e);
11 }
```

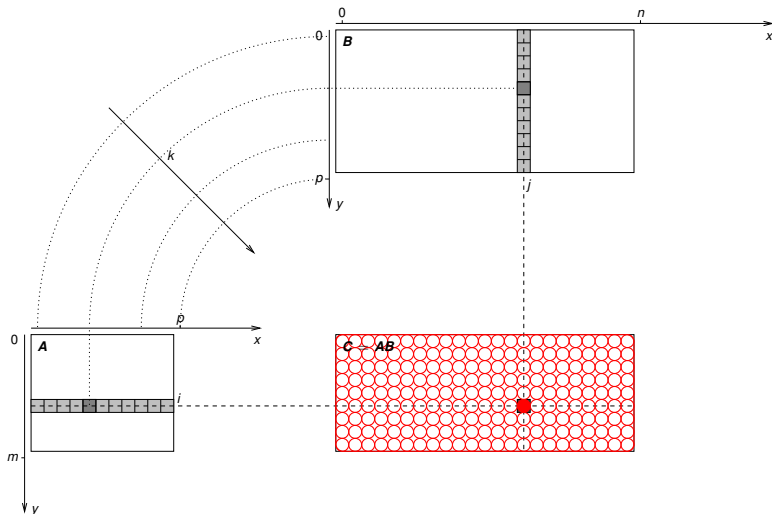
Sprawdzanie poprawności

Wydajność programu nie ma żadnego znaczenia, jeżeli daje on niepoprawne wyniki!

Jądro mnożenia macierzy w języku C – matmul0.c

```
1  #include "matmul.h"
2
3  void matmul(float *C, float *A, float *B,
4             int m, int p, int n, ...)
5  {
6      int i, j;
7      for (i = 0; i < m; i++)
8          for (j = 0; j < n; j++)
9              matmul_kernel(i, j, C, A, B, m, p, n);
10 }
11
12 static void matmul_kernel(int i, int j,
13                            float *C, float *A, float *B, int m, int p, int n)
14 {
15     int k;
16     float s = 0;
17     for (k = 0; k < p; k++)
18         s += A[i*p + k] * B[k*n + j];
19     C[i*n + j] = s;
20 }
```

Podstawowy schemat mnożenia macierzy – „ziarna”



Mnożenie macierzy w języku CUDA C – matmul1.cu

```
1 #include "helper_cuda.h" // -I$(NVCUDASAMPLES_ROOT)/common/inc
2 #include "matmul.h"
3
4 void matmul(float *C, float *A, float *B,
5             int m, int p, int n, ...)
6 {
7     // 1. Wybór urządzenia CUDA dla kolejnych funkcji API
8     checkCudaErrors(cudaSetDevice(0));
9
10    // 2. Alokacja pamięci urządzenia
11    float *dev_A, *dev_B, *dev_C;
12    checkCudaErrors(cudaMalloc(&dev_A, m*p*sizeof(float)));
13    checkCudaErrors(cudaMalloc(&dev_B, p*n*sizeof(float)));
14    checkCudaErrors(cudaMalloc(&dev_C, m*n*sizeof(float)));
15
16    // 3. Kopiowanie danych gospodarza do urządzenia
17    checkCudaErrors(cudaMemcpy(dev_A, A, m*p*sizeof(float),
18                                cudaMemcpyHostToDevice));
19    checkCudaErrors(cudaMemcpy(dev_B, B, p*n*sizeof(float),
20                                cudaMemcpyHostToDevice));
```

Mnożenie macierzy w języku CUDA C – matmul1.cu

```
1 // 4. Asynchroniczne uruchomienie jądra
2 dim3 dimBlock(n, m);
3
4 cudaEvent_t start, stop; // pomiar czasu wykonania jądra
5 checkCudaErrors(cudaEventCreate(&start));
6 checkCudaErrors(cudaEventCreate(&stop));
7 checkCudaErrors(cudaEventRecord(start, 0));
8
9     matmul_kernel<<<1, dimBlock>>>
10         (dev_C, dev_A, dev_B, m, p, n);
11
12 checkCudaErrors(cudaGetLastError());
13
14 checkCudaErrors(cudaEventRecord(stop, 0));
15 checkCudaErrors(cudaEventSynchronize(stop));
16 float elapsedTime;
17 checkCudaErrors(cudaEventElapsedTime(&elapsedTime,
18                                     start, stop));
19 checkCudaErrors(cudaEventDestroy(start));
20 checkCudaErrors(cudaEventDestroy(stop));
```

Mnożenie macierzy w języku CUDA C – matmul1.cu

```
1 // 5. Synchronizacja z jądrem (bariera)
2 checkCudaErrors(cudaDeviceSynchronize());
3
4 // 6. Kopiowanie wyników pracy z urządzenia do gospodarza
5 checkCudaErrors(cudaMemcpy(C, dev_C, m*n*sizeof(float),
6                             cudaMemcpyDeviceToHost));
7
8 // 7. Zwolnienie pamięci urządzenia
9 checkCudaErrors(cudaFree(dev_C));
10 checkCudaErrors(cudaFree(dev_B));
11 checkCudaErrors(cudaFree(dev_A));
12
13 checkCudaErrors(cudaDeviceReset()); // dla debuggera
14
15 printf("GPU_(kernel)_time_=%.3f_ms_(%6.3f_GFLOP/s)\n",
16        elapsedTime, 2e-6 * m * n * p / elapsedTime);
17 }
```

Jądro mnożenia macierzy w CUDA C – matmul1.cu

```
1  __global__
2  static void matmul_kernel(
3      float *C, float *A, float *B, int m, int p, int n)
4  {
5      int i = threadIdx.y;
6      int j = threadIdx.x;
7      int k;
8      float s = 0;
9      for (k = 0; k < p; k++)
10         s += A[i*p + k] * B[k*n + j];
11     C[i*n + j] = s;
12 }
```

Położenie „ziarna”, nad którym pracuje dany wątek. . .

. . . ustalamy na podstawie współrzędnych wątku w bloku i współrzędnych bloku w sieci.

Porównanie wydajności implementacji CPU i GPU

- Na początek prosty test dla $m = 13$, $p = 24$, $n = 35$.

- CPU (AMD Ryzen 9 7900)

CPU (total!) time = 0.005 ms (4.325 GFLOP/s)
max. abs. err. = 2.9e-06

- GPU (GeForce RTX 4070: 56 SM, 7168 rdzeni, $C_c = 8.9$)

GPU (kernel) time = 0.024 ms (0.927 GFLOP/s)
CPU (total!) time = 52.35 ms (0.000 GFLOP/s)

- Obie implementacje dały prawidłowe wyniki.
- Nawet czas wykonania samego jądra GPU (nie mówiąc o narzutach na transfer danych i uruchomienie jądra) jest znacznie dłuższy niż całkowity czas dla CPU.
- Wiadomo, że implementacja GPU nie jest efektywna – wykorzystuje tylko jeden blok, więc tylko jeden SM.
- Czasy są bardzo krótkie – zwiększymy rozmiar problemu, aby zmniejszyć narzuty i zwiększyć dokładność pomiaru.

Porównanie wydajności implementacji CPU i GPU

- Do dalszych testów przyjmujemy $m = p = n = 1021$.

- CPU

```
CPU (total!) time = 572.073 ms (3.72 GFLOP/s)
max. abs. err. = 2.1e-04
```

- GPU (GeForce RTX 4070)

```
CUDA error at ../matmul1.cu:52
code=9(cudaErrorInvalidConfiguration)
"cudaGetLastError() "
```

Błąd oznacza nieprawidłową konfigurację wykonania (*execution configuration*), opisywaną składnią `<<<...>>>`.

- Przyczyna niepowodzenia:

$$T_B = mn = 1021^2 > T_{B \max} = 1024$$

- Podzielimy wątki pomiędzy wiele bloków, odpowiadających „kafelkom” (*tile*) o rozmiarach $M \times N$ macierzy \mathbf{C} .

Parametry „geometryczne” CUDA dla różnych C_c

Parametr	cudaDeviceGetAttribute →cudaDevAttrMax...	1.x	2.x	3.x	5.x 6.x	7.x 8.x
$B_{G,x \text{ max}}$...GridDimX	$2^{16} - 1$		$2^{31} - 1$		
$B_{G,y \text{ max}}$...GridDimY	$2^{16} - 1$				
$B_{G,z \text{ max}}$...GridDimZ	1	$2^{16} - 1$			
$T_{B,x \text{ max}}$...BlockDimX	512	1024			
$T_{B,y \text{ max}}$...BlockDimY	512	1024			
$T_{B,z \text{ max}}$...BlockDimZ	64				
$T_B \text{ max}$...ThreadsPerBlock	512	1024			

Legenda (J_K oznacza liczbę J -ów na jeden K):

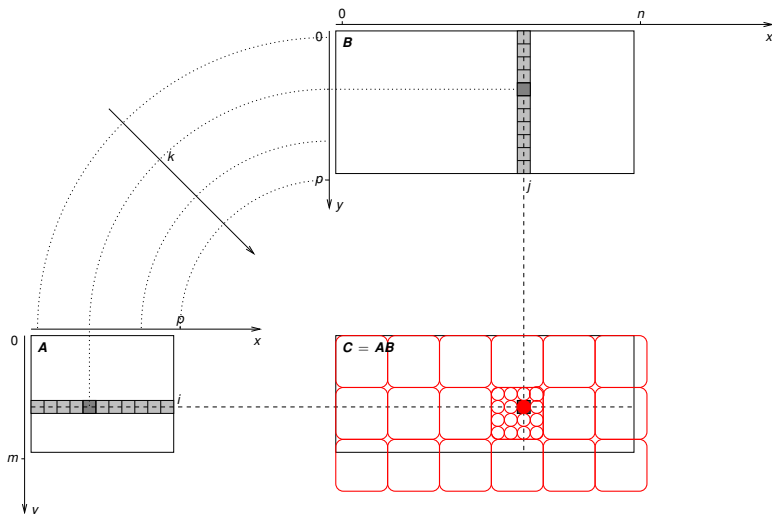
G sieć bloków (*grid*)

B blok wątków (*block*)

T wątek (*thread*)

x,y,z współrzędne „kierunkowe”

Schemat mnożenia macierzy z podziałem C na bloki



Jądro mnożenia macierzy w CUDA C – matmul2.cu

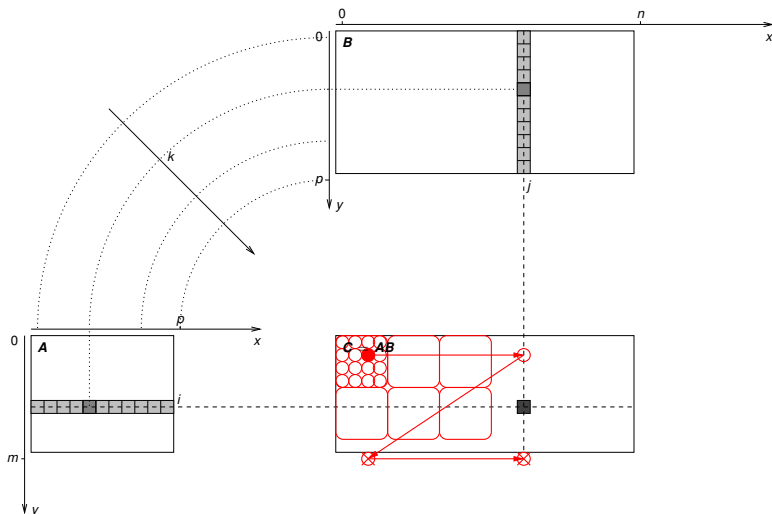
```
1 void matmul(float *C, float *A, float *B,
2             int m, int p, int n, ...)
3 { ...
4     dim3 dimBlock(N, M), dimGrid((n+N-1)/N, (m+M-1)/M);
5     matmul_kernel<<<dimGrid, dimBlock>>>
6         (dev_C, dev_A, dev_B, m, p, n);
7     ...
8 }
9 __global__
10 static void matmul_kernel(
11     float *C, float *A, float *B, int m, int p, int n)
12 {
13     int i = threadIdx.y + blockIdx.y*blockDim.y; if (i < m) {
14     int j = threadIdx.x + blockIdx.x*blockDim.x; if (j < n) {
15         float s = 0;
16         for (int k = 0; k < p; k++)
17             s += A[i*p + k] * B[k*n + j];
18         C[i*n + j] = s;
19     }}}
```

Porównanie wydajności implementacji CPU i GPU

- Do dalszych testów przyjmujemy $M = N = 32$. Tym razem spełnione jest już ograniczenie $T_B = MN \leq T_{B \max}$.
- GPU (GeForce RTX 4070)
 GPU (kernel) time = 1.039 ms (2047.927 GFLOP/s)
 CPU (total!) time = 48.761 ms (43.655 GFLOP/s)
- Jeśli m (albo n) nie jest wielokrotnością M (odpowiednio N), to musimy pokryć „kafelkami” macierz z nadmiarem...
- ...i dlatego wyniki są prawidłowe dopiero po ograniczeniu dziedziny obliczeń – inaczej dodajemy przypadkowe wartości i sięgamy poza przydzieloną pamięć.
- Bez trudu uzyskaliśmy przyspieszenie $550 \times$ lub $12 \times$.
- Kod nadal nie jest w pełni uniwersalny ze względu na ograniczenia $B_{G,x} \leq B_{G,x \max}$ i $B_{G,y} \leq B_{G,y \max}$, czyli:

$$mn \leq B_{G,x \max} B_{G,y \max} T_{B \max} \stackrel{3.0}{=} 1,4 \cdot 10^{17}.$$

Mnożenie macierzy o dużych (dowolnie!) rozmiarach



Jądro mnożenia macierzy w CUDA C – matmul3.cu

```
1 void matmul(float *C, float *A, float *B,
2             int m, int p, int n, ...)
3 { ...
4     dim3 dimBlock(N, M), dimGrid(gridSizeX, gridSizeY);
5     ...
6 }
7 __global__
8 static void matmul_kernel(
9     float *C, float *A, float *B, int m, int p, int n)
10 {
11     for (int i = threadIdx.y + blockIdx.y*blockDim.y;      i<m;
12           i += blockDim.y*blockDim.y) {
13         for (int j = threadIdx.x + blockIdx.x*blockDim.x;    j<n;
14               j += blockDim.x*blockDim.x) {
15             float s = 0;
16             for (int k = 0; k < p; k++)
17                 s += A[i*p + k] * B[k*n + j];
18             C[i*n + j] = s;
19         }
16     }
17 }
```

Analiza wydajności implementacji GPU

RTX 4070

- Program działa poprawnie dla dowolnych rozmiarów \mathbf{C} .
- Powróciliśmy do pierwotnego kodu zawierającego trzy zagnieżdżone pętle. Indeksy pętli zewnętrznych zwiększamy tym razem nie o 1, tylko o liczbę wątków wzdłuż danego wymiaru, równą $T_{B,x}B_{G,x}$ ($T_{B,y}B_{G,y}$).
- Każdy wątek oblicza wiele elementów \mathbf{C} – przyspieszenie.
- Wydajność jądra istotnie zależy od konfiguracji wykonania:

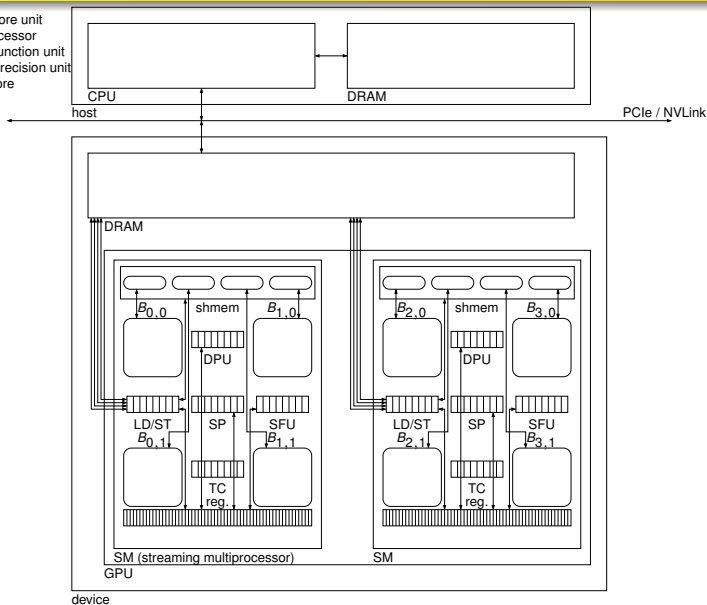
czas [ms]	$T_{B,x} = T_{B,y}$					
$B_{G,x} = B_{G,y}$	1	2	4	8	16	32
1				478.1	120.9	49.93
2			488.1	121.4	31.82	12.35
4		482.3	120.9	30.26	7.743	3.145
8	477.1	122.8	30.02	7.870	2.131	1.620
16	136.4	32.12	8.237	2.102	1.212	1.062
32	42.39	10.08	2.299	1.219	1.042	1.025

Analiza wydajności implementacji GPU

- Na RTX 4070 osiągnęliśmy wydajność ok. 2075.8 GFLOP/s, czyli przyspieszenie $S_u \approx 558$.
- Teoretyczna wydajność GPU wynosi ok. 35.48 TFLOP/s ($2 \text{ FLOP} \cdot \text{liczba rdzeni} \cdot \text{częstotliwość zegara}$) – $17 \times >!$
- Ograniczeniem są m.in. dostępy do pamięci. Algorytm bazuje na obliczaniu iloczynu skalarnego (CGMA tylko $\frac{2p-1}{2p+1} \approx 1$), co oznacza, że dane są niepotrzebnie wielokrotnie czytane z pamięci.
- Przepustowość GPU wynosi ok. 126 GT/s. Po pomnożeniu przez CGMA dostaniemy 126 GFLOP/s. Nasza $16,5 \times$ lepsza wydajność wynika z b. dużej pamięci *cache*.
- Sytuację można spróbować poprawić używając zamiast DDR DRAM dużo szybszej pamięci współdzielonej (SM),
- ... jest ona jednak znacznie mniejsza, więc trzeba będzie przetwarzać dane małymi porcjami kopiowanymi do niej.
- Kluczowe jest wykorzystanie hierarchii pamięci CUDA.

Hierarchiczny model sprzętowy systemu pamięci

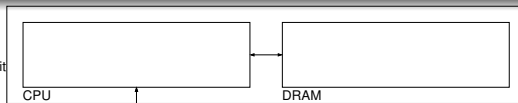
LD/ST - load/store unit
 SP - scalar processor
 SFU - special function unit
 DPU - double precision unit
 TC - „tensor” core



Hierarchiczny model programowy systemu pamięci

LD/ST - load/store unit
 SP - scalar processor
 SFU - special function unit
 DPU - double precision unit
 TC - „tensor” core

`cudaMemcpy(...)`



PCIe / NVLink

`__constant__`

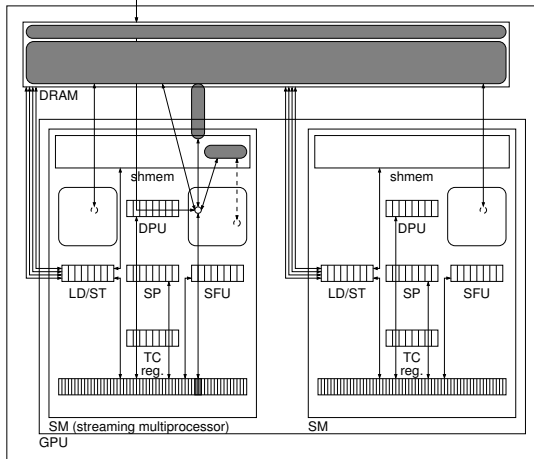
`__device__`

`cudaMalloc()`

`auto ≥ ...`

`__shared__`

`auto < ...`



pamięć stała

pamięć globalna

pamięć lokalna

pamięć współdzielona

rejstry

device

Hierarchiczny model programowy systemu pamięci

Pamięć*	Alokacja	Dostępność		Gdzie jest?
		Miejsce	Czas	
globalna	<code>__device__ v[DIM]</code> <code>cudaMalloc(&v, DIM)</code>	sieć	aplikacja	DRAM
stała [†]	<code>__constant__ v[DIM]</code>	sieć	aplikacja	DRAM
współdzielona	<code>__shared__ v[DIM]</code>	blok	jądro	SM
lokalna	<code>auto v[A_LOT]</code>	wątek	jądro	DRAM SM ^{2.0} (L1)
rejestry	<code>auto v[A_FEW]</code>	wątek	jądro	SM

Istnieją funkcje `__isGlobal()`, `__isConstant()`, `__isShared()`...

*Pomijamy pamięć tekstur i pamięć powierzchni^{2.0}, gdyż używane są one przede wszystkim do grafiki, a nie obliczeń. Jednak w niektórych aplikacjach można wykorzystać filtracyjne właściwości tekstur np. do interpolacji.

[†]Jest to mała (64 KiB) pamięć, buforowana (*cached*) na GPU. Dla $C_c \geq 2.0$ (gdzie GPU jest wyposażony w cache L2) jest rzadziej używana.

Rozmiary poziomów hierarchii systemu pamięci [KiB]

Pamięć	Parametr	2.0 2.1	3.0	3.2 3.5	5.0 5.3	6.1 7.0 8.6 8.9	6.0 6.2 7.5	8.0
globalna		∞ (pamięć zewnętrzna, rzędu 10 GiB)						
stała	C_{\max}	64						
współdz.	$S_{M \max}$	32×1,5		32×2	32×3	32×2	32×5	
lokalna	$A_T \max$	512						
rejestry (32 b)	$R_{M \max}$ [Ki]	32	64					
	$R_T \max$	63		255				

- Szybkie pamięci są małe, więc trzeba przetwarzać dane małymi porcjami, jawnie kopiowanymi do nich.
- Pamięć w rejestrach bywa większa niż w wolniejszej od nich pamięci współdzielonej – rejestry trzeba jak najlepiej wykorzystywać.

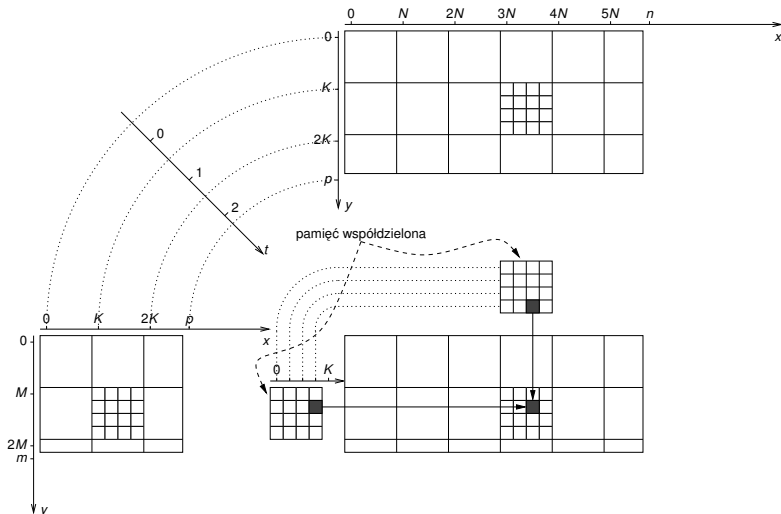
Mnożenie macierzy raz jeszcze – technika „kafelkowa”

- Przetwarzanie danych porcjami oznacza, że na „kafelki” musimy podzielić nie tylko macierz **C**, ale także **A** i **B**.
- Ze względu na łatwość zapisu i symetrię algorytmu przyjmujemy, kwadratowe „kafelki” $K \times K$ ($K = M = N$).
- Zmieniamy „ziarnistość” jądra (ale nie wątku!): algorytm opisujemy i realizujemy blokowo – dla każdego „kafelka” **C**.
- Obliczenie elementów „kafelka” **C** przebiega etapami, podczas których skokowo poruszamy się po „kafelkach” **A** (w prawo blokowego wiersza) i **B** (w dół blokowej kolumny):

Mnożenie macierzy techniką „kafelkową” – algorytm

- 1: **for all** „kafelek” macierzy **C** **do**
- 2: wyzeruj akumulator każdego elementu „kafelka” **C**
- 3: **for all** „kafelek” wiersza **A** i kolumny **B** **do**
- 4: pobierz kolejne kafelki **A** i **B** do pamięci współdzielonej
- 5:
- 6: dodaj do akumulatora każdego elementu „kafelka” **C**
 częstkowe iloczyny skalarne „kafelków” **A** i **B**
- 7: pobranych „kafelków” **A** i **B** zostaną dodane
- 8: **end for**
- 9: zapisz akumulator każdego elementu „kafelka” **C**
- 10: **end for**

Mnożenie macierzy techniką „kafelkową” – obliczenia



Zyski z mnożenia macierzy techniką „kafelkową”

- Dzięki identycznym i kwadratowym „kafelkom” każdy wątek ładuje do pamięci współdzielonej po jednym elemencie macierzy **A** i **B**, odpowiadającym mu „geograficznie”.
- Z raz pobranego elementu **A** (**B**) korzysta potem K wątków odpowiadających elementom danego wiersza (kolumny) **C**.
- Wymagana liczba odczytów danych z pamięci globalnej maleje K razy. Współczynnik pamięciowy złożoności obliczeniowej (CGMA) wynosi teraz w przybliżeniu K .
- Na obu etapach przetwarzania (kopiowanie i obliczenia) przetwarzamy tylko mały fragment danych wejściowych, uzyskując **lokalność** algorytmu w dostęпах do pamięci.
- W każdym obiegu pętli algorytmu (dla „kafelka” **C**) dokonujemy dwóch synchronizacji wątków na barierze (**poczekaj**). Ponieważ jest to możliwe tylko w ramach bloku, więc musimy utożsamić „kafelki” macierzy **C** z blokiem.
- **Warunkowa synchronizacja spowoduje zakleszczenie!**