# Large-Scale Scientific Computations

School of
Chemical Engineering

School of
Mechanical Engineering

School of
Electrical and Computer Engineering

**Computer Center**

**Parallel CFD
and Optimization Unit**

**Computing Systems
Laboratory**

Introduction to modern graphics processing units (GPU) architecture and programming in CUDA

Xenofon Trompoukis, Dr Mechanical engineer

School of Mechanical Engineering, NTUA,
Parallel CFD & Optimization Unit
email: xeftro@gmail.com

❑ Parallel processing units

❑ High floating-point operations rate (double and single precision arithmetic)

❑ GPU embedded, low latency, RAM

❑ Various programming environments

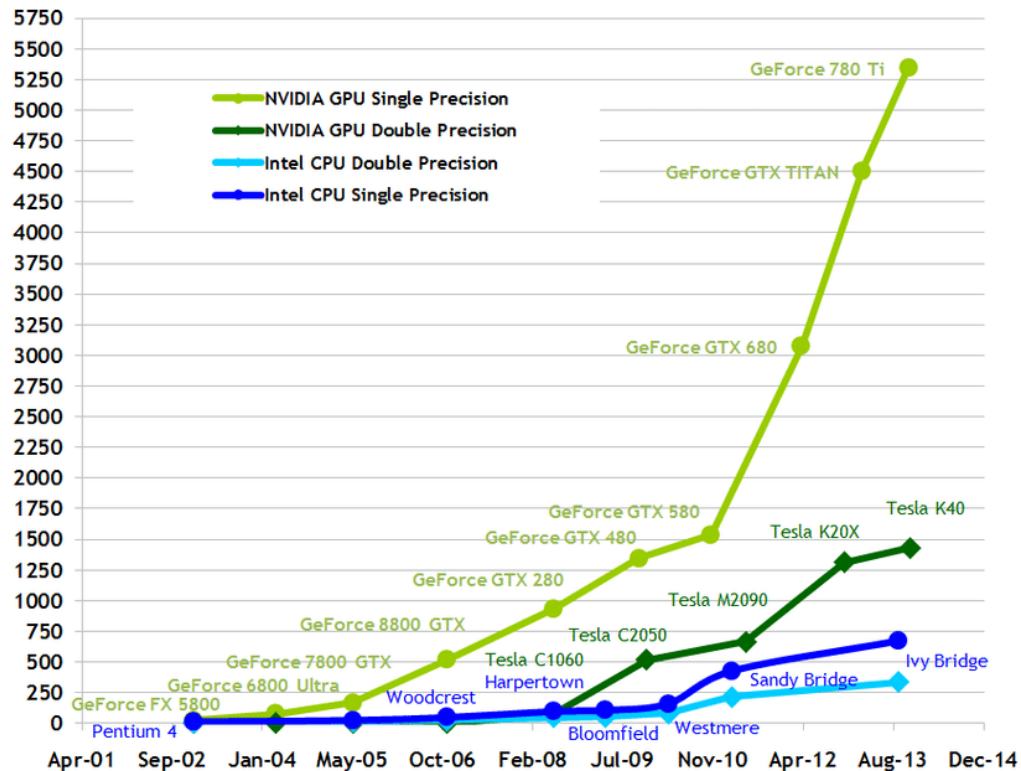❑ Low cost & energy consumption based on their computational power

# Why GPUs ?

☐ Parallel processing units

☐ High floating-point operations rate (double and single precision arithmetic)

☐ GPU embedded, low latency, RAM

☐ Various programming environments

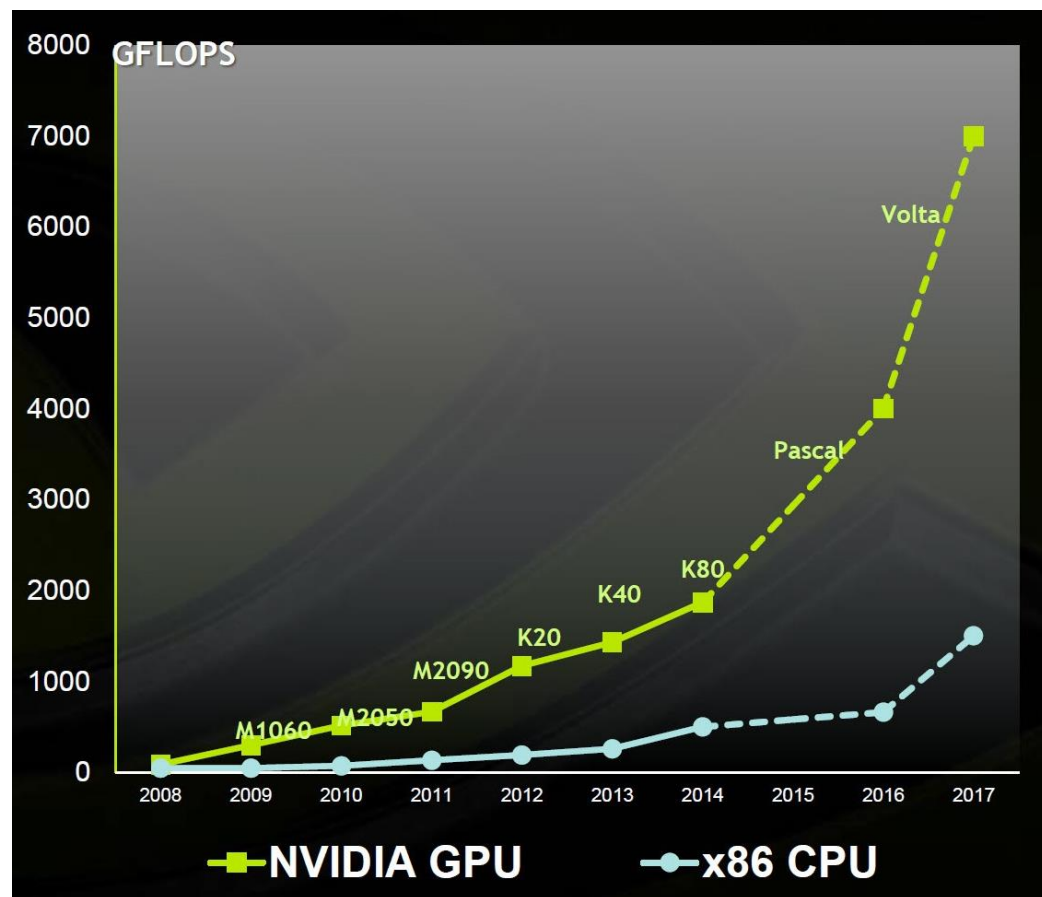☐ Low cost & energy consumption based on their computational power

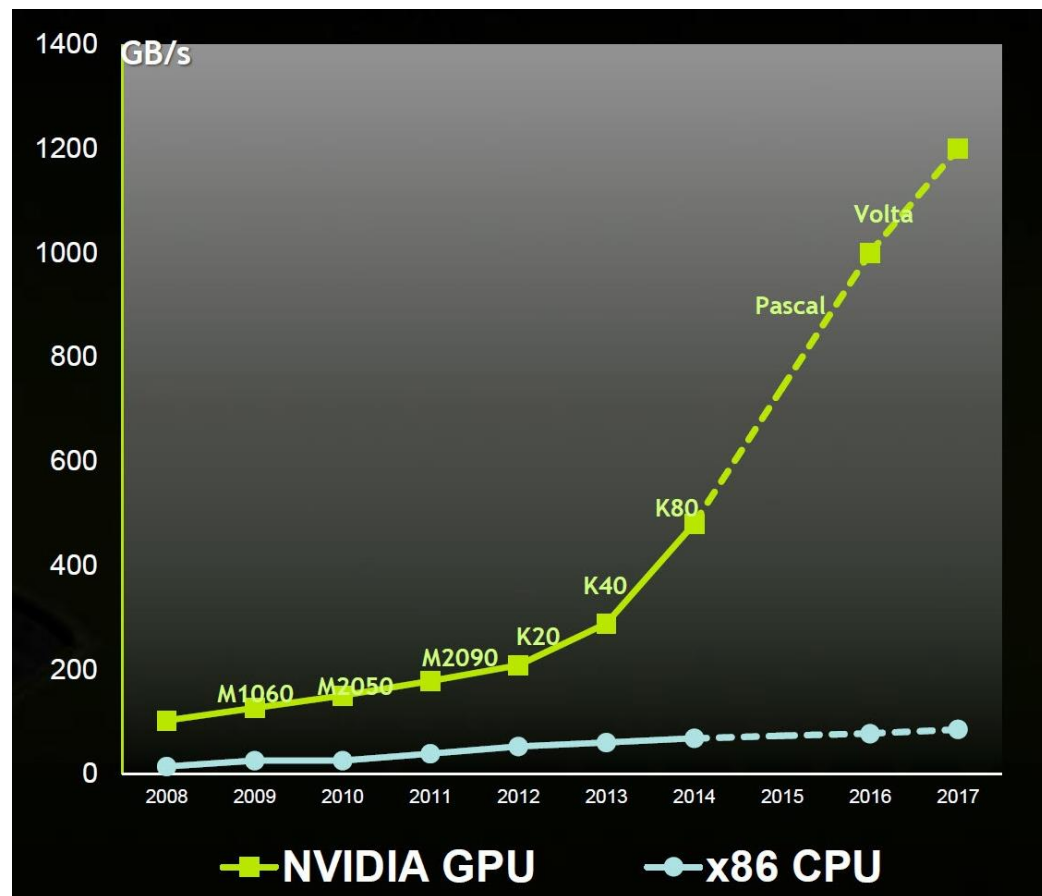**Theoretical GFLOP/s**

*CUDA Programming Guide v6.5*

Legend:
- NVIDIA GPU Single Precision
- NVIDIA GPU Double Precision
- Intel CPU Double Precision
- Intel CPU Single Precision

Data points (NVIDIA GPU Single Precision): GeForce FX 5800, GeForce 6800 Ultra, GeForce 7800 GTX, GeForce 8800 GTX, GeForce GTX 280, GeForce GTX 480, GeForce GTX 580, GeForce GTX 680, GeForce GTX TITAN, GeForce 780 Ti

Data points (NVIDIA GPU Double Precision): Tesla C1060, Tesla C2050, Tesla M2090, Tesla K20X, Tesla K40

Data points (Intel CPU): Pentium 4, Woodcrest, Harpertown, Bloomfield, Westmere, Sandy Bridge, Ivy Bridge

X-axis: Apr-01, Sep-02, Jan-04, May-05, Oct-06, Feb-08, Jul-09, Nov-10, Apr-12, Aug-13, Dec-14

Y-axis: 0, 250, 500, 750, 1000, 1250, 1500, 1750, 2000, 2250, 2500, 2750, 3000, 3250, 3500, 3750, 4000, 4250, 4500, 4750, 5000, 5250, 5500, 5750

# Why GPUs ?

☐ Parallel processing units

☐ High floating-point operations rate (double and single precision arithmetic)

☐ GPU embedded, low latency, RAM

☐ Various programming environments

☐ Low cost & energy consumption based on their computational power

# Why GPUs ?

☐ Parallel processing units

☐ High floating-point operations rate (double and single precision arithmetic)

☐ GPU embedded, low latency, RAM

☐ Various programming environments

☐ Low cost & energy consumption based on their computational power

# Why GPUs ?

❑ Parallel processing units

❑ High floating-point operations rate (double and single precision arithmetic)

❑ GPU embedded, low latency, RAM

❑ **Various programming environments**

❑ Low cost & energy consumption based on their computational power
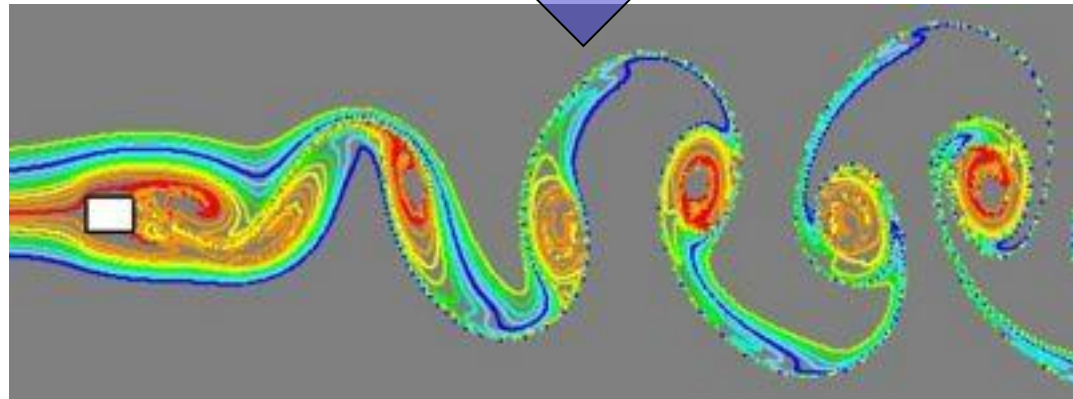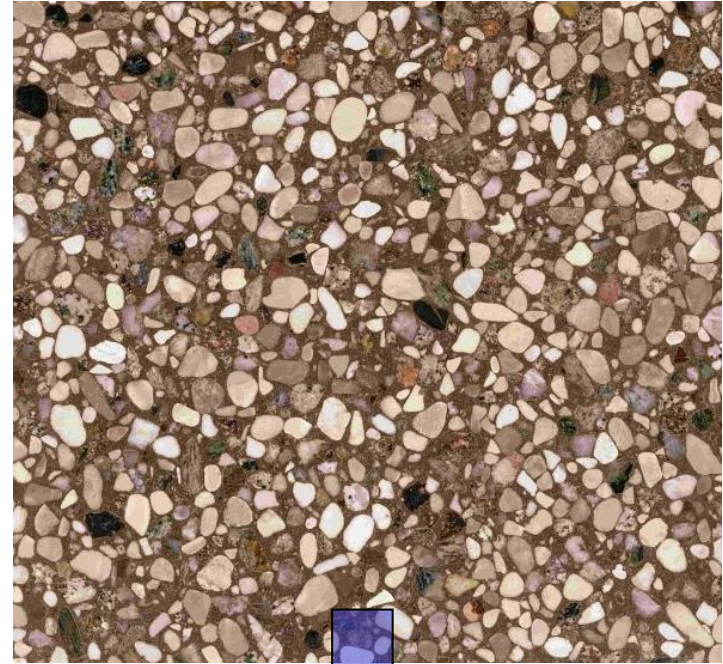
OpenCL: Cross platform implementation
- C++

CUDA: Developed by NVIDIA, specialized for NVIDIA GPUs
- C++
- FORTRAN
- Python

❏ Parallel processing units

❏ High floating-point operations rate (double and single precision arithmetic)

❏ GPU embedded, low latency, RAM

❏ Various programming environments

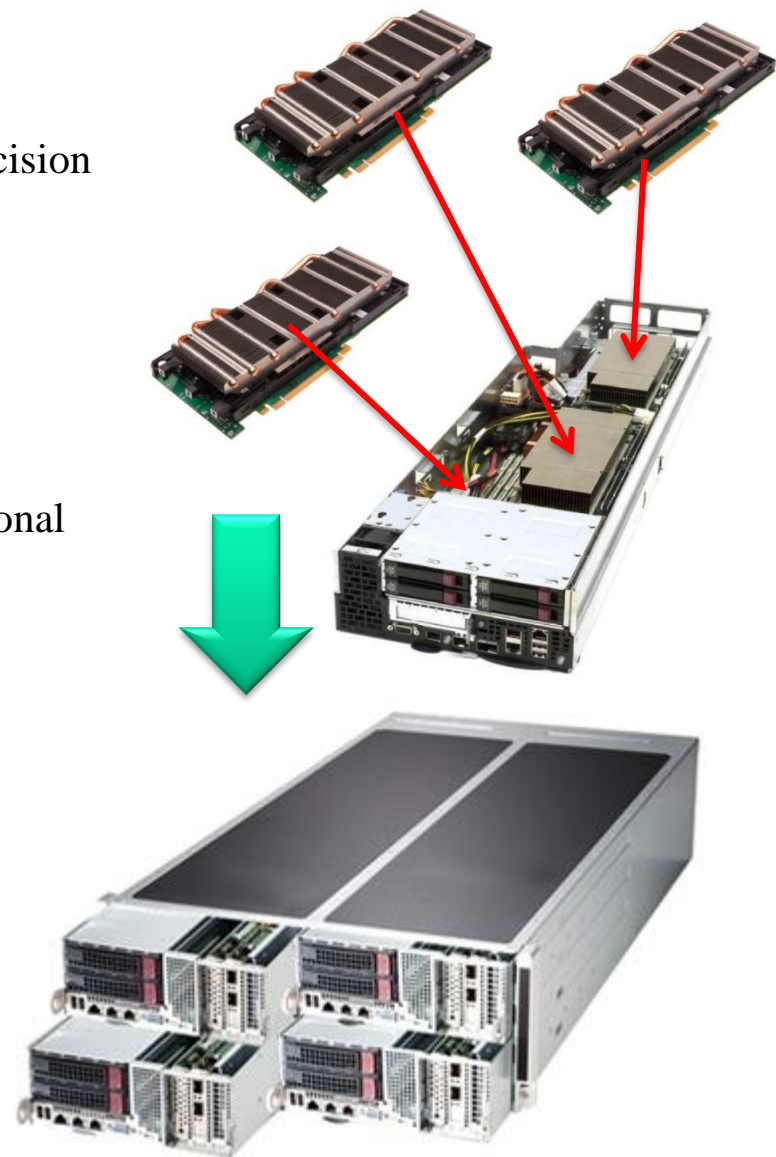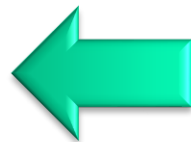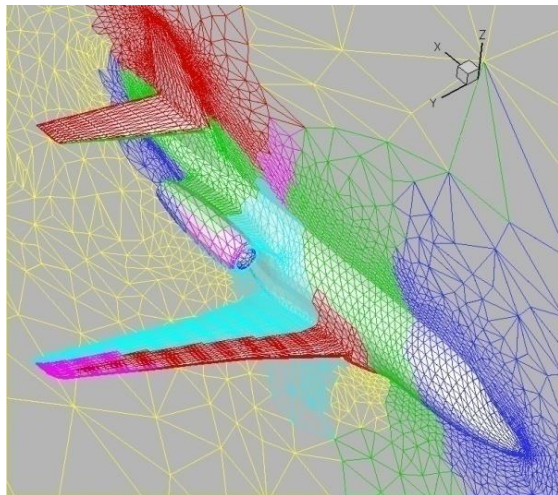❏ Low cost & energy consumption based on their computing power

# Why GPUs ?

❑ Parallel processing units

❑ High floating-point operations rate (double and single precision arithmetic)

❑ GPU embedded, low latency, RAM

❑ Various programming environments

❑ Low cost & energy consumption based on their computational power

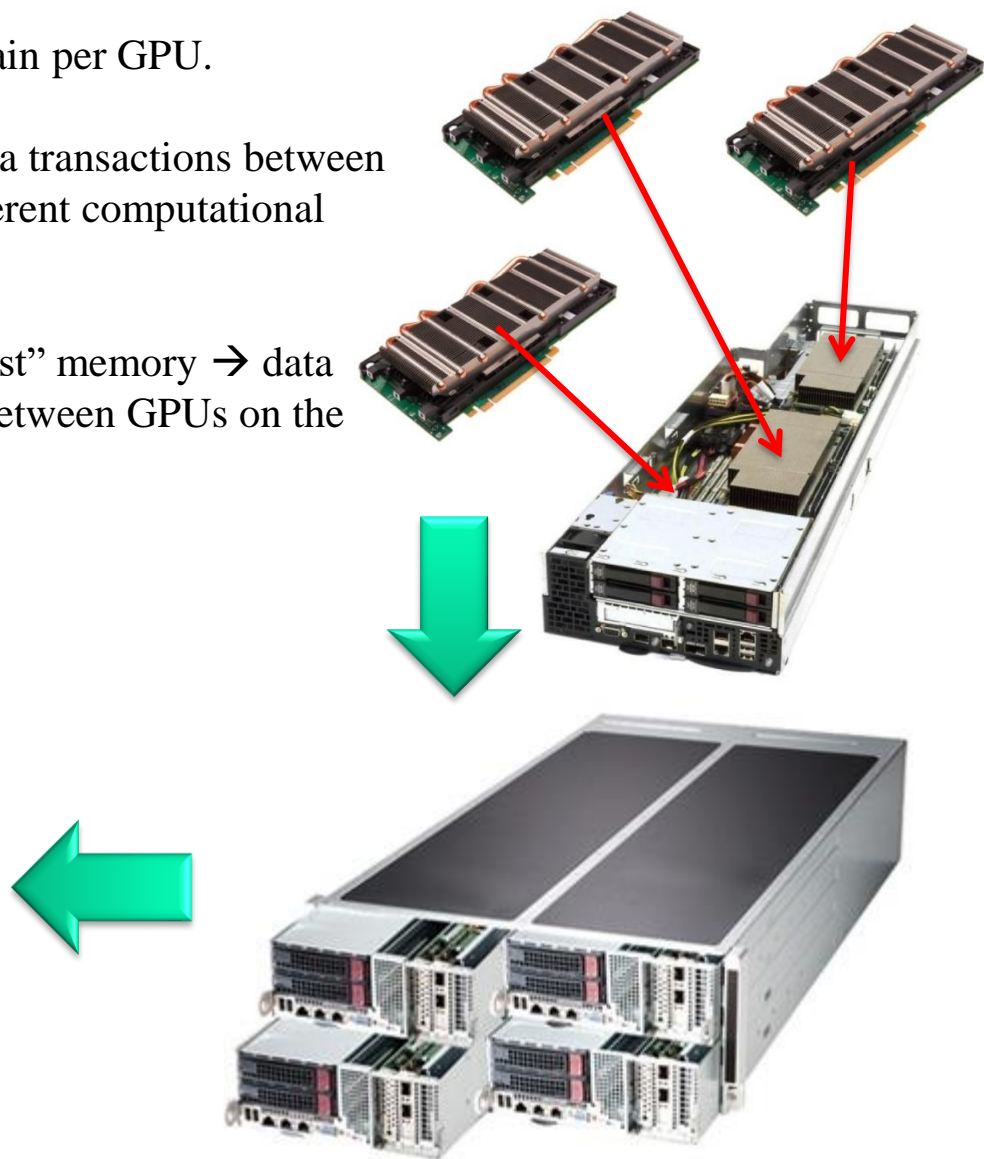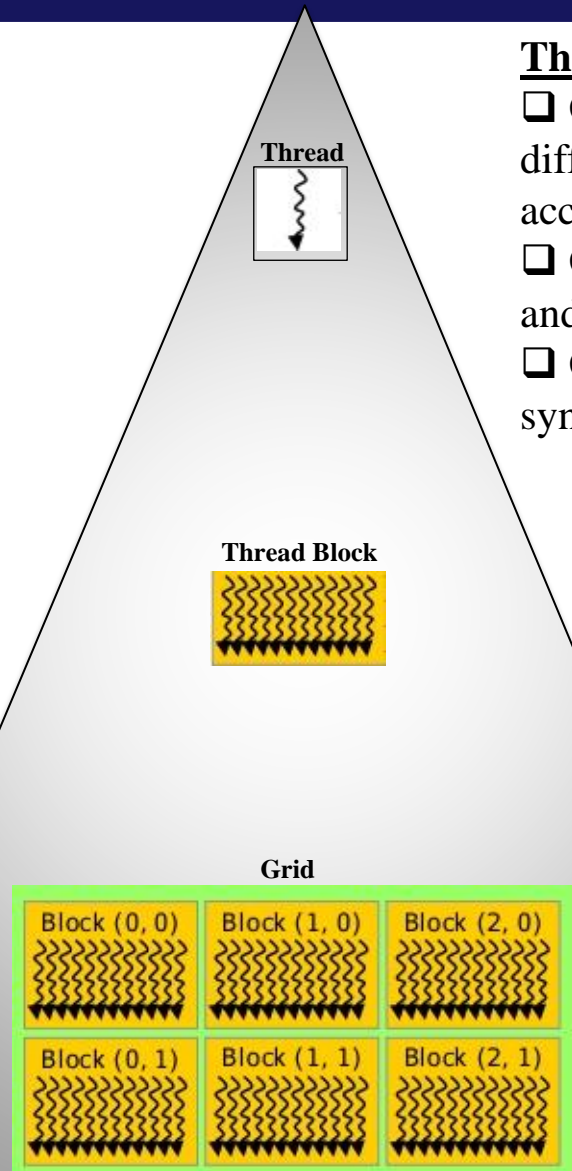❑ 1 sub-domain per GPU.

❑ MPI → data transactions between GPUs on different computational nodes.

❑ Shared "host" memory → data transactions between GPUs on the same node.

**GPUs** = Powerful, massively parallel CPU co-processors

# GPU Architecture

**Thread**

**Thread Block**

**Grid**

| Block (0, 0) | Block (1, 0) | Block (2, 0) |
| Block (0, 1) | Block (1, 1) | Block (2, 1) |

**Thread**: Fundamental computational unit
❑ GPU threads execute the same fragment of code (kernel) using different data (<u>SIMT</u>: Single Instruction Multiple Threads) accessing the same (device) memory.
❑ GPU threads are grouped in warps (i.e. group of 32 threads) and are executed at GPU multiprocessors.
❑ GPU threads of the same warp are executed in parallel in a synchronous manner.

**Block**: Cluster of warps
❑ Each multiprocessor can execute at least a thread block.
❑ GPU block threads, which belong to different warps, are executed in parallel and asynchronous manner.
❑ Synchronization and fast data transactions through shared memory

**Grid**: Cluster of thread blocks

✓ A FERMI GPU can execute up to 24,576 threads in parallel
✓ The programmer defines the thread block and grid dimensions

```cuda
 1 #include <cuda.h>
 2 #include <stdio.h>
 3 #include <string>
 4 #include <iostream>
 5
 6 // kernels :
 7 __global__ void helloGPU();
 8
 9 // host functions :
10 void Stop(std::string);
11
12 //
13 //
14 // ***********************************************************
15 int main()
16 // ***********************************************************
17 {
18     // kernel launch :
19     helloGPU<<< /*nbBlocks */ 1, /* nbThreads */ 1>>>();
20
21     // synchronize host/device :
22     cudaError_t err = cudaDeviceSynchronize();
23     if (err != cudaSuccess) Stop("");
24
25     return 0;
26 }
27
28 //
29 //
30 // ***********************************************************
31 void Stop(std::string error_message)
32 // ***********************************************************
33 {
34     std::cerr << error_message << std::endl;
35     exit(1);
36 }
37
38 //
39 //
40 // ***********************************************************
41 __global__ void helloGPU()
42 // ***********************************************************
43 {
44     printf("# hello world from thread %d in block %d\n",threadIdx.x,blockIdx.x);
45 }
```

**#include <cuda.h>**

**__global__ void helloGPU();**

**helloGPU<<<GridDim, BlockDim>>>();**

**Host-Device Synchronization**

NATIONAL TECHNICAL
UNIVERSITY OF ATHENS

EuroCC@Greece

**nvcc  –arch=sm_80  *.cu  –o  *.exe  –lcuda  – lcudart**

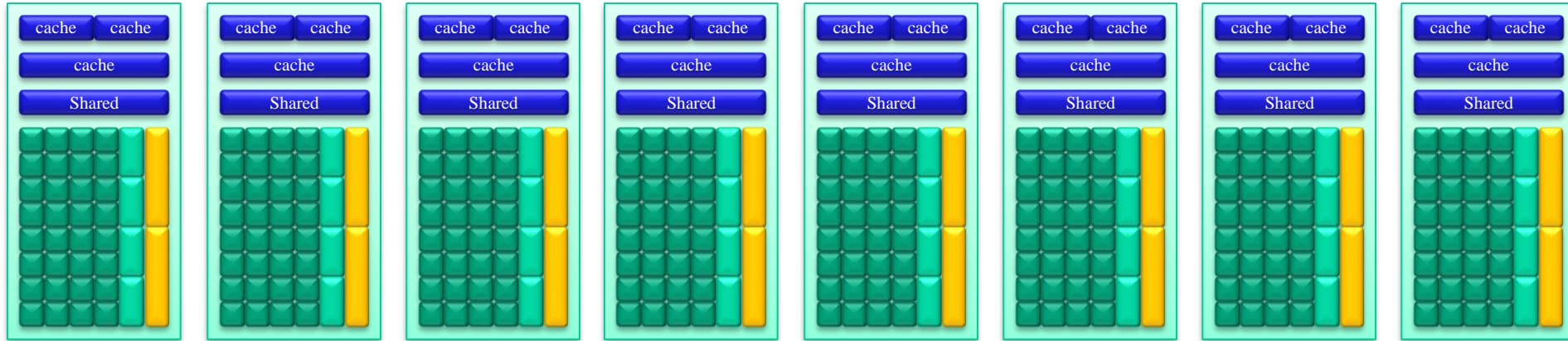compiler     GPU          Source     executable      libraries
             Compute      code
             Capability
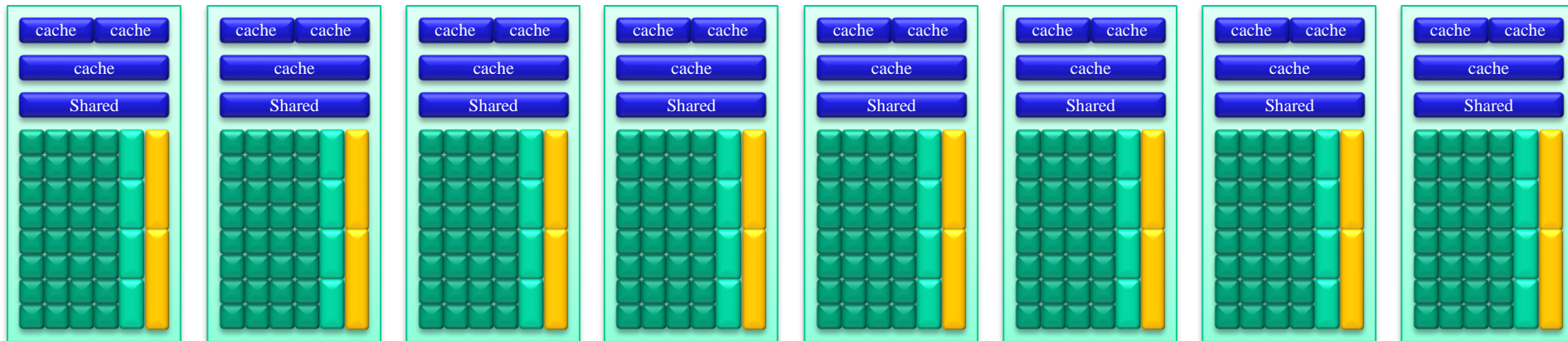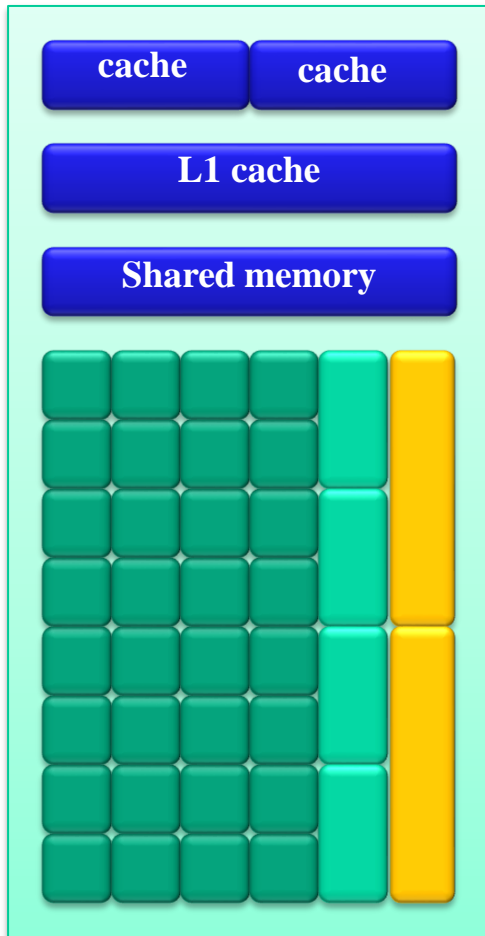
# GPU Architecture

Fermi: 16 multiprocessors



L2 cache

# GPU Architecture



1 FERMI multiprocessor consists of:

❑ 32 (CUDA) cores

❑ 4 Special Function Units (SFUs)

❑ 2 warp schedulers

❑ Shared memory

❑ cache memory (L1, constant & texture)

❑ 32768 32-bit registers

❑ Thread blocks are "split" into the multiprocessors based on kernel's requirements on registers and shared memory. Then, the warp schedulers of each multiprocessor organize block threads into warps.

❑ The best performing block size is related only with the GPU architecture and kernel requirements not with the application itself.

```cpp
21        // CPU allocations :
22        const int size = 10;
23        double  * A    = new double[size];
24        double  * B    = new double[size];
25        double  * C    = new double[size];
26
27        for (int i=0; i<size; i++) A[i] = (double)i;
28        for (int i=0; i<size; i++) B[i] = (double)i;
```

```cpp
53        // CPU deallocations :
54        if (A) delete[] A; A = NULL;
55        if (B) delete[] B; B = NULL;
56        if (C) delete[] C; C = NULL;
```

```
30      // GPU allocations :
31      double* _A = (double*)GPUalloc( (void*)A,size*sizeof(double),"A" );
32      double* _B = (double*)GPUalloc( (void*)B,size*sizeof(double),"B" );
33      double* _C = (double*)GPUalloc(           size*sizeof(double),"C" );
```

```
68  // ***************************************************************
69  void* GPUalloc(const int size, std::string error_message)
70  // ***************************************************************
71  {
72      void* devp = NULL;
73      cudaError_t err = cudaMalloc(&devp,size);
74      if (err != cudaSuccess) Stop("GPU allocation failed " + error_message);
75
76      return devp;
77  }
```

```
81  // ***************************************************************
82  void* GPUalloc(void* hostp, const int size, std::string error_message)
83  // ***************************************************************
84  {
85      void* devp = NULL;
86      cudaError_t err1 = cudaMalloc(&devp,size);
87      if (err1 != cudaSuccess) Stop("GPU allocation failed -" + error_message + "-\n"
88
89      cudaError_t err2 = cudaMemcpy(devp,hostp,size,cudaMemcpyHostToDevice);
90      if (err2 != cudaSuccess) Stop("invalid copy -" + error_message + "-\n");
91
92      return devp;
93  }
```

# Vector summation

```
35        // kernel launch :
36        const int nbThreads = 128;
37        const int nbBlocks  = (size+nbThreads-1)/nbThreads;      ←
38
39        vectorAddGPU<<<nbBlocks,nbThreads>>>(size,_A,_B,_C);
40
41        // "download" result :
42        cudaError_t err = cudaDeviceSynchronize();
43        if (err != cudaSuccess) Stop("error in kernel launch");
44
45        err = cudaMemcpy(C,_C,size*sizeof(double),cudaMemcpyDeviceToHost);
46        if (err != cudaSuccess) Stop("invalid copy -C-");
```

```
107 // ***********************************************************
108 __global__ void vectorAddGPU(const int size, double* _A, double* _B, double* _C)
109 // ***********************************************************
110 {
111      const int i = blockDim.x*blockIdx.x + threadIdx.x;      ←
112      if (i < size)                                            ←
113      {
114          _C[i] = _A[i] + _B[i];
115      }
116 }
```

```
58        // GPU deallocations :
59        cudaFree(_A); _A = NULL;                                ←
60        cudaFree(_B); _B = NULL;
61        cudaFree(_C); _C = NULL;
```

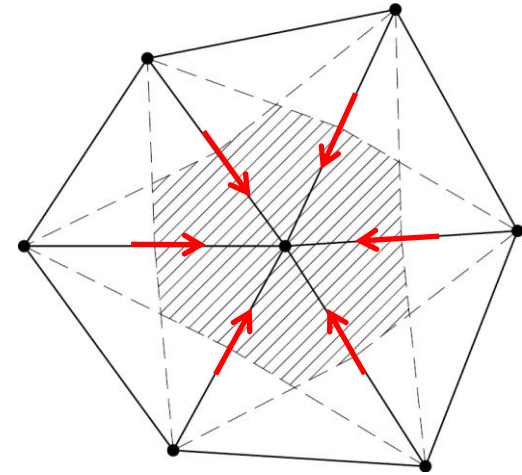❑ Avoid threads running in parallel to write at the same memory position (memory conflict).

❑ Threads from the same warp should access to the same device memory segment, since access to a 128-byte device memory segment can be performed within a single memory transaction.

Use Shared, constant and/or texture memory when possible.

❑ Be careful with *if statements* – avoid thread divergence.

❑ If it is possible, use single precision instead of double precision arithmetic. In Fermi GPUs, single precision operation rate is 2x higher than the double precision one.

❑ Use all the available resources (GPU + CPU).

# Piece of advice

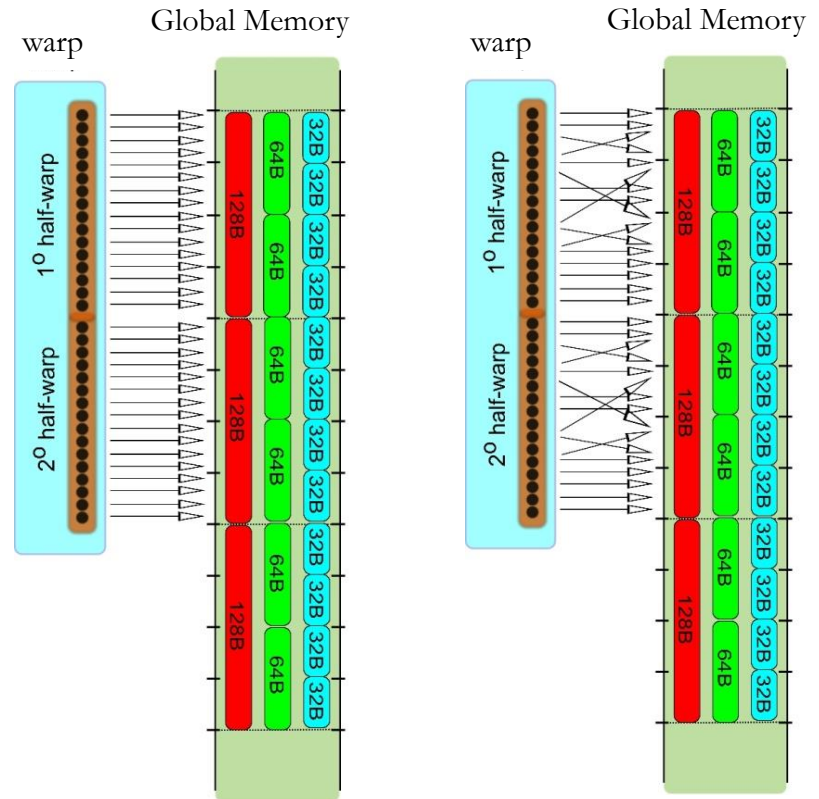❏ Avoid threads running in parallel to write at the same memory position (memory conflict).

❏ Threads from the same warp should access to the same device memory segment, since access to a 128-byte device memory segment can be performed within a single memory transaction.

Use Shared, constant and/or texture memory when possible.

❏ Be careful with *if statements* – avoid thread divergence.

❏ If it is possible, use single precision instead of double precision arithmetic. In Fermi GPUs, single precision operation rate is 2x higher than the double precision one.

❏ Use all the available resources (GPU + CPU).

# Piece of advice

❑ Avoid threads running in parallel to write to the same memory position (memory conflict).

❑ Threads from the same warp should access to the same device memory segment, since access to a 128-byte device memory segment can be performed within a single memory transaction.

Use Shared, constant and/or texture memory when possible.

❑ Be careful with *if statements* – avoid thread divergence.

❑ If it is possible, use single precision instead of double precision arithmetic. In Fermi GPUs, single precision operation rate is 2x higher than the double precision one.

❑ Use all the available resources (GPU + CPU).

```
 1 Some Instructions A
 2 if (logical_statement) {
 3         Some Instructions B
 4 }
 5 else if (another_logical_statement) {
 6         Some Instructions C
 7 }
 8 else {
 9         Some Instructions D
10 }
11 Some Instructions E
```

# Piece of advice

❑ Avoid threads running in parallel to write to the same memory position (memory conflict).

❑ Threads from the same warp should access to the same device memory segment, since access to a 128-byte device memory segment can be performed within a single memory transaction.

Use Shared, constant and/or texture memory when possible.

❑ Be careful with *if statements* – avoid thread divergence.

❑ If it is possible, use single precision instead of double precision arithmetic. In Fermi GPUs, single precision operation rate is 2x higher than the double precision one.

❑ Use all the available resources (GPU + CPU).

$$\overbrace{\frac{\partial \vec{R}}{\partial \vec{U}} \Delta \vec{U} = -R\left(\vec{U}\right)}^{\textbf{MPA}}$$

$$\underbrace{\frac{\partial \vec{R}}{\partial \vec{U}}}_{\textbf{SPA}} \qquad \underbrace{-R\left(\vec{U}\right)}_{\textbf{DPA}}$$

$$\vec{U}^{\,n+1} = \vec{U}^{\,n} + \Delta\vec{U}$$

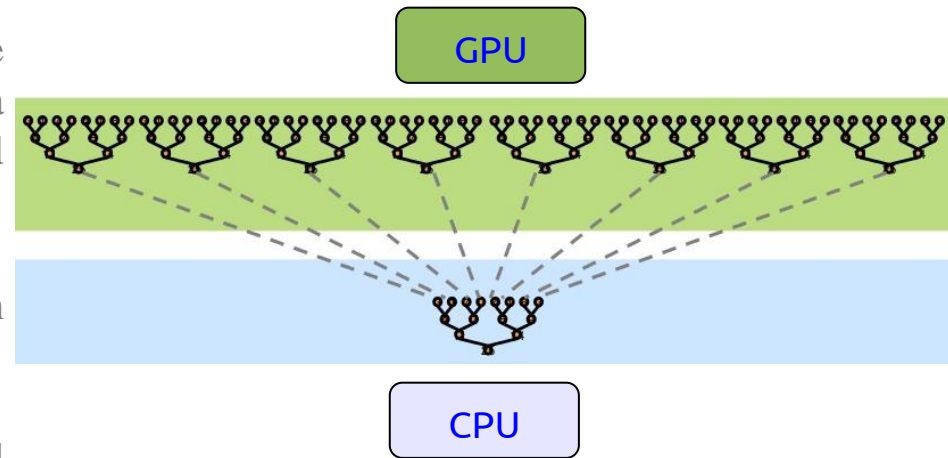❑ Avoid threads running in parallel to write to the same memory position (memory conflict).

❑ Threads from the same warp should access to the same device memory segment, since access to a 128-byte device memory segment can be performed within a single memory transaction.

Use Shared, constant and/or texture memory when possible.

❑ Be careful with *if statements* – avoid thread divergence.

❑ If it is possible, use single precision instead of double precision arithmetic. In Fermi GPUs, single precision operation rate is 2x higher than the double precision one.

❑ Use all the available resources (GPU + CPU).



GPU

CPU

# Dot product

```
  6 #define nbThreads 128
```

```
112 // *****************************************************************
113 __global__ void dotProductGPU(const int size, double* _A, double* _B, double* _C)
114 // *****************************************************************
115 {
116     __shared__ double dot[nbThreads]; dot[threadIdx.x] = 0.;
117 ▮
118     const int i = blockDim.x*blockIdx.x + threadIdx.x;
119     if (i < size)
120     {
121         dot[threadIdx.x] = _A[i] * _B[i];
122     }
123
124     block_reduction2(dot);
125
126     if (threadIdx.x == 0)
127     {
128         _C[blockIdx.x] = dot[0];
129     }
130 }
```

```
134 // *****************************************************************
135 __device__ void block_reduction1(double* dot)
136 // *****************************************************************
137 {
138     __syncthreads();
139
140     if (threadIdx.x == 0)
141     {
142         for (int it=1; it<nbThreads; it++) dot[0] += dot[it];
143     }
144 }
```

# Dot product

```
148  // **************************************************************
149  __device__ void block_reduction2(double* dot)
150  // **************************************************************
151  {
152      if (nbThreads != 128)
153      {
154          printf(" *** W A R N I N G : block_reduction2 works only for 128 threads per block\n");
155      }
156
157      __syncthreads(); if (threadIdx.x < 64) dot[threadIdx.x] += dot[threadIdx.x + 64];
158      __syncthreads(); if (threadIdx.x < 32) dot[threadIdx.x] += dot[threadIdx.x + 32];
159      __syncthreads(); if (threadIdx.x < 16) dot[threadIdx.x] += dot[threadIdx.x + 16];
160      __syncthreads(); if (threadIdx.x <  8) dot[threadIdx.x] += dot[threadIdx.x +  8];
161      __syncthreads(); if (threadIdx.x <  4) dot[threadIdx.x] += dot[threadIdx.x +  4];
162      __syncthreads(); if (threadIdx.x <  2) dot[threadIdx.x] += dot[threadIdx.x +  2];
163      __syncthreads(); if (threadIdx.x <  1) dot[threadIdx.x] += dot[threadIdx.x +  1];
164  }
```

# Summary

- ❏ \_\_global\_\_ : GPU function launched by the host (kernel)
- ❏ \_\_device\_\_ : GPU function launched by the device
- ❏ \_\_host\_\_    : CPU function launched by the host

- ❏ \_\_shared\_\_ : Variable in the shared memory

- ❏ \_\_syncthreads() : Block thread synchronization
- ❏ cudaDeviceSynchronize() : CPU-GPU synchronization

- ❏ cudaError_t cudaMalloc(void** ptr, size_t size);

- ❏ cudaError_t cudaFree(void* ptr);

- ❏ cudaError_t cudaMemcpy(void* destination, void* source, size_t size, cudaMemcpyKind kind);

# Matrix-matrix multiplication

```
82        // matrix-matrix multiplication on GPU :
83        dim3 dimBlock(blockSize,blockSize);
84        dim3 dimGrid (NI,NJ);
85
86        matmulGPU_2<<<dimGrid, dimBlock>>>(_A,_B,_C);
```

```
10 // constant variables :
11   __constant__ __device__ int _ni;
12   __constant__ __device__ int _nj;
13   __constant__ __device__ int _nk;
```
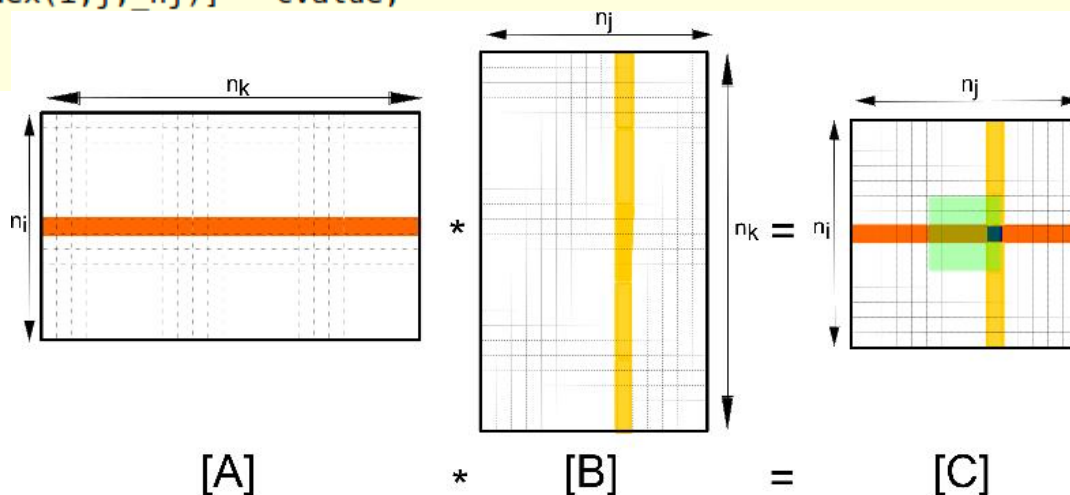
```
48        cudaError_t err1 = cudaMemcpyToSymbol(_ni, &ni, sizeof(int));
49        cudaError_t err2 = cudaMemcpyToSymbol(_nj, &nj, sizeof(int));
50        cudaError_t err3 = cudaMemcpyToSymbol(_nk, &nk, sizeof(int));
```

```
223 // *********************************************************************
224 __global__ void matmulGPU_1(double* _A, double* _B, double* _C)
225 // *********************************************************************
226 {
227         const int i = blockIdx.x*blockDim.x + threadIdx.x;
228         const int j = blockIdx.y*blockDim.y + threadIdx.y;
229
230         if (i < _ni && j < _nj)
231         {
232                 double cvalue = 0.;
233                 for (int k=0; k<_nk; k++)
234                 {
235                         cvalue += _A[index(i,k,_nk)] * _B[index(k,j,_nj)];
236                 }
237                 _C[index(i,j,_nj)] = cvalue;
238         }
239 }
```
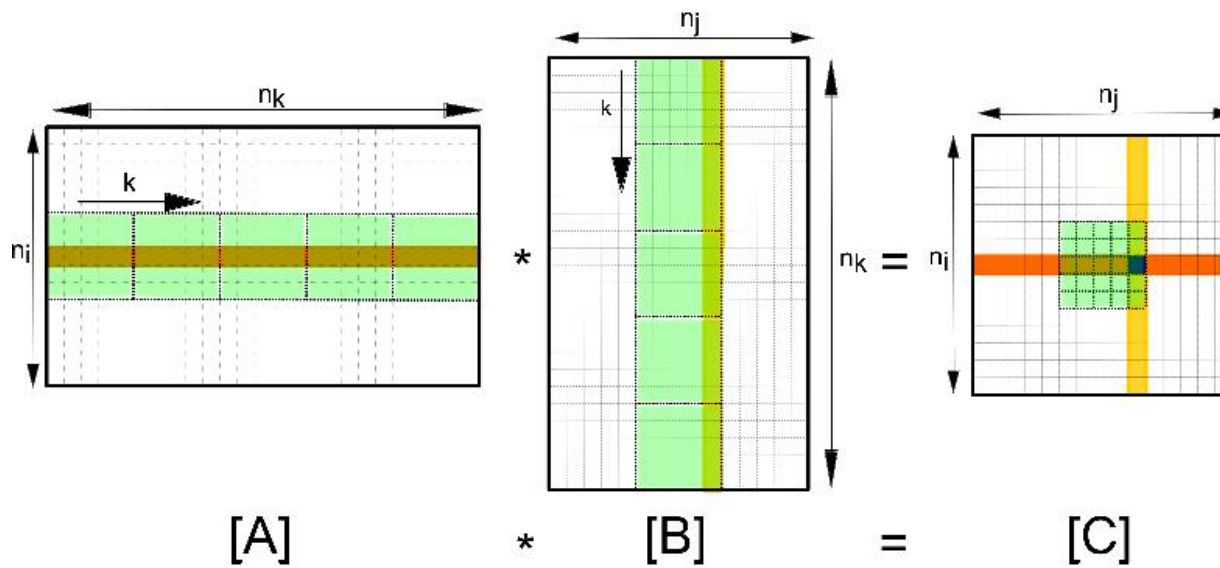


[A]    *    [B]    =    [C]

```
214 // *********************************************************************
215 __host__ __device__ int index (const int i, const int j, const int nj)
216 // *********************************************************************
217 {
218         return i*nj + j;
219 }
```

[A]    *    [B]    =    [C]

```
243 // *************************************************************
244 __global__ void matmulGPU_2(double* _A, double* _B, double* _C)
245 // *************************************************************
246 {
247         __shared__ double A[blockSize][blockSize];
248         __shared__ double B[blockSize][blockSize];
249
250      const int igl = blockIdx.x*blockDim.x + threadIdx.x;
251      const int jgl = blockIdx.y*blockDim.y + threadIdx.y;
252
253      double cvalue = 0.;
254      for (int k=0; k<_NK; k++)
255      {
256            A[threadIdx.x][threadIdx.y] = 0.;
257            B[threadIdx.x][threadIdx.y] = 0.;
258
259            const int iloc = k*blockSize + threadIdx.x;
260            const int jloc = k*blockSize + threadIdx.y;
261
262            if (jloc < _nk && igl < _ni) A[threadIdx.x][threadIdx.y] = _A[index(igl ,jloc,_nk)];
263            if (iloc < _nk && jgl < _nj) B[threadIdx.x][threadIdx.y] = _B[index(iloc,jgl ,_nj)];
264            __syncthreads();
265
266            for (int m=0; m<blockDim.y; m++) cvalue += A[threadIdx.x][m] * B[m][threadIdx.y];
267      }
268
269      if (igl < _ni && jgl < _nj)
270      {
271            _C[index(igl,jgl,_nj)] = cvalue;
272      }
273 }
```