# Mining Massive Datasets

## Lecture 14

**Artur Andrzejak**

[http://pvs.ifi.uni-heidelberg.de](http://pvs.ifi.uni-heidelberg.de)

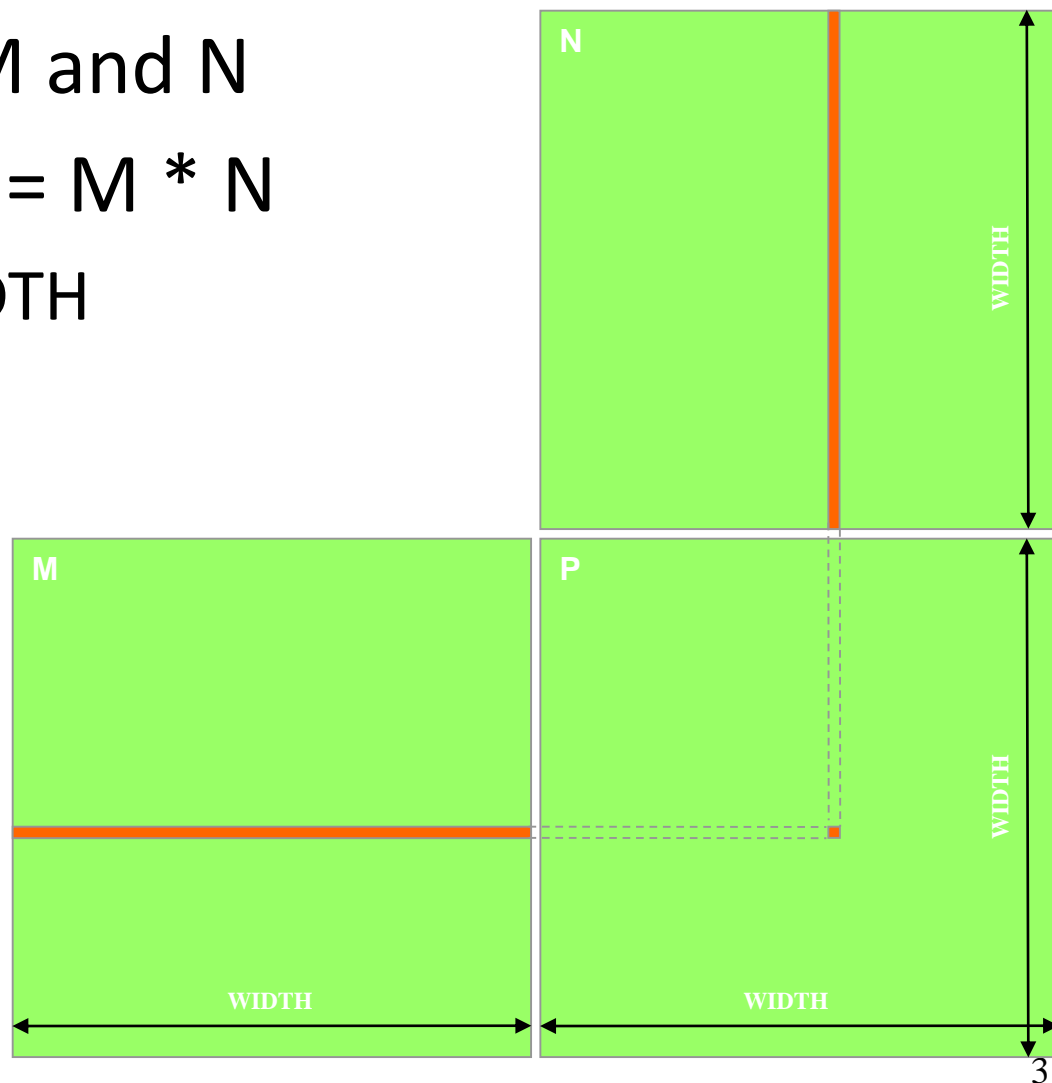# Introduction to GPU-Programming with CUDA:
## Motivation

Slides (partially modified) by:
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
ECE 498AL Spring 2010, University of Illinois, Urbana-Champaign

# Square Matrix Multiplication Example

- Input: matrices M and N

- Output: matrix P = M * N
  - Size WIDTH x WIDTH

# Each Matrix is Stored in a 1D-Array

$$\begin{array}{|c|c|c|c|}
\hline
M_{0,0} & M_{0,1} & M_{0,2} & M_{0,3} \\
\hline
M_{1,0} & M_{1,1} & M_{1,2} & M_{1,3} \\
\hline
M_{2,0} & M_{2,1} & M_{2,2} & M_{2,3} \\
\hline
M_{3,0} & M_{3,1} & M_{3,2} & M_{3,3} \\
\hline
\end{array}$$

M

$$\begin{array}{|c|c|c|c|c|c|c|c|c|c|c|c|c|c|c|c|}
\hline
M_{0,0} & M_{0,1} & M_{0,2} & M_{0,3} & M_{1,0} & M_{1,1} & M_{1,2} & M_{1,3} & M_{2,0} & M_{2,1} & M_{2,2} & M_{2,3} & M_{3,0} & M_{3,1} & M_{3,2} & M_{3,3} \\
\hline
\end{array}$$

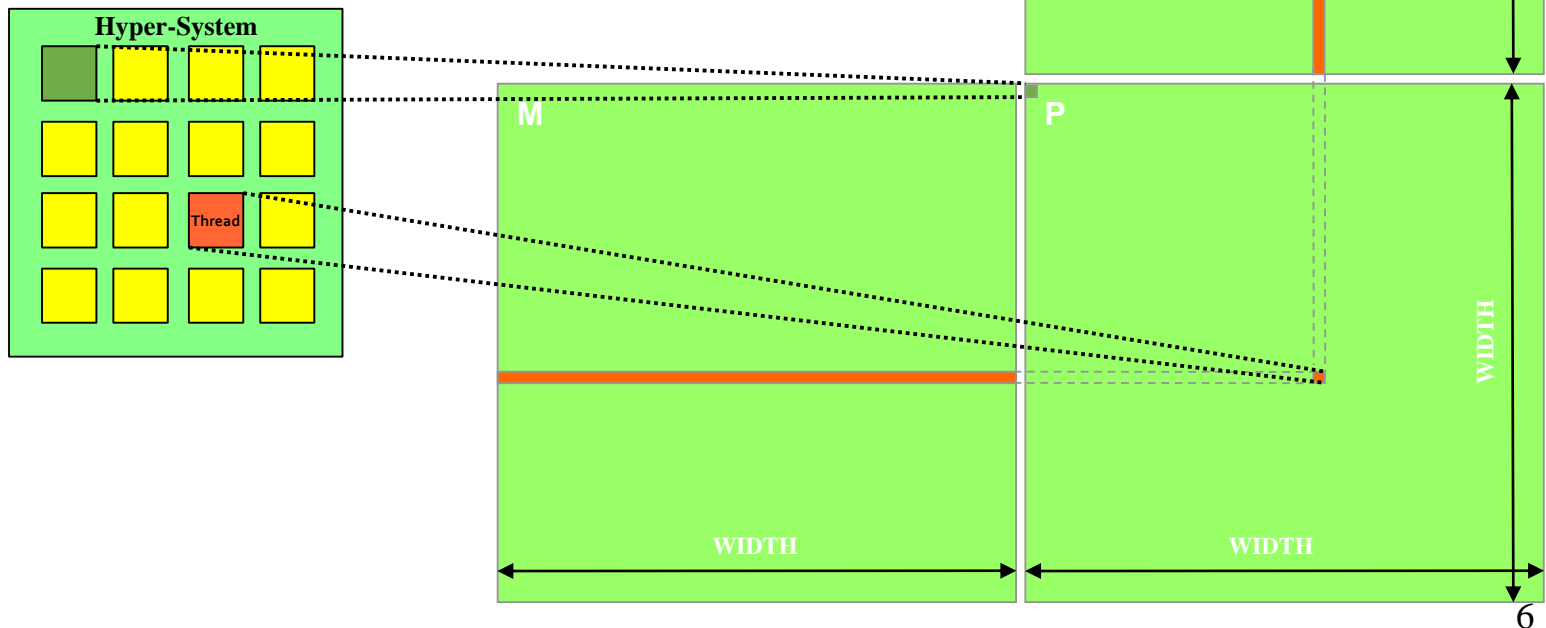=> Element M[i,j] is adressed via M[i*WIDTH + j]

# A Simple Host (CPU) Version in C

```c
void MatrixMulOnHost(float* M, float* N, float* P, int Width)
{
    for (int i = 0; i < Width; ++i)
        for (int j = 0; j < Width; ++j) {
            double sum = 0;
            for (int k = 0; k < Width; ++k) {
                double a = M[i * width + k];
                double b = N[k * width + j];
                sum += a * b;
            }
            P[i * Width + j] = sum;
        }
}
```



N

k

j

WIDTH

M

i

k

WIDTH

P

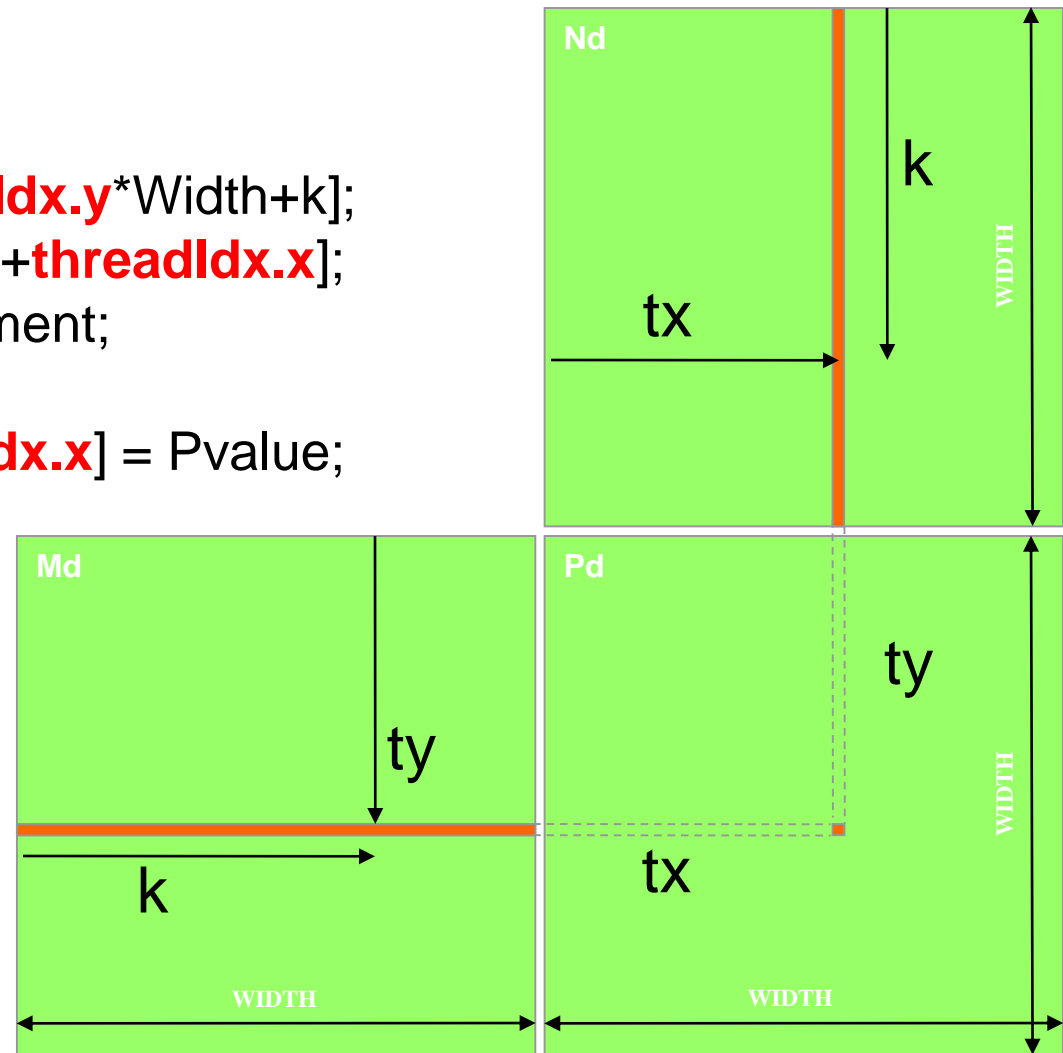Result matrix

WIDTH

WIDTH

# Square Matrix Multiplication Example

- Assume we have a system with WIDTH x WIDTH many cores
- Obvious parallelization:
    - Each thread runs on separate core and computes element P[i * Width + j]
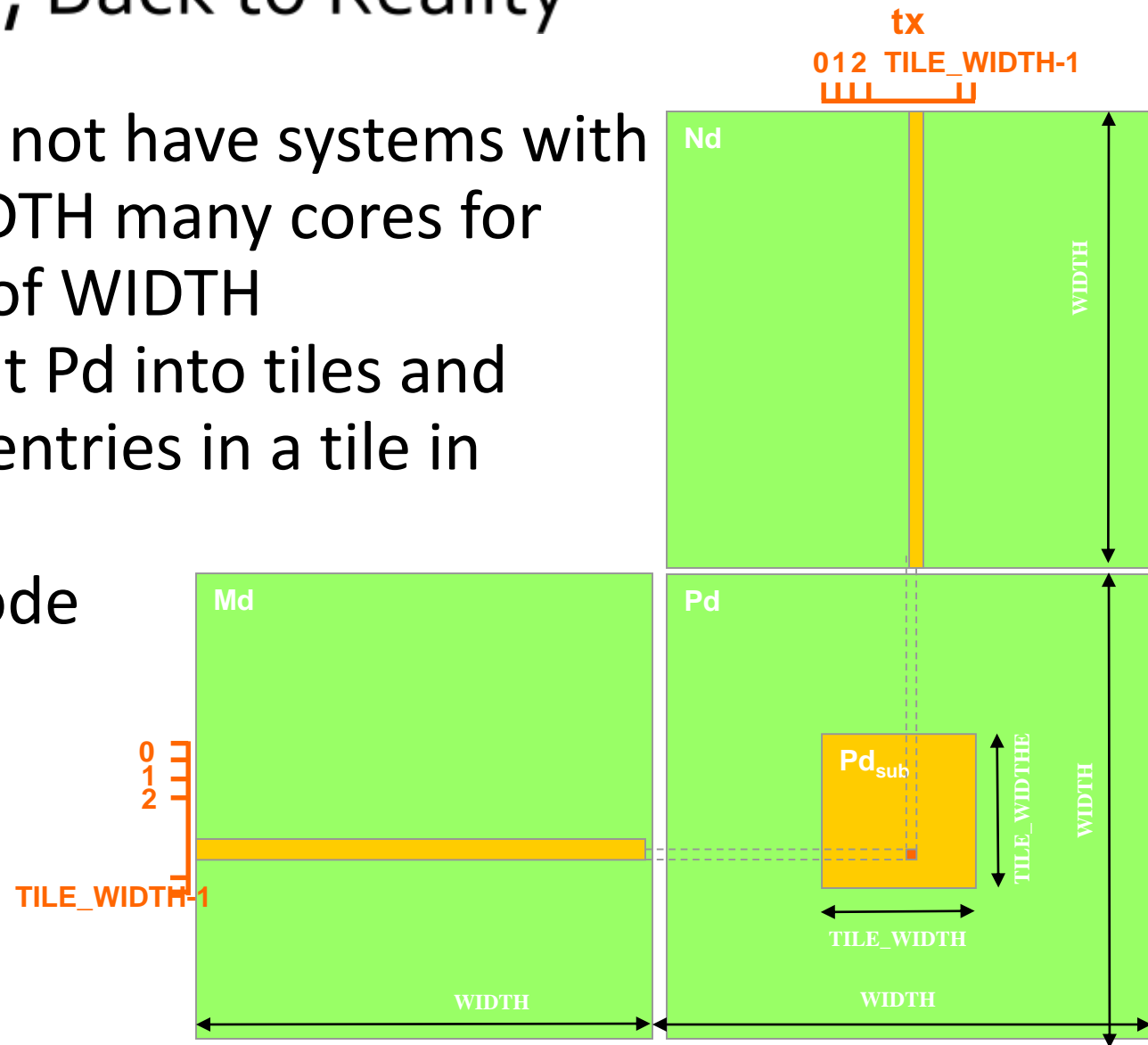
# Each Thread Runs Following Code

```
void MatrixMulKernel (float* Md, float* Nd, float* Pd, int Width)
{
float Pvalue = 0;
for (int k = 0; k < Width; ++k) {
      float Melement = Md[threadIdx.y*Width+k];
      float Nelement = Nd[k*Width+threadIdx.x];
      Pvalue += Melement * Nelement;
 }
  Pd[threadIdx.y*Width+threadIdx.x] = Pvalue;
}
```

- Where do **threadIdx.x** and **threadIdx.y** come from?
  - Their values are "magically" assigned by hardware
  - Different pairs of values for each thread

# Back to Life, Back to Reality

- <u>Problem</u>: we not have systems with WIDTH x WIDTH many cores for large values of WIDTH
- <u>Solution</u>: Split Pd into tiles and compute all entries in a tile in parallel
- => How to code this easily?

# Introduction to GPU-Programming with CUDA:
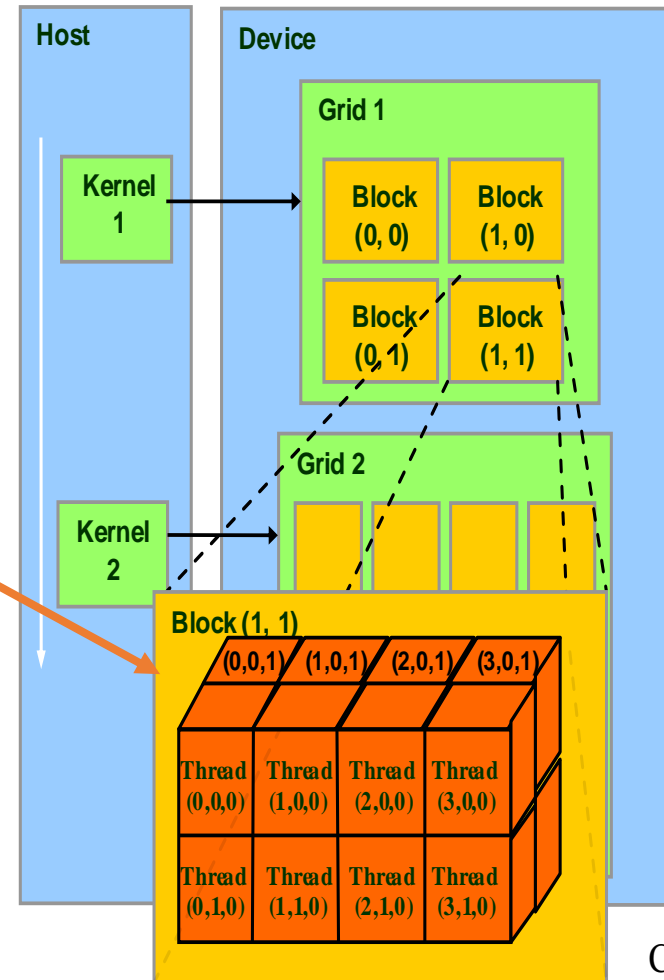## CUDA Architecture
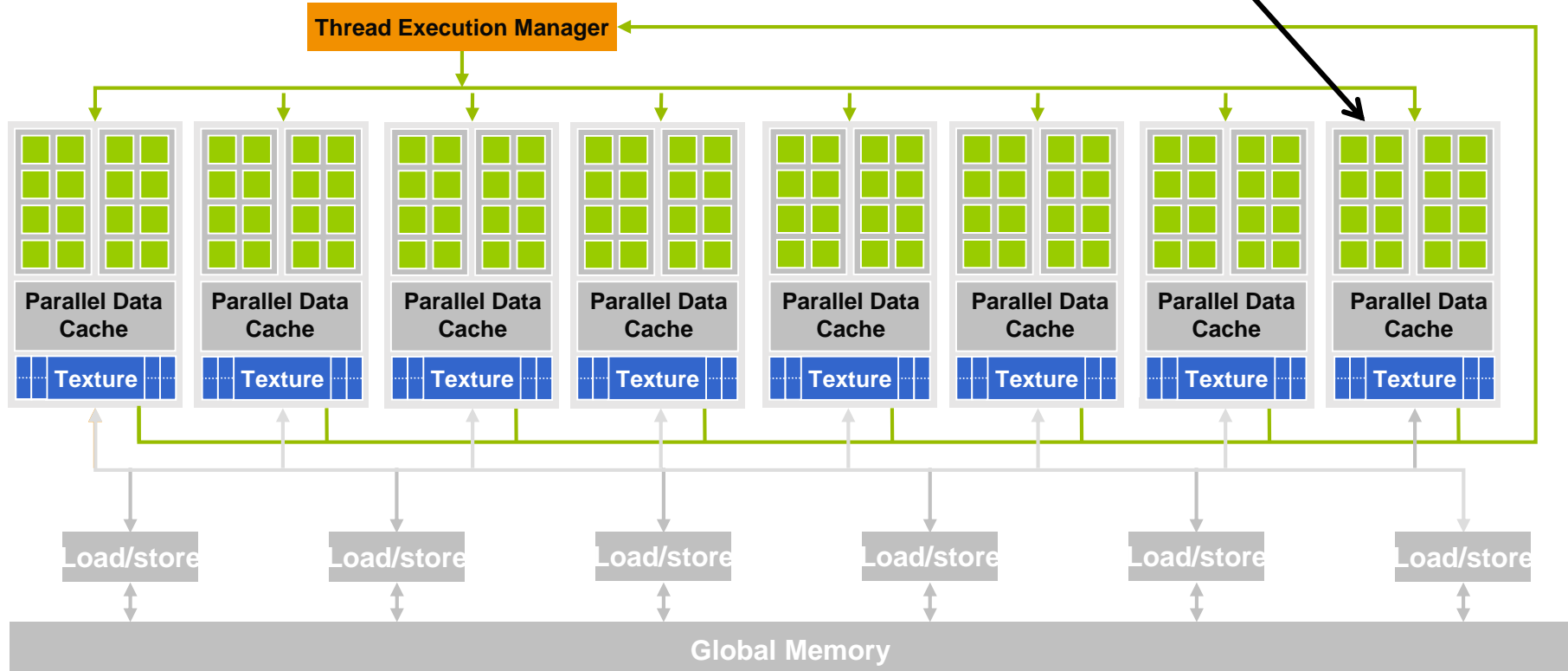
# Blocks of Threads in CUDA

- The **CUDA** programming model allows very large number of (virtual) threads
- These are grouped into **blocks**
  - Each block contains multiple threads, arranged as a 1D, 2D, or 3D array; here: 3D array

- Blocks correspond to units of hardware (**streaming multiprocesors**, **SM**)



Courtesy: NDVIA

# G8o CUDA mode – A **Device** Example

- Each block (or a part of it, a **warp**) runs on a separate streaming multiprocessor SM
- Example: the G80 CUDA GPU

# Threads within Blocks

- ## Each SM has many **streaming processors** (**SP**s)

  - Each SP executes "own" thread, all SPs the <u>same code</u>

  - Threads within a block cooperate via **shared memory**, **atomic operations** and **barrier synchronization**

  - Threads in different blocks <u>cannot</u> cooperate

# Grid of Blocks

- Blocks are organized into a 1D or 2D **Grid**
- Grid is the <u>unit of execution</u>
  - => A program is always executed as a grid of blocks

- BUT: We cannot specify the <u>order</u> of execution of the blocks
  - => Each block is completely independent of other blocks



Courtesy: NDVIA

# Back to the Matrix Example

- Each 2D thread <u>block</u> computes a (TILE_WIDTH)$^2$ sub-matrix (<u>tile</u>)
  - Each has (TILE_WIDTH)$^2$ threads
- The whole matrix is split into a 2D grid of tiles (or, equivalently blocks)
- => Because of CUDA's grid of blocks <u>we don't need nested loops over tiles</u>!

You still need to put a loop around the kernel call for cases where WIDTH/ TILE_WIDTH is greater than max grid size (64K)!

# Matrix Multiplication Using Multiple Blocks

- **Break-up Pd into tiles**
- **Each block calculates one tile**
  - Each thread calculates one element
- **Block size equal tile size**

# A Small Example

Block(0,0)      Block(1,0)

$$
\begin{array}{|cc|cc|}
\hline
P_{0,0} & P_{1,0} & P_{2,0} & P_{3,0} \\
P_{0,1} & P_{1,1} & P_{2,1} & P_{3,1} \\
\hline
P_{0,2} & P_{1,2} & P_{2,2} & P_{3,2} \\
P_{0,3} & P_{1,3} & P_{2,3} & P_{3,3} \\
\hline
\end{array}
$$

TILE_WIDTH = 2

Block(0,1)      Block(1,1)

16

# A Small Example: Multiplication

# Revised Matrix Multiplication Kernel using Multiple Blocks

```
__global__ void MatrixMulKernel(float* Md, float* Nd, float* Pd, int Width)
{
// Calculate the row index of the Pd element and M
int Row = blockIdx.y*TILE_WIDTH + threadIdx.y;
// Calculate the column idenx of Pd and N
int Col = blockIdx.x*TILE_WIDTH + threadIdx.x;

float Pvalue = 0;
// each thread computes one                           matrix
for (int k = 0; k < Width; ++k)
  Pvalue += Md[Row*Width+k] * Nd[k*Width+Col];

Pd[Row*Width+Col] = Pvalue;
}
```

Values provided „magically" by the CUDA-environment

# Introduction to GPU-Programming with CUDA:
## Blocks, Threads and Warps

Slides (partially modified) by:
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
ECE 498AL Spring 2010, University of Illinois, Urbana-Champaign

# Grid of Blocks and Block of Threads

- Threads are organized in CUDA in two levels:
  - **Block of threads**: a container of threads with 1D, 2D or 3D indexing
  - **Grid of blocks**: a container of blocks with 1D or 2D indexing
- A single <u>identical</u> program (**kernel**) is launched on a grid of blocks
  - Each thread receives "magically" different block-in-grid and thread-in-block indices



20

# Why Execution over a <u>Grid</u> of Blocks?

- 1. Our code doesn't need one or two outer loops for iterating over blocks (= tiles for MM)
- 2. Scalability: If a (future) GPU has more processors, more blocks can be executed in parallel (=> Good to have many blocks!)



Each block can execute in any order relative to other blocks.

# G80 Example: Executing Thread Blocks



**Blocks**

- Blocks are assigned to **Streaming Multiprocessors**

  - G80: Max. **8** blocks per SM

  - In total, max. **768** threads per SM

  - Use max #threads per SM (later more):

    - **Good**: 256 (threads/block) * 3 blocks => 768 threads

    - **Bad**: 512 (threads/block) * 1 block => 512 threads

▸ HD keeps track of block id's, thread id's and their status => limitations

# Thread Scheduling/Execution

- Each block is divided in 32-thread **warps**
  - 4 cycles per warp needed on G80
- How many warps per SM?
  - 3 blocks are assigned to a SM
  - Each block has 256 threads
  - How many warps are there in an SM?
    - Each block is divided into 256/32 = 8 Warps
    - There are 8 * 3 = 24 Warps
- At any point in time, <u>only one</u> of the 24 warps will be selected for instruction fetch and execution

G8o: 8 SPs per MP

Block 1 Warps
...
t0 t1 t2 ... t31

Block 2 Warps
...
t0 t1 t2 ... t31

**Streaming Multiprocessor**

| Instruction L1 | Data L1 |
|---|---|

**Instruction Fetch/Dispatch**

**Shared Memory**

| SP | | SP | |
|---|---|---|---|
| SP | SFU | SP | SFU |
| SP | | SP | |
| SP | | SP | |

# SM Warp Scheduling

**SM multithreaded Warp scheduler**

**time**

warp 8 instruction 11

warp 1 instruction 42

warp 3 instruction 95

warp 8 instruction 12

warp 3 instruction 96

- Switching between warps allows for (memory) **latency hiding**
  - While a warp waits for data from memory, other warps (with ready data) are executed
  - No overhead: hardware implements zero-overhead Warp scheduling
- It is important to <u>have many threads per SM</u> (see "Use max #threads per SM") on slide 8
  - More threads => more warps per SM => higher probability to find a warp with ready data (operands) => better latency hiding

# Thread Divergence at Branching

- Main performance concern with branching is **divergence**
  - = Threads <u>within a single warp</u> take different paths

  A;
  if B then
      **C**
  else
      **D**;
  E;

- G80: Different execution paths are <u>serialized</u>

  Warp

# Introduction to GPU-Programming with CUDA:
## Memory Efficiency

Slides (partially modified) by:
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
ECE 498AL Spring 2010, University of Illinois, Urbana-Champaign

# CUDA Memory Model Overview

- Global memory
  - Main means of communicating R/W data between host and device
  - Contents visible to all threads
  - Long latency access
- We will focus on optimizing access to global memory only

# Optimizing Memory Access of a Single Warp

- When accessing global memory, peak performance utilization occurs when all threads in a half warp read or write continuous memory locations
- Such (good) access pattern is called **coalesced** access pattern (dt. vereinigt)

Half Warp = 16 threads

# Coalesced Access and Compute Capability



- 16 mem. transactions (each) at compute capability 1.0 or 1.1

- 1 memory transaction at compute capability 1.2

# Optimizing Matrix Access /1

| $M_{0,0}$ | $M_{0,1}$ | $M_{0,2}$ | $M_{0,3}$ |
|-----------|-----------|-----------|-----------|
| $M_{1,0}$ | $M_{1,1}$ | $M_{1,2}$ | $M_{1,3}$ |
| $M_{2,0}$ | $M_{2,1}$ | $M_{2,2}$ | $M_{2,3}$ |
| $M_{3,0}$ | $M_{3,1}$ | $M_{3,2}$ | $M_{3,3}$ |

N

M    P

Thread 1
Thread k

Layout:

| $M_{0,0}$ | $M_{0,1}$ | $M_{0,2}$ | $M_{0,3}$ | $M_{1,0}$ | $M_{1,1}$ | $M_{1,2}$ | $M_{1,3}$ | $M_{2,0}$ | $M_{2,1}$ | $M_{2,2}$ | $M_{2,3}$ | $M_{3,0}$ | $M_{3,1}$ | $M_{3,2}$ | $M_{3,3}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Matrices are stored in the global memory

30

# Optimizing Matrix Access /2

| | | | |
|---|---|---|---|
| $M_{0,0}$ | $M_{0,1}$ | $M_{0,2}$ | $M_{0,3}$ |
| $M_{1,0}$ | $M_{1,1}$ | $M_{1,2}$ | $M_{1,3}$ |
| $M_{2,0}$ | $M_{2,1}$ | $M_{2,2}$ | $M_{2,3}$ |
| $M_{3,0}$ | $M_{3,1}$ | $M_{3,2}$ | $M_{3,3}$ |

Md: Good or bad?    Nd: Good or bad?

Md

Nd

Thread 1
Thread k

T1 Tk

Layout:

| $M_{0,0}$ | $M_{0,1}$ | $M_{0,2}$ | $M_{0,3}$ | $M_{1,0}$ | $M_{1,1}$ | $M_{1,2}$ | $M_{1,3}$ | $M_{2,0}$ | $M_{2,1}$ | $M_{2,2}$ | $M_{2,3}$ | $M_{3,0}$ | $M_{3,1}$ | $M_{3,2}$ | $M_{3,3}$ |

Matrices are stored in the global memory

# Nd: Good or bad?

Access direction in Kernel code



$M_{0,0}$ $M_{1,0}$ $M_{2,0}$ $M_{3,0}$

$M_{0,1}$ $M_{1,1}$ $M_{2,1}$ $M_{3,1}$

$M_{0,2}$ $M_{1,2}$ $M_{2,2}$ $M_{3,2}$

$M_{0,3}$ $M_{1,3}$ $M_{2,3}$ $M_{3,3}$

$T_1$ $T_2$ $T_3$ $T_4$

Waiting for a memory line: Other warps are executing (latency hiding)

Coalesced memory access

| Sched. Period 1 | Sched. Period 2 |
|---|---|
| $T_1$ $T_2$ $T_3$ $T_4$ | $T_1$ $T_2$ $T_3$ $T_4$ |

...

M

| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ | $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ | $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ | $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |

# Md: Good or bad?

Access direction in Kernel code

| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ | $T_1$ |
| $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ | $T_2$ |
| $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ | $T_3$ |
| $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ | $T_4$ |

$T_1...T_4$ of a single warp access „wide" memory addresses => not coalesced

**Scheduling Period 2**

$T_1$   $T_2$   $T_3$   $T_4$

...

**Scheduling Period 1**

$T_1$   $T_2$   $T_3$   $T_4$

M

| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ | $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ | $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ | $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |

# Thank you.

Questions?

# Additional Slides

# Introduction to GPU-Programming with CUDA:
## Clustering on GPUs

# Was sollen wir parallelisieren?

- Zur Vereinfachung werden i.A. nur Teile eines Programmes auf GPU programmiert
  - Aufwändiges wird auf GPU parallelisiert, Rest seriell auf CPU ausgeführt
- Was parallelisieren wir bei k-Means?

  - Berechnung der Abstände zwischen Punkten (Samples) und Centroiden und ihre Zuordnung zu den Centroiden

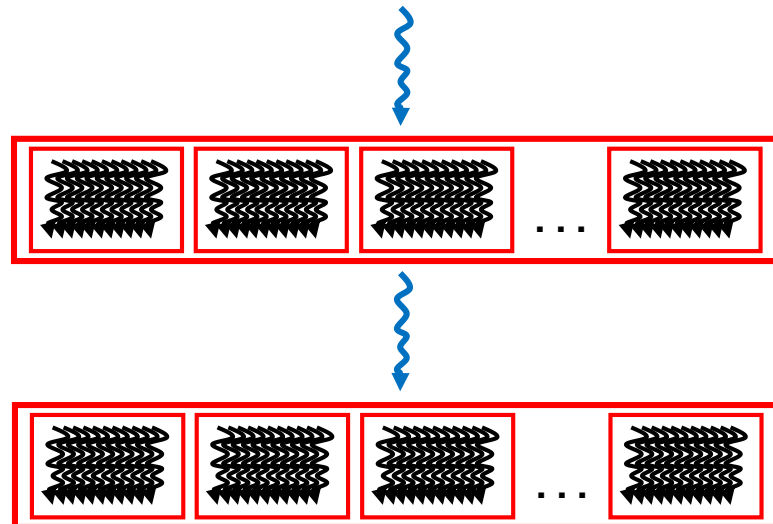Serieller Code (CPU)

Paralleler Code (GPU)

KernelA<<< nBlk, nTid >>>(args);

Serieller Code (CPU)

Paralleler Code (GPU)

KernelA<<< nBlk, nTid >>>(args);

# Kernel Routine - Übersicht

- Sei nts die Gesamtanzahl der GPU-Threads
- Der Thread t (t aus {0,..,nts-1}) bearbeitet Samples mit Indices t, t+nts, t+2*nts,..., t+p*nts (< n, #Samples)

- Für jeden dieser Samples $s_i$:

  - Für jeden Centroid $C_j$ (der k aktuellen Centroide):

    - Berechne die quadrierte Entfernung $D^2(s_i, c_j)$

    - Falls $D^2(s_i, c_j)$ kleiner ist als bisherige min. Centroid-Entfernung:

      - Setze $C_j$ als den nächsten (nearest) Centroid für $s_i$ fest

- Warte, bis alle Threads fertig sind

# Kernel – Parameter und Eingaben

- Wir haben die Anzahl der Blöcke (gridDim, ein 1D-Vektor) und die Anzahl der Threads pro Block (blockDim, auch ein 1D-Vektor) festgelegt
  - Abhängig von der GPU und Datengröße
- Jeder Thread bekommt „magisch" die Werte:
  - gridDim, blockDim, blockIdx, threadIdx

- Wir haben als weitere Eingabe für den Kernel:
  - dataSize (= n), numDims (= d), numClusters (= k)
  - data[][] – Matrix: n Samples mit jeweils d Koordinaten
  - centroids[][] – Matrix: k aktuelle Centroide, jeweils d El.
  - membership[] – Vektor von n Integers: Zuordnung von Samples zu den Centroiden bzw. Clustern
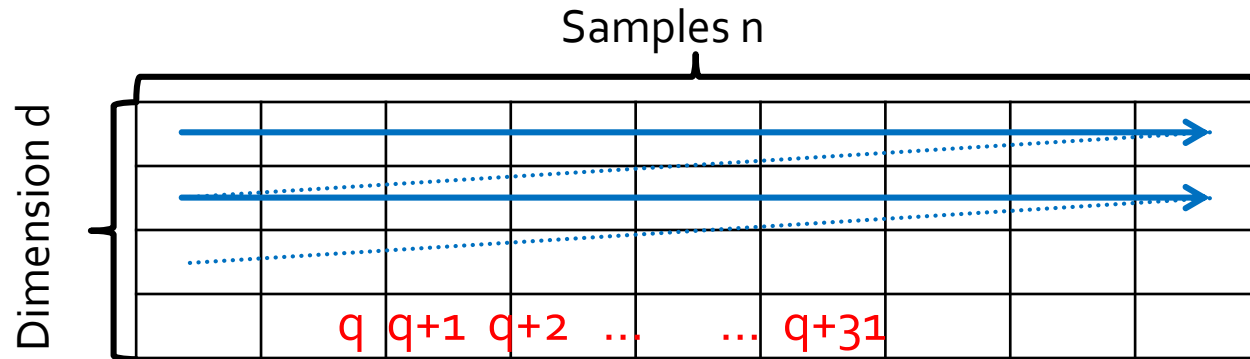
# Kernel - Pseudocode

threadId = blockDim*blockIdx + threadIdx
nts = blockDim*gridDim                        // total # of threads
**for** i = threadId **to** dataSize-1 **step** nts:      // get sample $s_i$

   minDistance = MAX_DOUBLE

   clusterIdx = 0

   **for** j = 0 **to** numClusters-1:                  // for each centroid $C_j$ …

      distance = 0

      **for** dim = 0 **to** numDims-1: // compute $D^2(s_i, c_j)$

         distance = distance + (data[dim][i] – centroids[j][dim])^2

      **if** distance < minDistance **then**:

         minDistance = distance

         clusterIdx = j

   membership[i] = clusterIdx
synchronize threads

# Code für Sample $S_i$ und Centroid $C_j$

```
for dim = 0 to numDims-1: // compute D²(sᵢ, cⱼ)
    distance = distance + (data[dim][i] – centroids[j][dim])^2
if distance < minDistance then:
    minDistance = distance
    clusterIdx = j
```

- Roter Code: Berechnet distance = $D^2(s_i, c_j)$
- Violetter Code: Überprüft, ob distance ein neues Minumum ist, ggf. aktualisiert clusterIdx
- Der rote Code ist kritisch für die Leistung
  - Warum? Was ist da interessant?

# Zugriff auf Matrix data



- **Matrix data ist im Speicher transponiert**
  - Angrenzende Speicherzellen halten (i.A.) die gleiche Koordinate von zwei verschiedenen Samples – warum?
- distance = distance + (**data[dim][i]** – centroids[j][dim])^2
  - Diese Anweisung wird <u>gleichzeitig</u> für eine ganze Warp ausgeführt, mit i-Werten: q, q+1,q+2, …, q+31
  - D.h. Wir bekommen einen „coalesced"-Zugriff auf data!

# Referenz und Ergebnisse

- Kapitel 5 aus „Scaling Up Machine Learning"
  - Meichun Hsu, Ren Wu and Bin Zhang: „Uniformly Fine-Grained Data-Parallel Computing for Machine Learning Algorithms" (Achtung: viele grobe Fehler!)

| Datensatz | | | | Zeiten (s) | | | Speedups (CPU vs. GPU) | |
|---|---|---|---|---|---|---|---|---|
| $N$ | $D$ | $k$ | $M$ | MineBench | HPLC | HPLG | MineBench | HPLC |
| 2M | 2 | 100 | 50 | 154 | 36 | 1.45 | 106 | 25 |
| 2M | 2 | 400 | 50 | 563 | 118 | 2.16 | 261 | 55 |
| 2M | 8 | 100 | 50 | 314 | 99 | 2.48 | 127 | 40 |
| 2M | 8 | 400 | 50 | 1214 | 354 | 4.53 | 268 | 78 |
| 4M | 2 | 100 | 50 | 308 | 73 | 2.88 | 107 | 25 |
| 4M | 2 | 400 | 50 | 1128 | 236 | 4.36 | 259 | 54 |
| 4M | 8 | 100 | 50 | 629 | 197 | 4.95 | 127 | 40 |
| 4M | 8 | 400 | 50 | 2429 | 709 | 9.03 | 269 | 79 |

- MineBench = open source kmeans; HPLC = optimiertes kmeans auf CPU; HPLG = kmeans auf GPU
- CPU: Xeon 5345@ 2.33 GHz; GPU = GeForce GTX 280 (1GB)