GPU Teaching Kit

Accelerated Computing

ILLINOIS

# Memory and Data Locality

CUDA Memories

# Objective

– To learn to effectively use the CUDA memory types in a parallel program
  – Importance of memory access efficiency
  – Registers, shared memory, global memory
  – Scope and lifetime

# Review: Image Blur Kernel.

```
// Get the average of the surrounding 2xBLUR_SIZE x 2xBLUR_SIZE box
for(int blurRow = -BLUR_SIZE; blurRow < BLUR_SIZE+1; ++blurRow) {
    for(int blurCol = -BLUR_SIZE; blurCol < BLUR_SIZE+1; ++blurCol) {

        int curRow = Row + blurRow;
        int curCol = Col + blurCol;
        // Verify we have a valid image pixel
        if(curRow > -1 && curRow < h && curCol > -1 && curCol < w) {
            pixVal += in[curRow * w + curCol];
            pixels++; // Keep track of number of pixels in the accumu
        }
    }
}

// Write our new pixel value out
out[Row * w + Col] = (unsigned char)(pixVal / pixels);
```
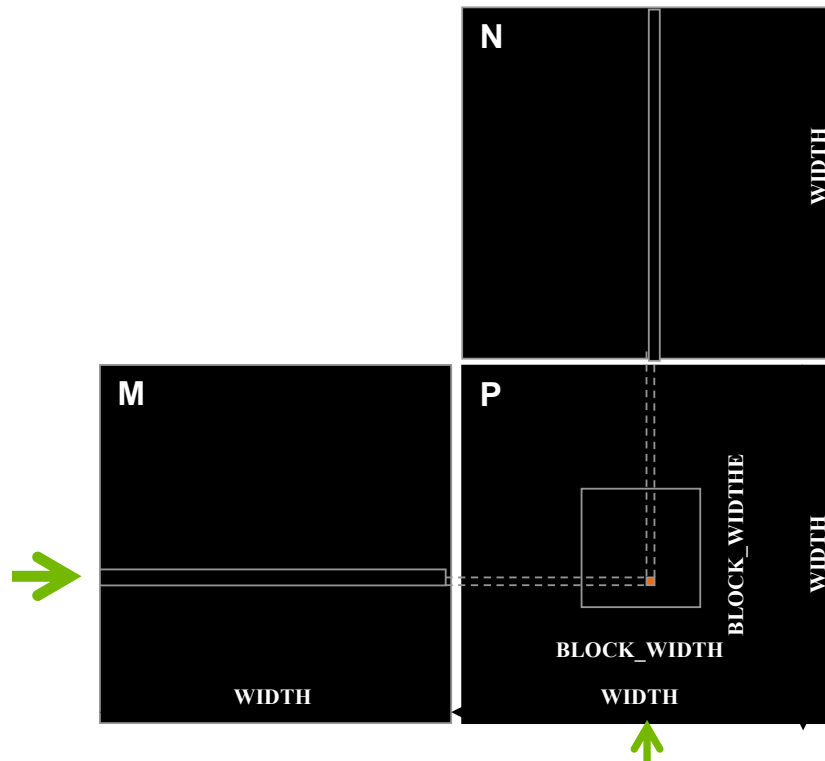
# How about performance on a GPU

- All threads access global memory for their input matrix elements
  - One memory accesses (4 bytes) per floating-point addition
  - 4B/s of memory bandwidth/FLOPS
- Assume a GPU with
  - Peak floating-point rate 1,600 GFLOPS with 600 GB/s DRAM bandwidth
  - 4*1,600 = 6,400 GB/s required to achieve peak FLOPS rating
  - The 600 GB/s memory bandwidth limits the execution at 150 GFLOPS

- This limits the execution rate to 9.3% (150/1600) of the peak floating-point execution rate of the device!

- Need to drastically cut down memory accesses to get close to the1,600 GFLOPS

# Example – Matrix Multiplication

# A Basic Matrix Multiplication

```
__global__ void MatrixMulKernel(float* M, float* N, float* P, int Width) {

  // Calculate the row index of the P element and M
  int Row = blockIdx.y*blockDim.y+threadIdx.y;

  // Calculate the column index of P and N
  int Col = blockIdx.x*blockDim.x+threadIdx.x;

  if ((Row < Width) && (Col < Width)) {
    float Pvalue = 0;
    // each thread computes one element of the block sub-matrix
    for (int k = 0; k < Width; ++k) {
      Pvalue += M[Row*Width+k]*N[k*Width+Col];
    }
    P[Row*Width+Col] = Pvalue;
  }

}
```
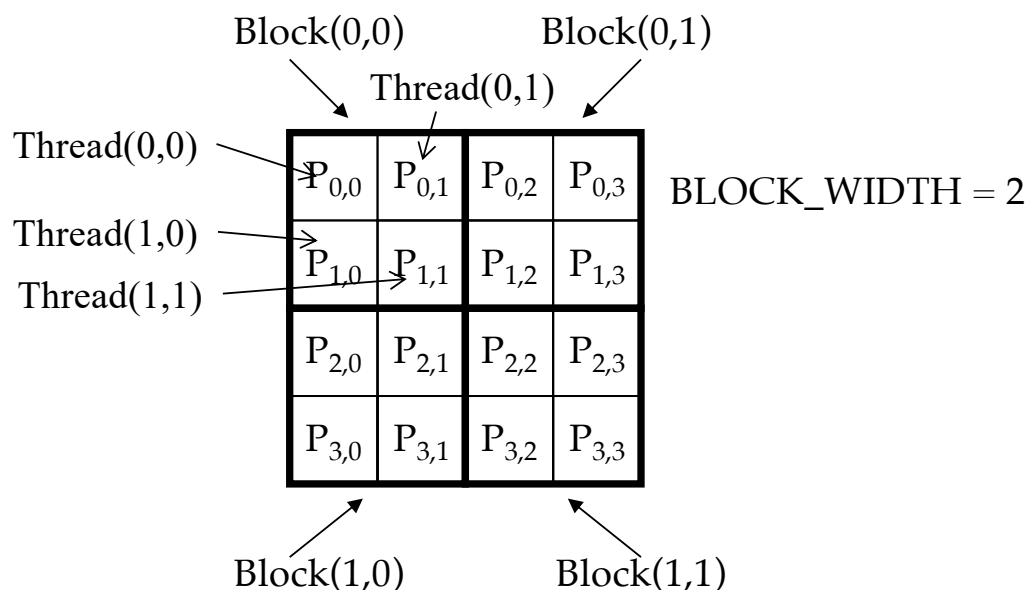
# Example – Matrix Multiplication

```
__global__ void MatrixMulKernel(float* M, float* N, float* P, int Width) {

  // Calculate the row index of the P element and M
  int Row = blockIdx.y*blockDim.y+threadIdx.y;

  // Calculate the column index of P and N
  int Col = blockIdx.x*blockDim.x+threadIdx.x;

  if ((Row < Width) && (Col < Width)) {
    float Pvalue = 0;
    // each thread computes one element of the block sub-matrix
    for (int k = 0; k < Width; ++k) {
      Pvalue += M[Row*Width+k]*N[k*Width+Col];
    }
    P[Row*Width+Col] = Pvalue;
  }

}
```

# A Toy Example: Thread to P Data Mapping



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

Thread(0,1)

Thread(0,0)

Thread(1,0)

Thread(1,1)

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

$BLOCK\_WIDTH = 2$

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

# Calculation of $P_{0,0}$ and $P_{0,1}$

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

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

| | | | |
|---|---|---|---|
| $P_{0,0}$ | $P_{0,1}$ | | |
| $P_{1,0}$ | $P_{1,1}$ | | |
| | | | |
| | | | |

# Memory and Registers in the Von-Neumann Model



Memory

I/O

Processing Unit

Reg File

ALU

Control Unit

PC          IR

# Programmer View of CUDA Memories

# Declaring CUDA Variables

| Variable declaration | Memory | Scope | Lifetime |
|---|---|---|---|
| int LocalVar; | register | thread | thread |
| __device__ __shared__ int SharedVar; | shared | block | block |
| __device__ int GlobalVar; | global | grid | application |
| __device__ __constant__ int ConstantVar; | constant | grid | application |

- __device__ is optional when used with __shared__, or __constant__
- Automatic variables reside in a register
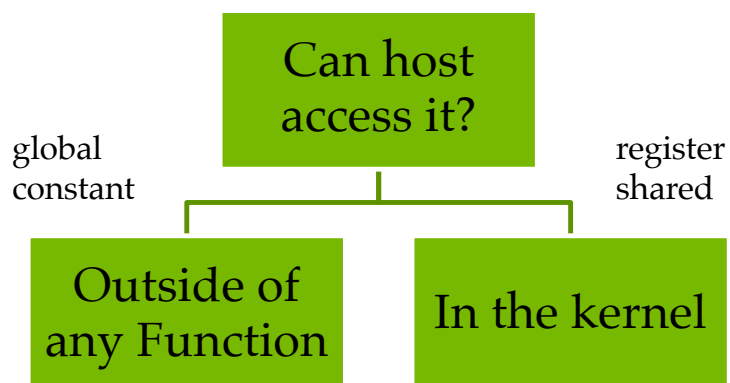  - Except per-thread arrays that reside in global memory

# Example:
# Shared Memory Variable Declaration

```
void blurKernel(unsigned char * in, unsigned char * out,
int w, int h)
{
    __shared__  float ds_in[TILE_WIDTH][TILE_WIDTH];

 …
}
```

# Where to Declare Variables?

global
constant

register
shared

Can host
access it?

Outside of
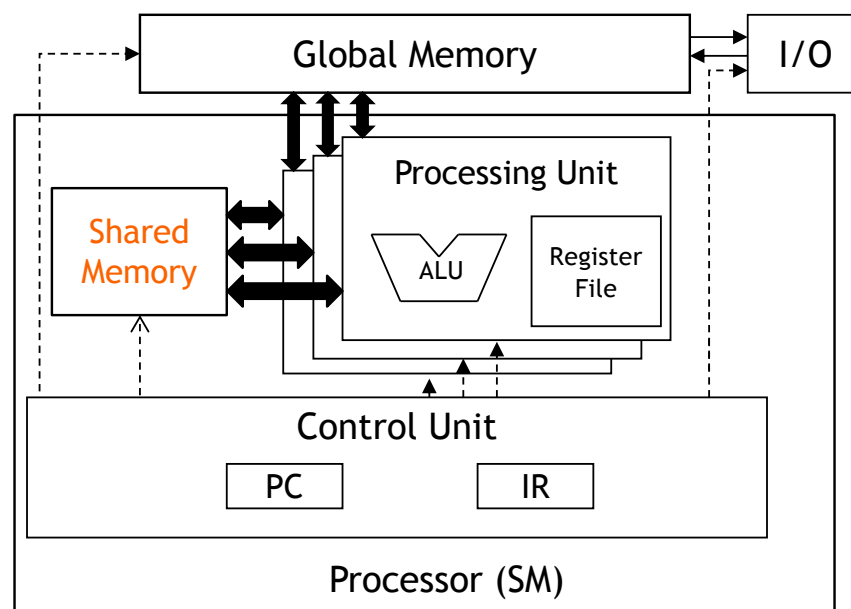any Function

In the kernel

# Shared Memory in CUDA

– A special type of memory whose contents are explicitly defined and used in the kernel source code
  – One in each SM
  – Accessed at much higher speed (in both latency and throughput) than global memory
  – Scope of access and sharing - thread blocks
  – Lifetime – thread block, contents will disappear after the corresponding thread finishes terminates execution
  – Accessed by memory load/store instructions
  – A form of scratchpad memory in computer architecture

# Hardware View of CUDA Memories

GPU Teaching Kit

Accelerated Computing

Module 4.2 – Memory and Data Locality

Tiled Parallel Algorithms

## Objective

– To understand the motivation and ideas for tiled parallel algorithms
  – Reducing the limiting effect of memory bandwidth on parallel kernel performance
  – Tiled algorithms and barrier synchronization

# Global Memory Access Pattern
## of the Basic Matrix Multiplication Kernel

Global Memory

# Tiling/Blocking - Basic Idea

Global Memory

On-chip Memory



Divide the global memory content into tiles

Focus the computation of threads on one or a small number of tiles at each point in time

# Tiling/Blocking - Basic Idea

Global Memory

On-chip Memory

Thread 1

Thread 2

• • •

# Basic Concept of Tiling

– In a congested traffic system, significant reduction of vehicles can greatly improve the delay seen by all vehicles
  – Carpooling for commuters
  – Tiling for global memory accesses
    – drivers = threads accessing their memory data operands
    – cars = memory access requests

## Some Computations are More Challenging to Tile

- Some carpools may be easier than others
    - Car pool participants need to have similar work schedule
    - Some vehicles may be more suitable for carpooling
- Similar challenges exist in tiling

---

## Carpools need synchronization.

- Good: when people have similar schedule

| Worker A | sleep | work | dinner |
|----------|-------|------|--------|
| Time | | | |
| Worker B | sleep | work | dinner |

# Carpools need synchronization.

– Bad: when people have very different schedule

Worker A    party            sleep          work

time

Worker B       sleep      work         dinner

---

# Same with Tiling

– Good: when threads have similar access timing

Thread 1

Time

Thread 2

• • •

Thread 1

Time

Thread 2

– Bad: when threads have very different timing

# Barrier Synchronization for Tiling



Thread 0
Thread 1
Thread 2
Thread 3
Thread 4

Time

Thread N-3
Thread N-2
Thread N-1

---

# Outline of Tiling Technique

– Identify a tile of global memory contents that are accessed by multiple threads

– Load the tile from global memory into on-chip memory

– Use barrier synchronization to make sure that all threads are ready to start the phase

– Have the multiple threads to access their data from the on-chip memory

– Use barrier synchronization to make sure that all threads have completed the current phase

– Move on to the next tile

GPU Teaching Kit

Accelerated Computing

ILLINOIS

# Memory Model and Locality

Tiled Matrix Multiplication

# Objective

– To understand the design of a tiled parallel algorithm for matrix multiplication
  – Loading a tile
  – Phased execution
  – Barrier Synchronization

# Matrix Multiplication

– Data access pattern
  – Each thread - a row of M and a column of N
  – Each thread block – a strip of M and a strip of N

# Tiled Matrix Multiplication

– Break up the execution of each thread into phases
– so that the data accesses by the thread block in each phase are focused on one tile of M and one tile of N
– The tile is of BLOCK_SIZE elements in each dimension

# Loading a Tile

- All threads in a block participate
  - Each thread loads one M element and one N element in tiled code

# Phase 0 Load for Block (0,0)

# Phase 0 Use for Block (0,0) (iteration 0)

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

Shared Memory

| $N_{0,0}$ | $N_{0,1}$ |
|---|---|
| $N_{1,0}$ | $N_{1,1}$ |

Shared Memory

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

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

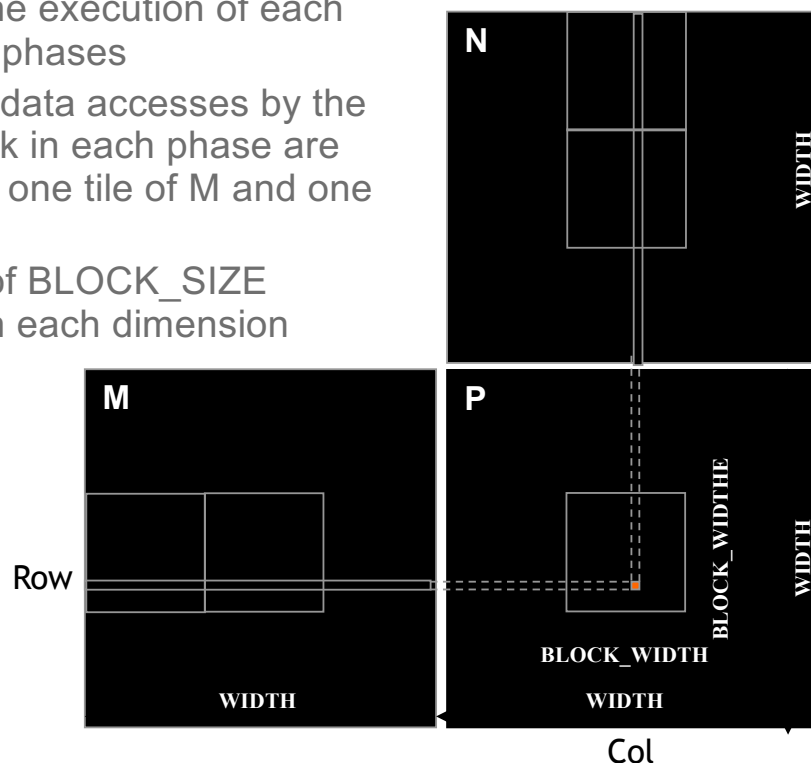| $P_{0,0}$ | $P_{0,1}$ | $P_{0,2}$ | $P_{0,3}$ |
|---|---|---|---|
| $P_{1,0}$ | $P_{1,1}$ | $P_{1,2}$ | $P_{1,3}$ |
| $P_{2,0}$ | $P_{2,1}$ | $P_{2,2}$ | $P_{2,3}$ |
| $P_{3,0}$ | $P_{3,1}$ | $P_{3,2}$ | $P_{3,3}$ |

# Phase 0 Use for Block (0,0) (iteration 1)

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

Shared Memory

| $N_{0,0}$ | $N_{0,1}$ |
|---|---|
| $N_{1,0}$ | $N_{1,1}$ |

Shared Memory

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

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

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

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

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

Shared Memory

Shared Memory

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

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

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

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

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

Shared Memory

Shared Memory

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

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

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

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

$N_{2,0}$ $N_{2,1}$
$N_{3,0}$ $N_{3,1}$   Shared Memory

Shared Memory

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

# Execution Phases of Toy Example

| | Phase 0 | | | Phase 1 | | |
|---|---|---|---|---|---|---|
| thread$_{0,0}$ | $M_{0,0}$ ↓ $Mds_{0,0}$ | $N_{0,0}$ ↓ $Nds_{0,0}$ | PValue$_{0,0}$ += $Mds_{0,0}*Nds_{0,0}$ + $Mds_{0,1}*Nds_{1,0}$ | $M_{0,2}$ ↓ $Mds_{0,0}$ | $N_{2,0}$ ↓ $Nds_{0,0}$ | PValue$_{0,0}$ += $Mds_{0,0}*Nds_{0,0}$ + $Mds_{0,1}*Nds_{1,0}$ |
| thread$_{0,1}$ | $M_{0,1}$ ↓ $Mds_{0,1}$ | $N_{0,1}$ ↓ $Nds_{1,0}$ | PValue$_{0,1}$ += $Mds_{0,0}*Nds_{0,1}$ + $Mds_{0,1}*Nds_{1,1}$ | $M_{0,3}$ ↓ $Mds_{0,1}$ | $N_{2,1}$ ↓ $Nds_{0,1}$ | PValue$_{0,1}$ += $Mds_{0,0}*Nds_{0,1}$ + $Mds_{0,1}*Nds_{1,1}$ |
| thread$_{1,0}$ | $M_{1,0}$ ↓ $Mds_{1,0}$ | $N_{1,0}$ ↓ $Nds_{1,0}$ | PValue$_{1,0}$ += $Mds_{1,0}*Nds_{0,0}$ + $Mds_{1,1}*Nds_{1,0}$ | $M_{1,2}$ ↓ $Mds_{1,0}$ | $N_{3,0}$ ↓ $Nds_{1,0}$ | PValue$_{1,0}$ += $Mds_{1,0}*Nds_{0,0}$ + $Mds_{1,1}*Nds_{1,0}$ |
| thread$_{1,1}$ | $M_{1,1}$ ↓ $Mds_{1,1}$ | $N_{1,1}$ ↓ $Nds_{1,1}$ | PValue$_{1,1}$ += $Mds_{1,0}*Nds_{0,1}$ + $Mds_{1,1}*Nds_{1,1}$ | $M_{1,3}$ ↓ $Mds_{1,1}$ | $N_{3,1}$ ↓ $Nds_{1,1}$ | PValue$_{1,1}$ += $Mds_{1,0}*Nds_{0,1}$ + $Mds_{1,1}*Nds_{1,1}$ |

time →

# Execution Phases of Toy Example (cont.)

| | Phase 0 | | | Phase 1 | | |
|---|---|---|---|---|---|---|
| thread$_{0,0}$ | $M_{0,0}$ ↓ Mds$_{0,0}$ | $N_{0,0}$ ↓ Nds$_{0,0}$ | PValue$_{0,0}$ += Mds$_{0,0}$*Nds$_{0,0}$+ Mds$_{0,1}$*Nds$_{1,0}$ | $M_{0,2}$ ↓ Mds$_{0,0}$ | $N_{2,0}$ ↓ Nds$_{0,0}$ | PValue$_{0,0}$ += Mds$_{0,0}$*Nds$_{0,0}$+ Mds$_{0,1}$*Nds$_{1,0}$ |
| thread$_{0,1}$ | $M_{0,1}$ ↓ Mds$_{0,1}$ | $N_{0,1}$ ↓ Nds$_{1,0}$ | PValue$_{0,1}$ += Mds$_{0,0}$*Nds$_{0,1}$+ Mds$_{0,1}$*Nds$_{1,1}$ | $M_{0,3}$ ↓ Mds$_{0,1}$ | $N_{2,1}$ ↓ Nds$_{0,1}$ | PValue$_{0,1}$ += Mds$_{0,0}$*Nds$_{0,1}$+ Mds$_{0,1}$*Nds$_{1,1}$ |
| thread$_{1,0}$ | $M_{1,0}$ ↓ Mds$_{1,0}$ | $N_{1,0}$ ↓ Nds$_{1,0}$ | PValue$_{1,0}$ += Mds$_{1,0}$*Nds$_{0,0}$+ Mds$_{1,1}$*Nds$_{1,0}$ | $M_{1,2}$ ↓ Mds$_{1,0}$ | $N_{3,0}$ ↓ Nds$_{1,0}$ | PValue$_{1,0}$ += Mds$_{1,0}$*Nds$_{0,0}$+ Mds$_{1,1}$*Nds$_{1,0}$ |
| thread$_{1,1}$ | $M_{1,1}$ ↓ Mds$_{1,1}$ | $N_{1,1}$ ↓ Nds$_{1,1}$ | PValue$_{1,1}$ += Mds$_{1,0}$*Nds$_{0,1}$+ Mds$_{1,1}$*Nds$_{1,1}$ | $M_{1,3}$ ↓ Mds$_{1,1}$ | $N_{3,1}$ ↓ Nds$_{1,1}$ | PValue$_{1,1}$ += Mds$_{1,0}$*Nds$_{0,1}$+ Mds$_{1,1}$*Nds$_{1,1}$ |

time →

Shared memory allows each value to be accessed by multiple threads

---

# Barrier Synchronization

– Synchronize all threads in a block
  – __syncthreads()

– All threads in the same block must reach the __syncthreads() before any of the them can move on

– Best used to coordinate the phased execution tiled algorithms
  – To ensure that all elements of a tile are loaded at the beginning of a phase
  – To ensure that all elements of a tile are consumed at the end of a phase

# GPU Teaching Kit

Accelerated Computing