Now, this is just a simulation of what the blocks will look like once they've assembled
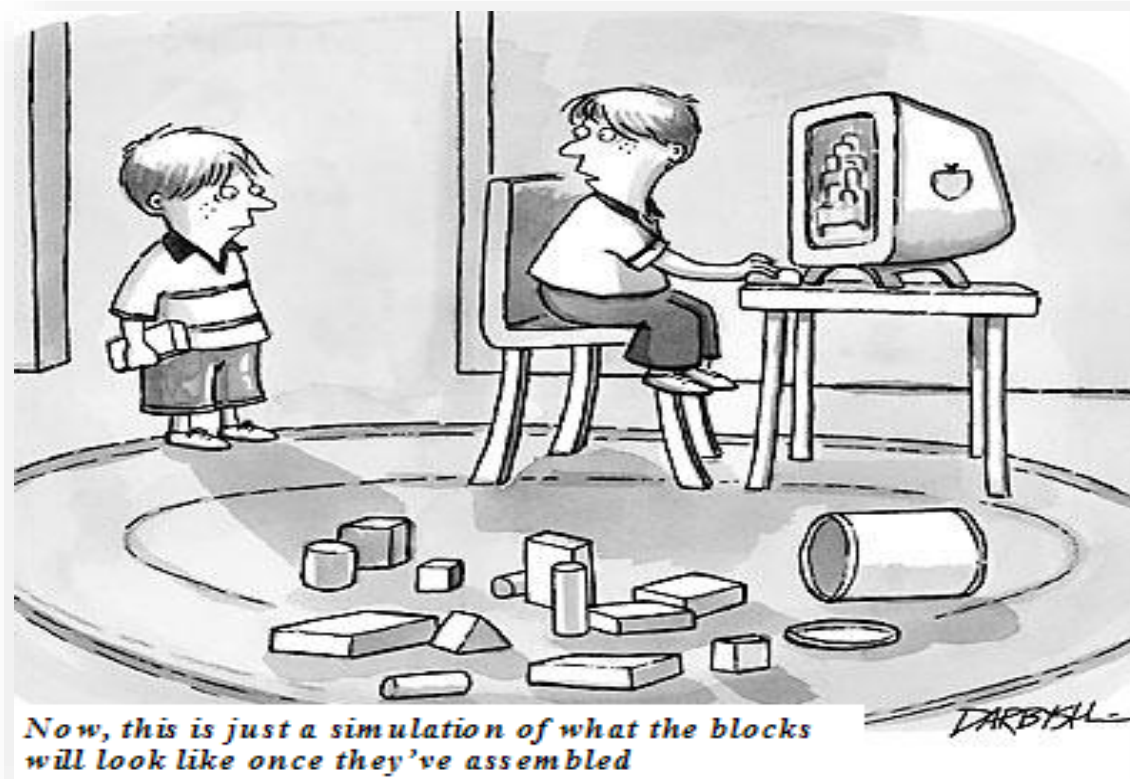
- Tiling for Matrix Multiplication
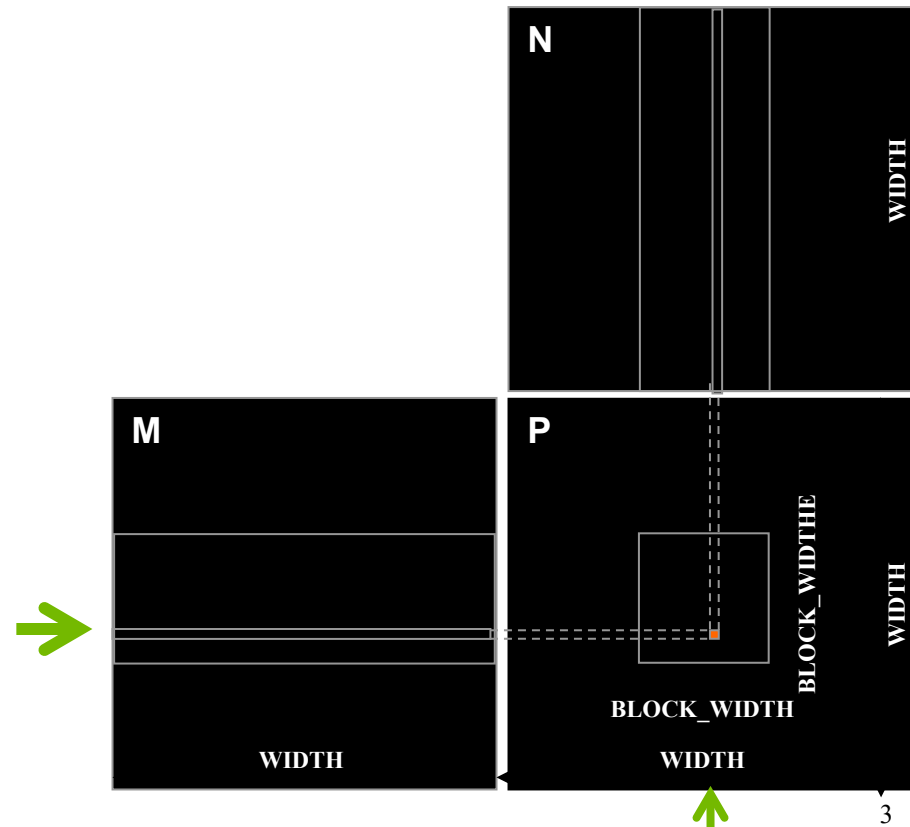
# Outline of Tiling Technique

- 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
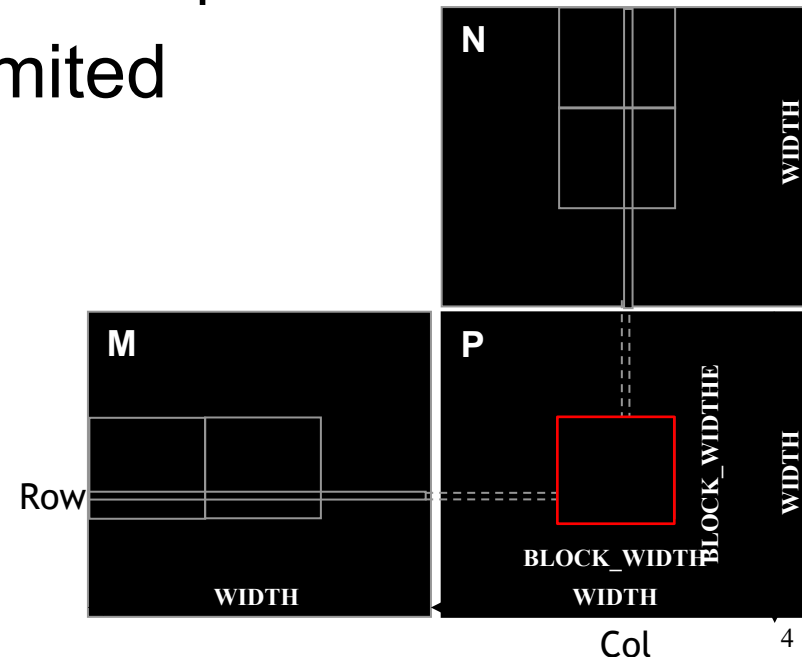
- Move on to the next tile

# 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

- Divide the M and N matrices into smaller tiles.
    - Breaks up the execution of each thread into phases focusing on one tile of M and one tile of N
    - Threads collaboratively load subsets of the M and N elements into the shared memory
        - The tile is of BLOCK_SIZE elements in each dimension
    - Threads use the elements in their dot product calculation.
- Size of the shared memory limited
    - Do not exceed the capacity
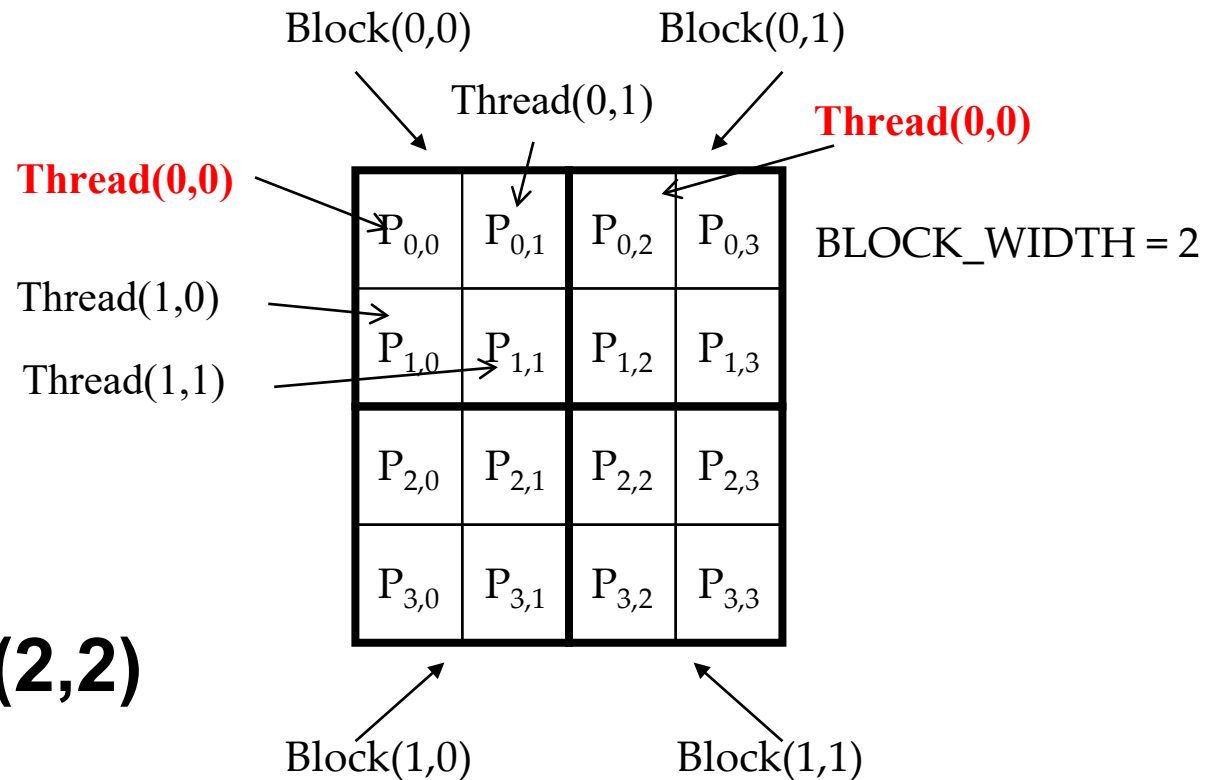    - Choose tile size carefully

# A Toy Example: Thread to P[4x4] Data Mapping

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

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

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

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

Thread(0,1)

**Thread(0,0)**

**Thread(0,0)**

BLOCK_WIDTH = 2

Thread(1,0)

Thread(1,1)

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

- **Grid(2,2), Block(2,2)**

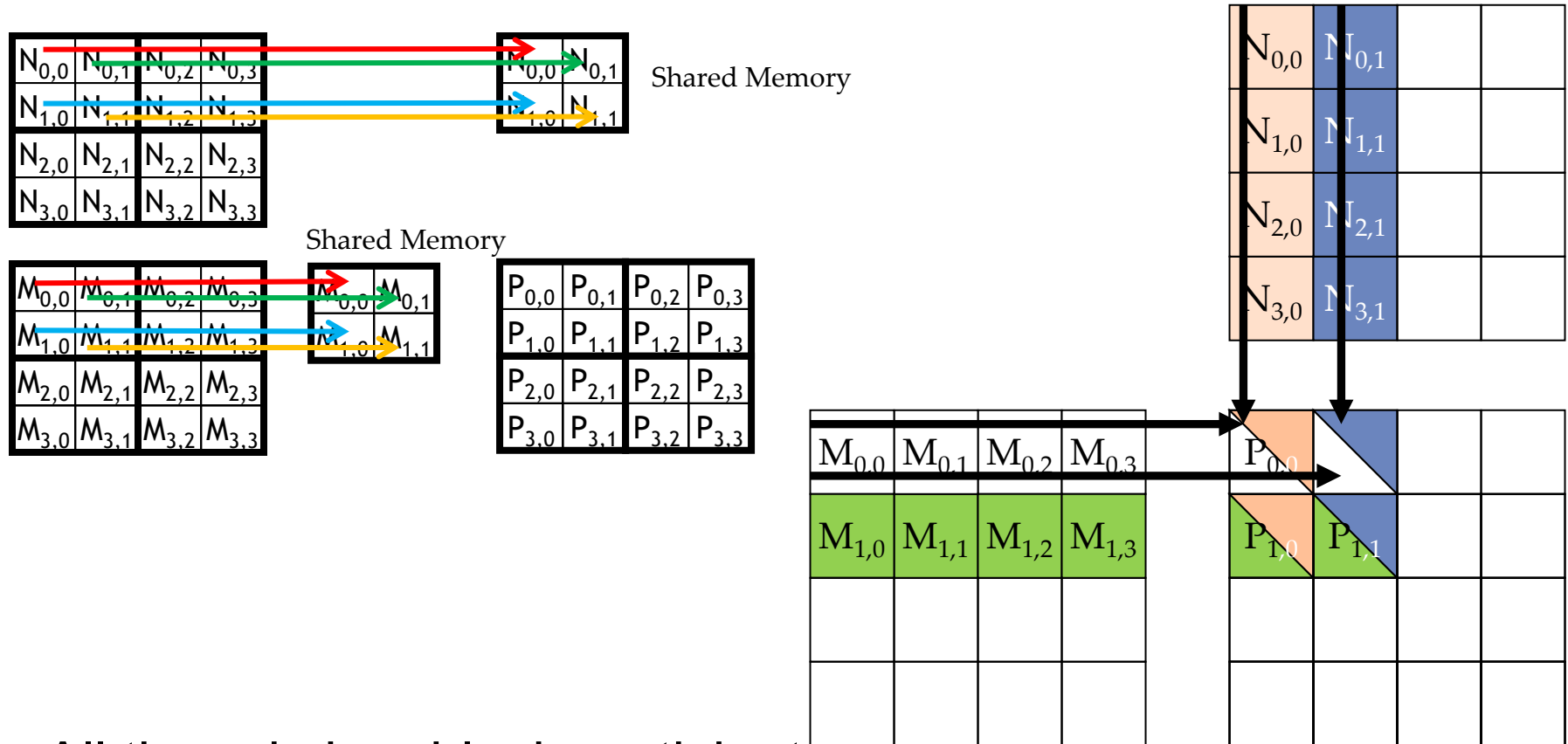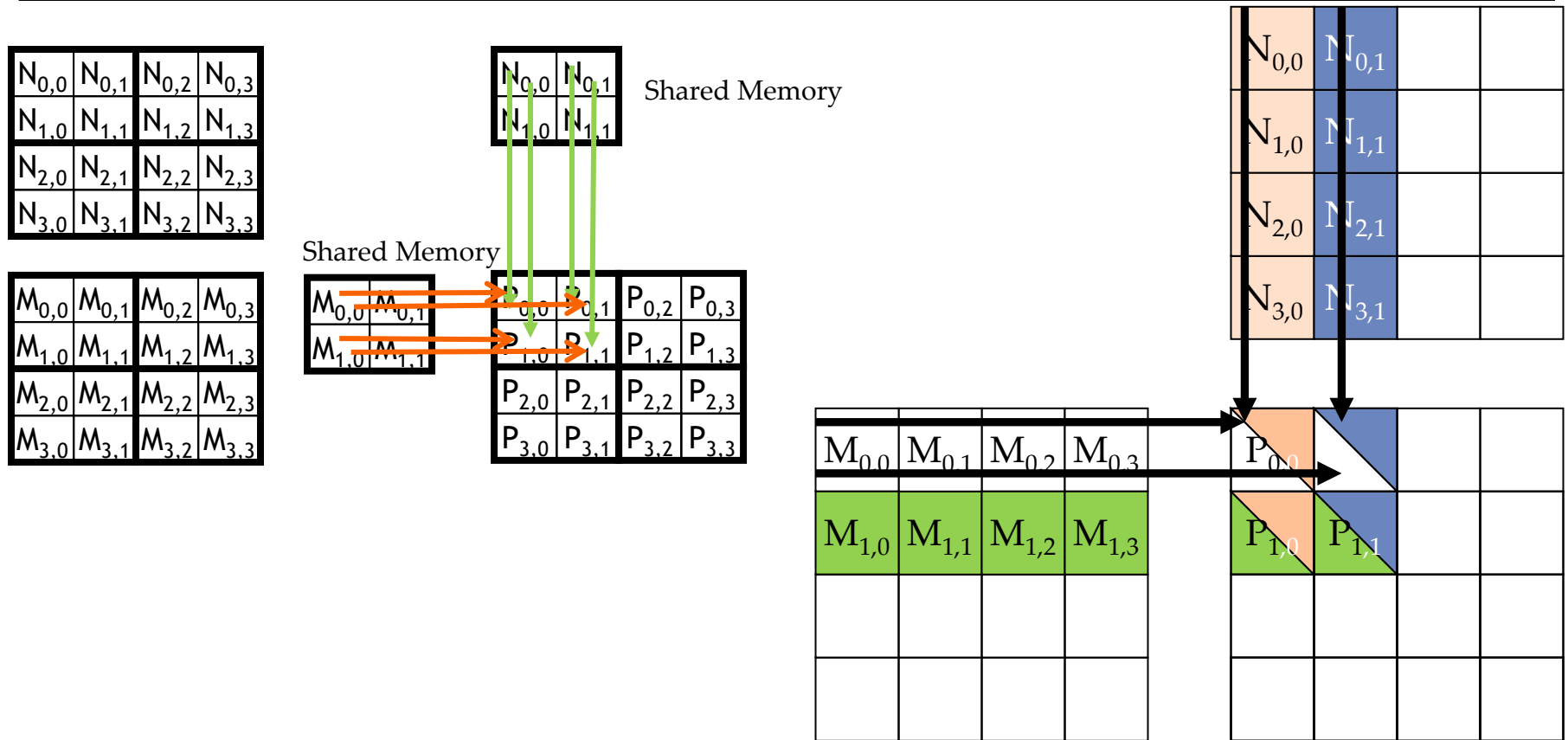Block(1,0)  Block(1,1)

# Phase 0 Load for Block (0,0)



- All threads in a block participate
- Each thread loads one M element and one N element in tiled code
  - four threads of block0,0 collaboratively load a tile of M and N into shared memory

# Phase 0 Load for Block (0,0)



- All threads in a block participate
- Each thread loads one M element and one N element in tiled code
    - four threads of block0,0 collaboratively load a tile of M and N into shared memory
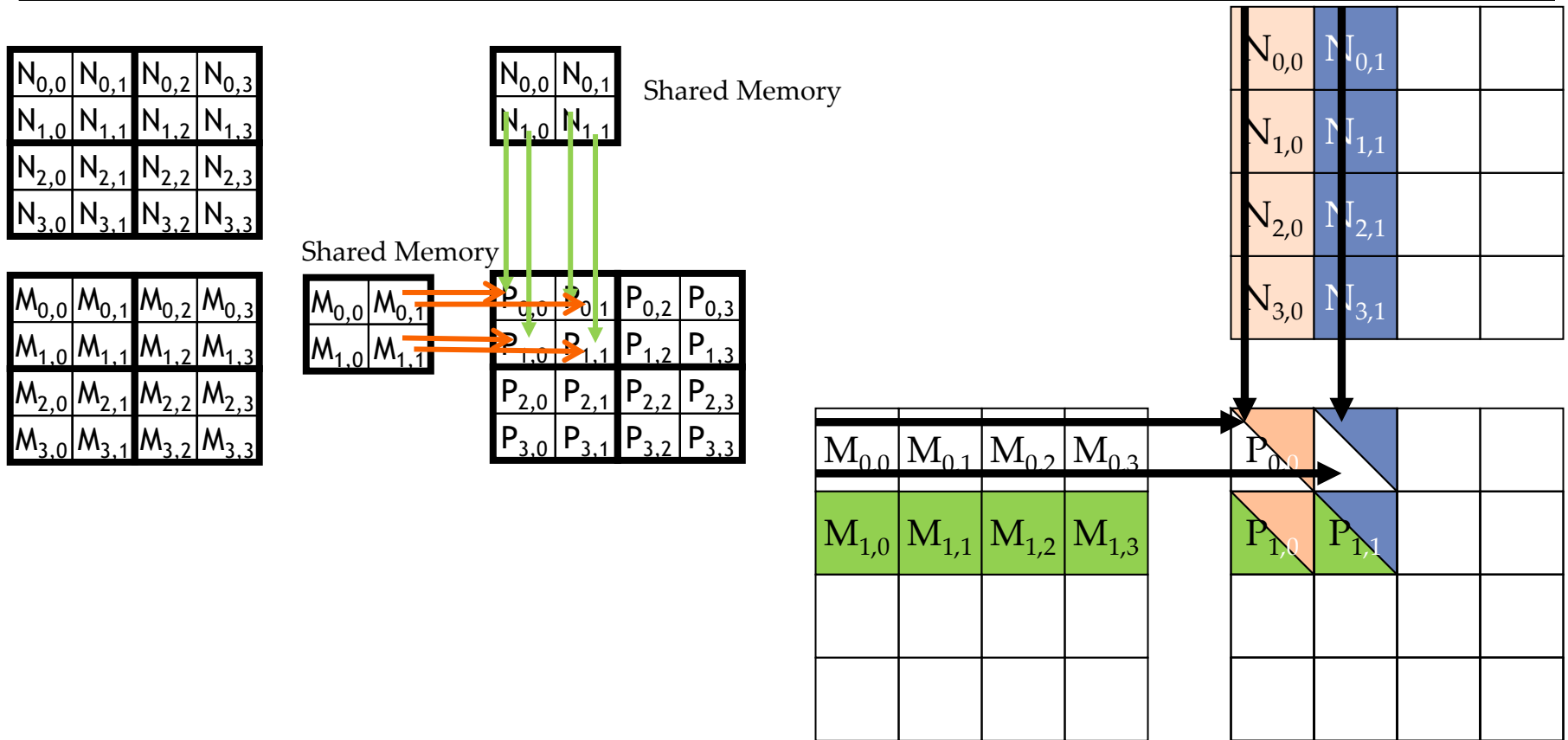
# Phase 0 Use for Block (0,0) (iteration 0)



– After tiles of M and N are in the shared memory, these elements are used in the calculation of the dot product
  - Two iterations
– Note that each value in the shared memory is used twice.
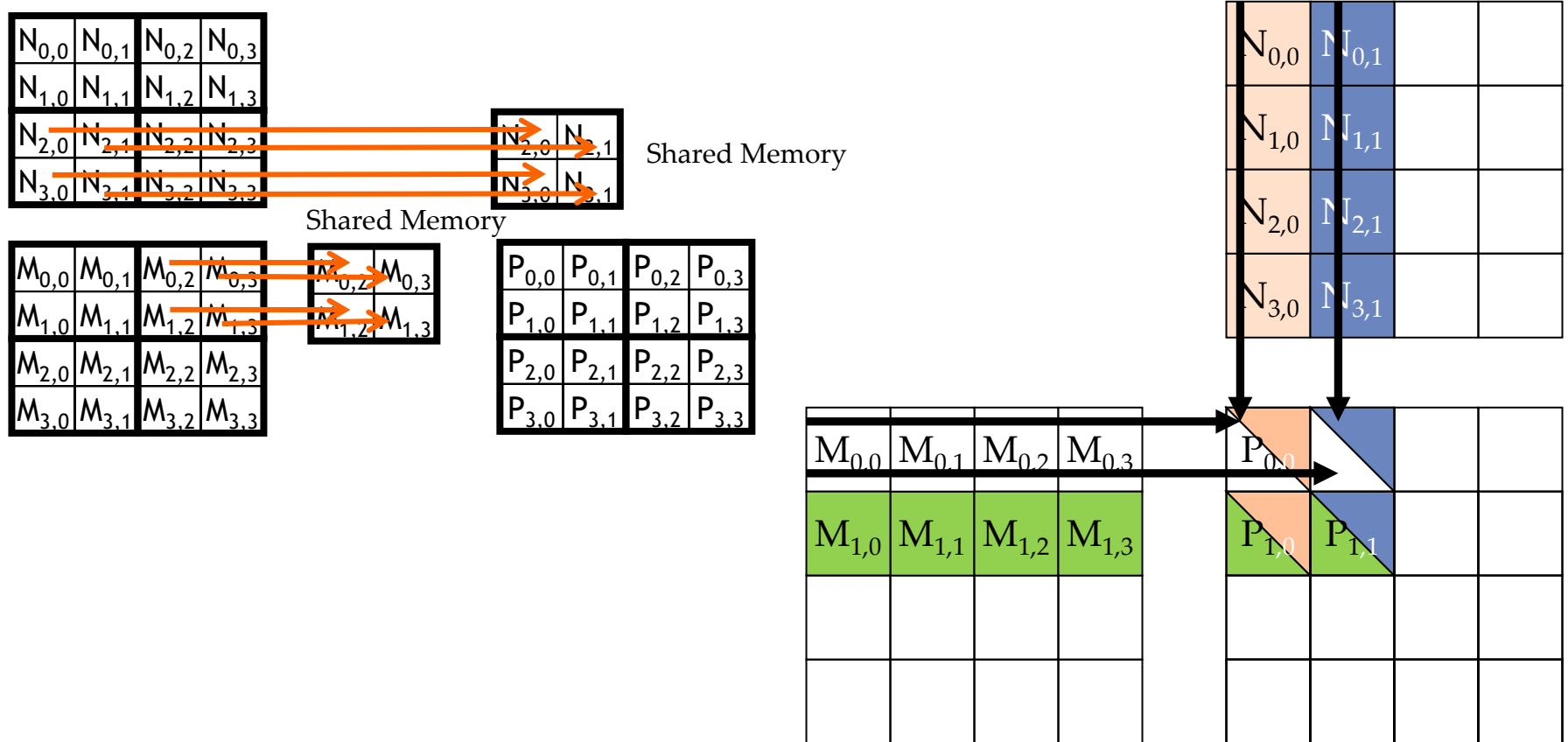  - Reduces global memory access by half

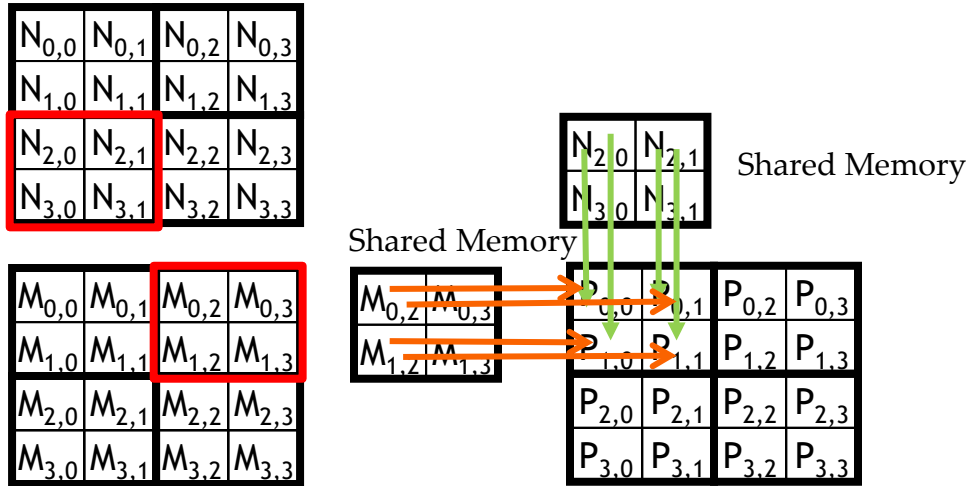# Phase 0 Use for Block (0,0) (iteration 1)

Shared Memory

Shared Memory

Partial accumulated dot product  is a private value generated for each thread

# Phase 1 Load for Block (0,0)

Shared Memory

Shared Memory
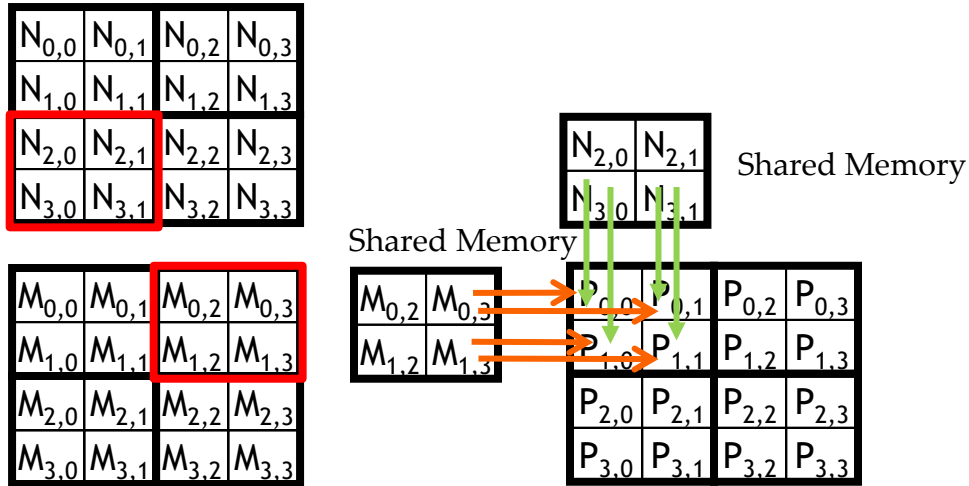
Shared Memory

Shared Memory

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



Shared Memory

Shared Memory

# 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 →

Mds/Nds: shared memory array for the M/N elements.

**Shared memory allows each value to be accessed by multiple threads**

# Memory Bandwidth Utilization

What is the throughput achieved when tile size is set to 16x16 for a square matrix multiplication assuming that memory GPU device delivers 720GB/sec memory bandwidth and peak throughput of 9300GFlops?

*Hint:* You need to consider number of global memory accesses required in each phase and number of floating point operations are executed per tile.

- ❏ 360 GFLOPS
- ❏ 720GFLOPS
- ❏ 1440 GFLOPS
- ❏ 2880 GFLOPS

# Next

- **CUDA implementation of Tiling**