

ECE408/CS483/CSE408 Spring 2019

Applied Parallel Programming

Lecture 22: Performance Considerations

© David Kirk/NVIDIA and Wen-mei Hwu, 2007-2018
ECE408/CS483/ECE498a, University of Illinois, Urbana-Champaign

1

Objective

- To learn more about CUDA memory coalescing
- To learn more about barrier synchronization

© David Kirk/NVIDIA and Wen-mei Hwu, 2007-2018
ECE408/CS483/ECE498a, University of Illinois, Urbana-Champaign

2

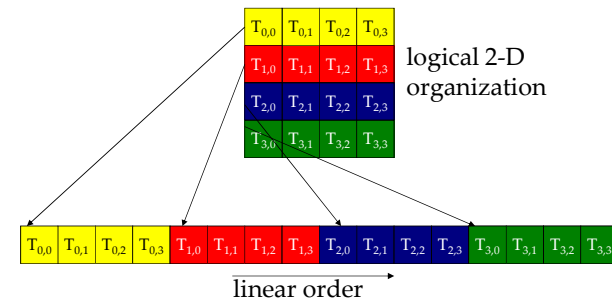
Memory Coalescing

- After row-major linearization, adjacent threads in each warp access adjacent memory locations
 - Accesses for threads in the warp consolidated into one (or two) DRAM access (burst)
- The point is efficiency in using DRAM bursts
 - Current DRAM burst size is 64-128 bytes
 - Any distance between locations accessed by adjacent threads in a warp reduces DRAM efficiency
- When adjacent threads in a warp access identical locations, accesses for threads in the warp consolidated into one DRAM access

© David Kirk/NVIDIA and Wen-mei Hwu, 2007-2018
ECE408/CS483/ECE498a, University of Illinois, Urbana-Champaign

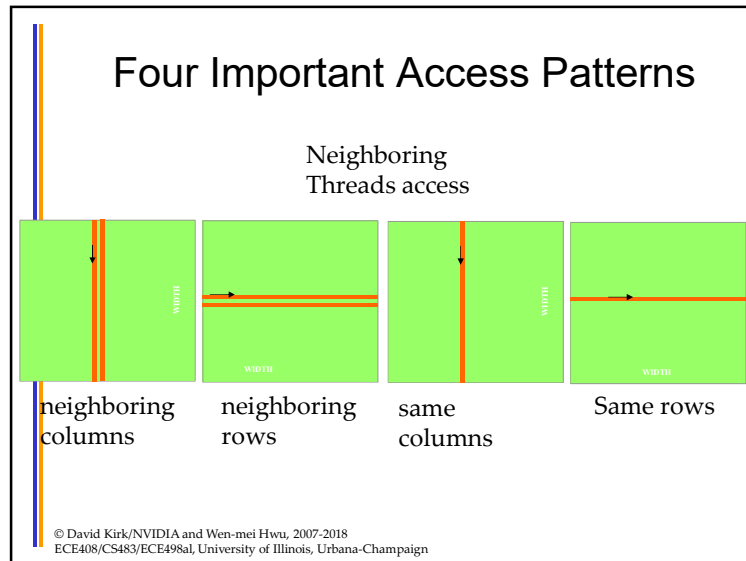
3

2D Threads in Linear Order (Row Major Layout)

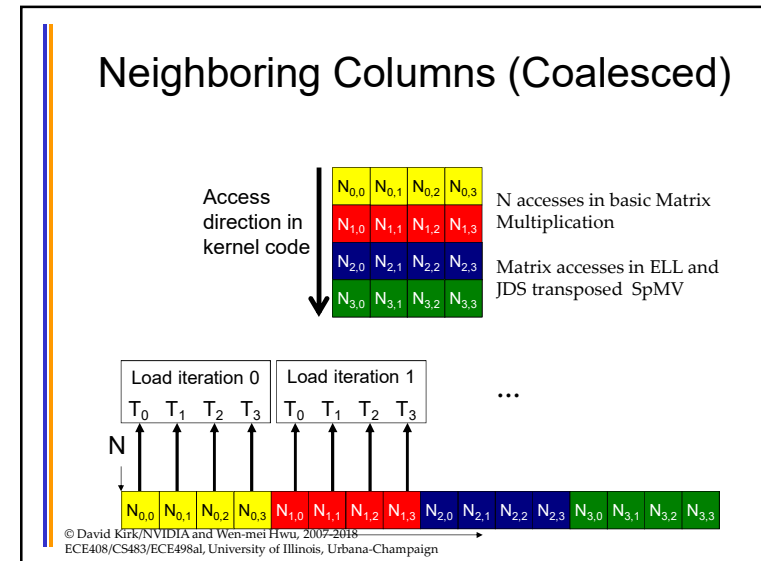


© David Kirk/NVIDIA and Wen-mei Hwu, 2007-2018
ECE408/CS483/ECE498a, University of Illinois, Urbana-Champaign

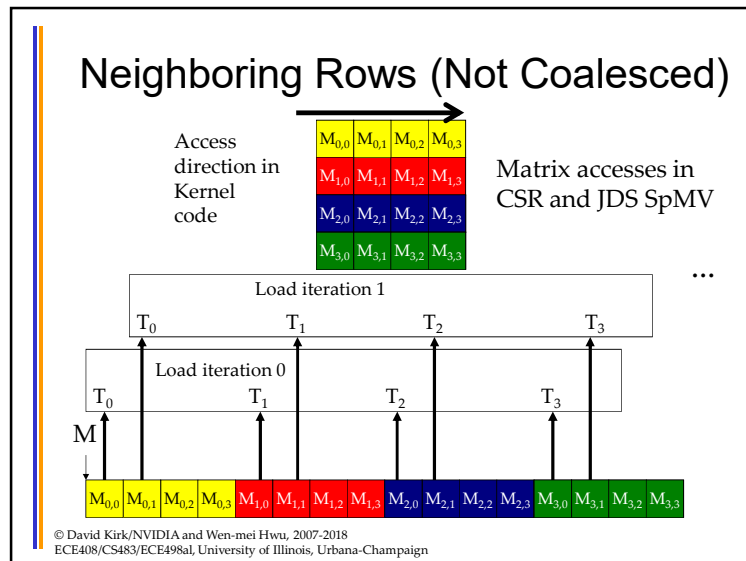
4



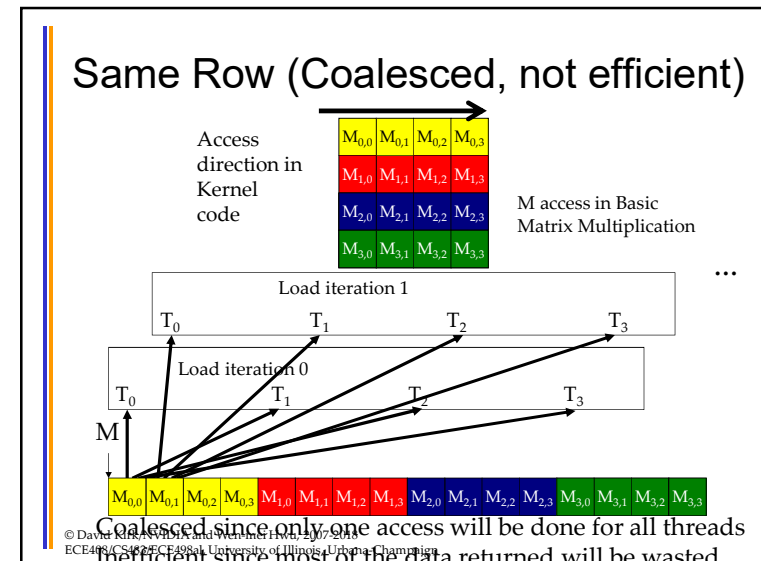
5



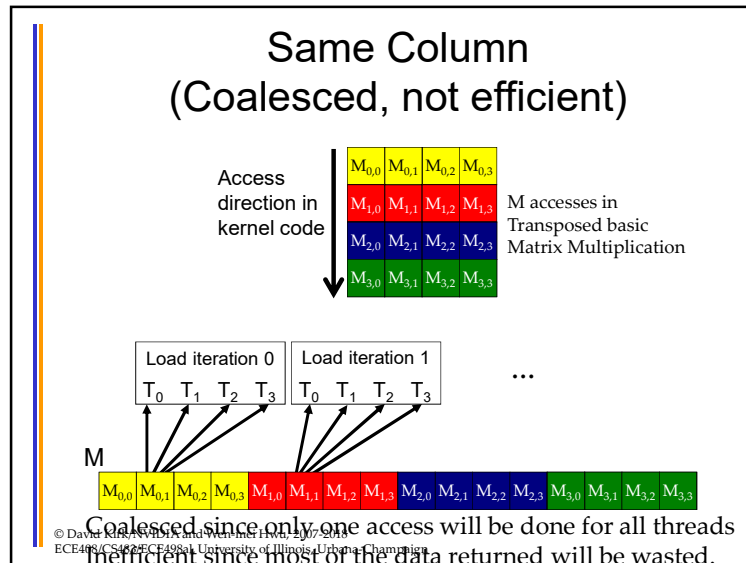
6



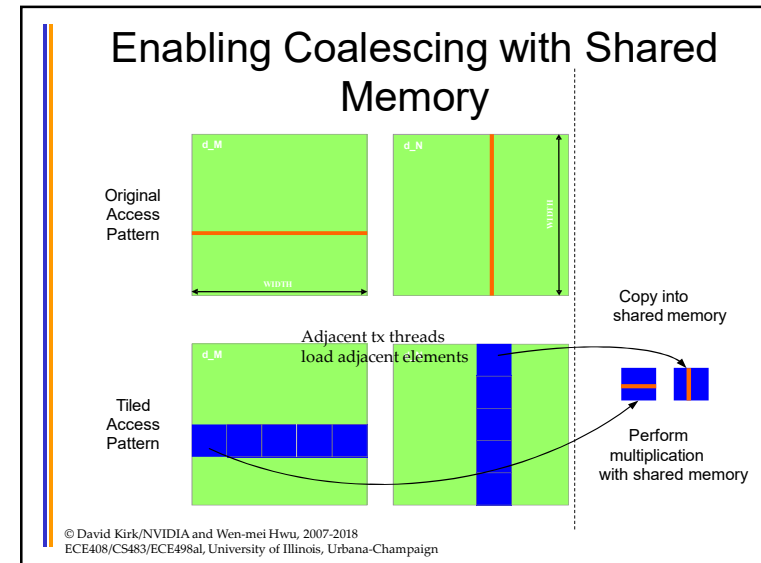
7



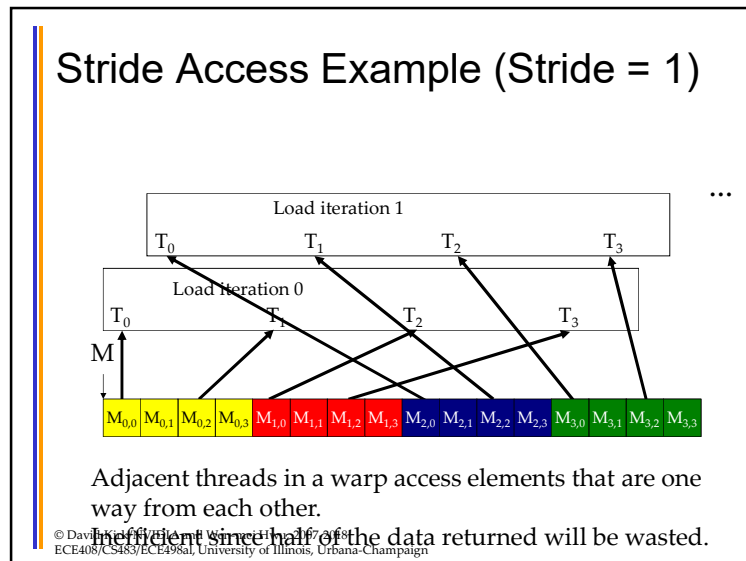
8



9



10



11

- ### Memory Coalescing - Summary
- After row-major linearization, adjacent threads in each warp access adjacent memory locations
 - Accesses for threads in the warp consolidated into one (or two) DRAM access (burst)
 - The point is efficiency in using DRAM bursts
 - Current DRAM burst size is 64-128 bytes
 - Any distance between locations accessed by adjacent threads in a warp reduces DRAM efficiency
 - When adjacent threads in a warp access identical locations, accesses for threads in the warp consolidated into one DRAM access
- © David Kirk/NVIDIA and Wen-mei Hwu, 2007-2018
ECE408/CS483/ECE498a, University of Illinois, Urbana-Champaign

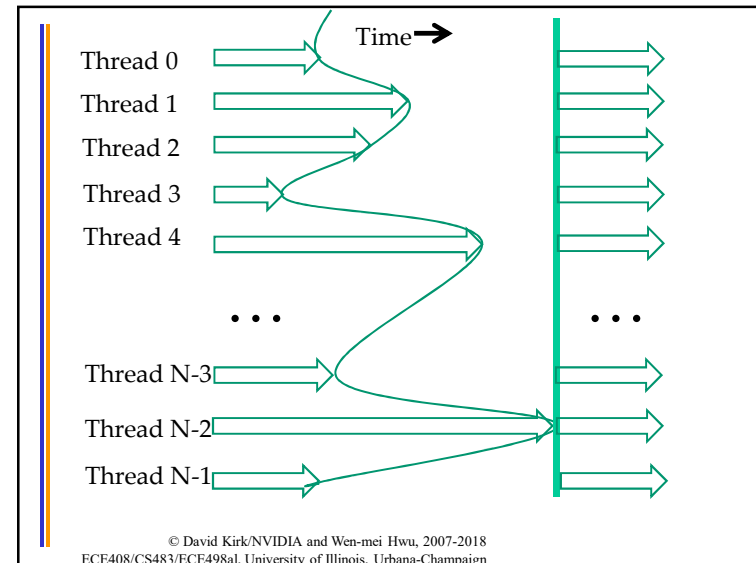
12

Barrier Synchronization

- An API function call in CUDA
 - `__syncthreads()`
- All threads in the same block must reach the `__syncthreads()` before any can move on
- Best used to coordinate tiled algorithms
 - To ensure that all elements of a tile are loaded
 - To ensure that all elements of a tile are consumed

© David Kirk/NVIDIA and Wen-mei Hwu, 2007-2018
ECE408/CS483/ECE498a, University of Illinois, Urbana-Champaign

13



14

The Reduction Steps

```
// Stride is distance to the next value being
// accumulated into the threads mapped position
// in the partialSum[] array
for (unsigned int stride = 1;
     stride <= blockDim.x; stride *= 2)
{
    __syncthreads();
    if (t % stride == 0)
        __syncthreads();    <- Why is this a bad idea?
    partialSum[2*t] += partialSum[2*t+stride];
}
```

© David Kirk/NVIDIA and Wen-mei Hwu, 2007-2018
ECE408/CS483/ECE498a, University of Illinois, Urbana-Champaign

15

READ CHAPTER 5

© David Kirk/NVIDIA and Wen-mei Hwu, 2007-2018
ECE408/CS483/ECE498a, University of Illinois, Urbana-Champaign

16