

ECE408 Spring 2020

Applied Parallel Programming

Lecture 17

Atomic Operations and
Histogramming

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

1

1

Objective

- To understand atomic operations
 - Read-modify-write in parallel computation
 - A primitive form of “critical regions” in parallel programs
 - Use of atomic operations in CUDA
 - Why atomic operations reduce memory system throughput
 - How to avoid atomic operations in some parallel algorithms
- To learn practical histogram programming techniques
 - Basic histogram algorithm using atomic operations
 - Atomic operation throughput
 - Privatization

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

2

2

A Common Collaboration Pattern

- Multiple bank tellers count the total amount of cash in the safe
- Each grab a pile and count
- Have a central display of the running total
- Whenever someone finishes counting a pile, add the subtotal of the pile to the running total
- A bad outcome
 - Some of the piles were not accounted for

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

3

3

A Common Arbitration Pattern

- Multiple customers booking air tickets
- Each
 - Brings up a flight seat map
 - Decides on a seat
 - Update the the seat map, mark the seat as taken
- A bad outcome
 - Multiple passengers ended up booking the same seat

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

4

4

Read-Modify-Write Operations

thread1: Old \leftarrow Mem[x] thread2: Old \leftarrow Mem[x]
 New \leftarrow Old + 1 New \leftarrow Old + 1
 Mem[x] \leftarrow New Mem[x] \leftarrow New

If Mem[x] was **initially 0**, what would the value of Mem[x] be after threads 1 and 2 have completed?

– What does each thread get in their Old variable?

The answer may vary due to data races. To avoid data races, you should use atomic operations

© David Kirk/NVIDIA and Wen-mei W. Hwu
 ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

5

5

Timing Scenario #1

Time	Thread 1	Thread 2
1	(0) Old \leftarrow Mem[x]	
2	(1) New \leftarrow Old + 1	
3	(1) Mem[x] \leftarrow New	
4		(1) Old \leftarrow Mem[x]
5		(2) New \leftarrow Old + 1
6		(2) Mem[x] \leftarrow New

- Thread 1 Old = 0
- Thread 2 Old = 1
- Mem[x] = 2 after the sequence

© David Kirk/NVIDIA and Wen-mei W. Hwu
 ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

6

6

Timing Scenario #2

Time	Thread 1	Thread 2
1		(0) Old \leftarrow Mem[x]
2		(1) New \leftarrow Old + 1
3		(1) Mem[x] \leftarrow New
4	(1) Old \leftarrow Mem[x]	
5	(2) New \leftarrow Old + 1	
6	(2) Mem[x] \leftarrow New	

- Thread 1 Old = 1
- Thread 2 Old = 0
- Mem[x] = 2 after the sequence

© David Kirk/NVIDIA and Wen-mei W. Hwu
 ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

7

7

Timing Scenario #3

Time	Thread 1	Thread 2
1	(0) Old \leftarrow Mem[x]	
2	(1) New \leftarrow Old + 1	
3		(0) Old \leftarrow Mem[x]
4	(1) Mem[x] \leftarrow New	
5		(1) New \leftarrow Old + 1
6		(1) Mem[x] \leftarrow New

- Thread 1 Old = 0
- Thread 2 Old = 0
- Mem[x] = 1 after the sequence

© David Kirk/NVIDIA and Wen-mei W. Hwu
 ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

8

8

Timing Scenario #4

Time	Thread 1	Thread 2
1		(0) Old \leftarrow Mem[x]
2		(1) New \leftarrow Old + 1
3	(0) Old \leftarrow Mem[x]	
4		(1) Mem[x] \leftarrow New
5	(1) New \leftarrow Old + 1	
6	(1) Mem[x] \leftarrow New	

- Thread 1 Old = 0
- Thread 2 Old = 0
- Mem[x] = 1 after the sequence

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

9

9

Atomic Operations Prevent Interleaving

```
thread1: Old  $\leftarrow$  Mem[x]
        New  $\leftarrow$  Old + 1
        Mem[x]  $\leftarrow$  New
```

```
thread2: Old  $\leftarrow$  Mem[x]
        New  $\leftarrow$  Old + 1
        Mem[x]  $\leftarrow$  New
```

Or

```
thread1: Old  $\leftarrow$  Mem[x]
        New  $\leftarrow$  Old + 1
        Mem[x]  $\leftarrow$  New
```

```
thread2: Old  $\leftarrow$  Mem[x]
        New  $\leftarrow$  Old + 1
        Mem[x]  $\leftarrow$  New
```

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

10

10

Without Atomic Operations

Mem[x] initialized to 0

```
thread1: Old  $\leftarrow$  Mem[x]
        New  $\leftarrow$  Old + 1
        Mem[x]  $\leftarrow$  New

thread2: Old  $\leftarrow$  Mem[x]
        New  $\leftarrow$  Old + 1
        Mem[x]  $\leftarrow$  New
```

- Both threads receive 0
- Mem[x] becomes 1

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

11

11

Needed When Threads Write to Same Location

When **two threads**

- **may write to the same memory** location,
- the program may **need atomic operations**.

Sharing is **not always easy to recognize...**

- Do two insertions into a hash table share data?
- What about two graph node updates based on all of the nodes' neighbors?
- What if nodes are on same side of bipartite graph?

© Steven S. Lumetta
ECE408/CS483/ECE498aI, University of Illinois, 2009-2020

12

12

What Exactly is “Atomic?”

To a high-energy photon, atoms are not.

Atomicity is ALWAYS with respect to something.

Two sections of code

- that execute atomically with respect to one another
- appear to the software as though
- the programs' execution did not interleave at all.

© Steven S. Lumetta
ECE408/CS483/ECE498aI, University of Illinois, 2009-2020

13

13

What Can Go Wrong?

Common failure mode:

- Programmer *thinks* operations are independent.
- Hasn't considered input data for which they are not.
- Or another programmer reuses code without understanding assumptions that imply independence.

Also: atomicity does not constrain relative order.

© Steven S. Lumetta
ECE408/CS483/ECE498aI, University of Illinois, 2009-2020

14

14

Implementing Atomic Operations

- Many ISAs offer synchronization primitives,
 - instructions with one (or more) address operands
 - that execute atomically with respect to one another when used on the same address.
- Mostly read, modify, write operations
 - Bit test and set
 - Compare and swap / exchange
 - Swap / exchange
 - Fetch and increment / add

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

15

15

Atomicity Enforced by Microarchitecture

When synchronization primitives execute,

- hardware ensures that no other thread
- accesses the location until the operation is complete.

Other threads that access the location

- are typically stalled or held in a queue until their turn.
- Threads perform atomic operations serially.

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

16

16

Atomic Operations in CUDA

- Function calls that are translated into single ISA instructions (a.k.a. *intrinsic*s)
 - Atomic add, sub, inc, dec, min, max, exch (exchange), CAS (compare and swap)
 - Read CUDA C programming Guide for more details

- Atomic Add

int atomicAdd(int address, int val);*

reads the 32-bit word **old** pointed to by **address** in global or shared memory, computes (**old** + **val**), and stores the result back to memory at the same address. The function returns **old**.

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

17

17

More Atomic Adds in CUDA

- Unsigned 32-bit integer atomic add
unsigned int atomicAdd(unsigned int address, unsigned int val);*
- Unsigned 64-bit integer atomic add
unsigned long long int atomicAdd(unsigned long long int address, unsigned long long int val);*
- Single-precision floating-point atomic add (capability > 2.0)
 - *float atomicAdd(float* address, float val);*

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

18

18

Histogramming

- A method for extracting notable features and patterns from large data sets
 - Feature extraction for object recognition in images
 - Fraud detection in credit card transactions
 - Correlating heavenly object movements in astrophysics
 - ...
- Basic histograms - for each element in the data set, use the value to identify a “bin” to increment

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

19

19

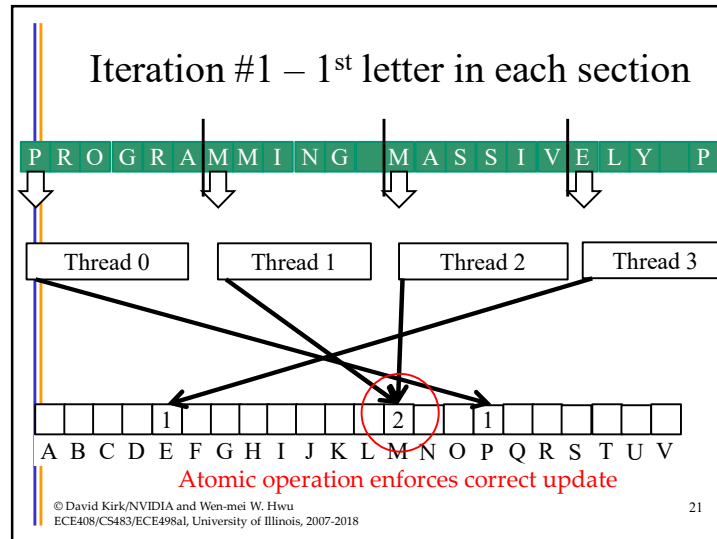
A Histogram Example

- In sentence “Programming Massively Parallel Processors” build a histogram of frequencies of each letter
- A(4), C(1), E(1), G(1), ...
- How do you do this in parallel?

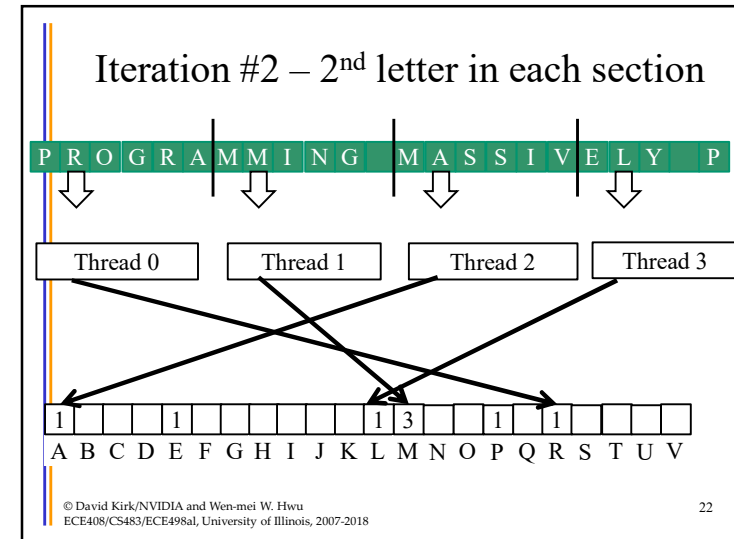
© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

20

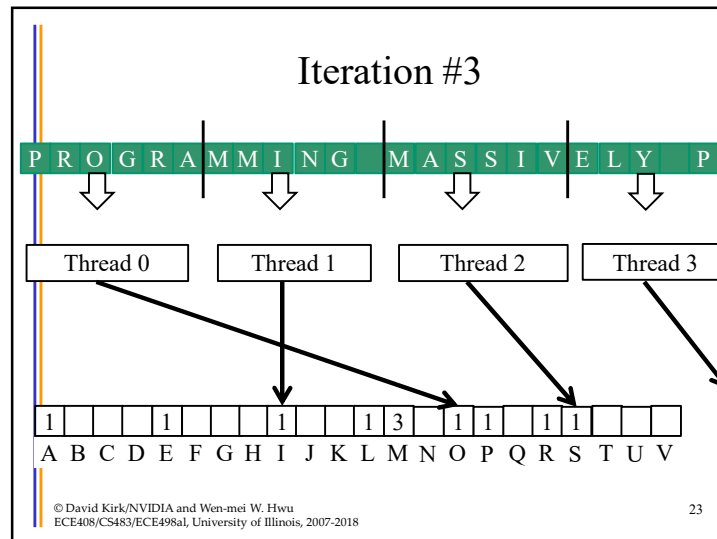
20



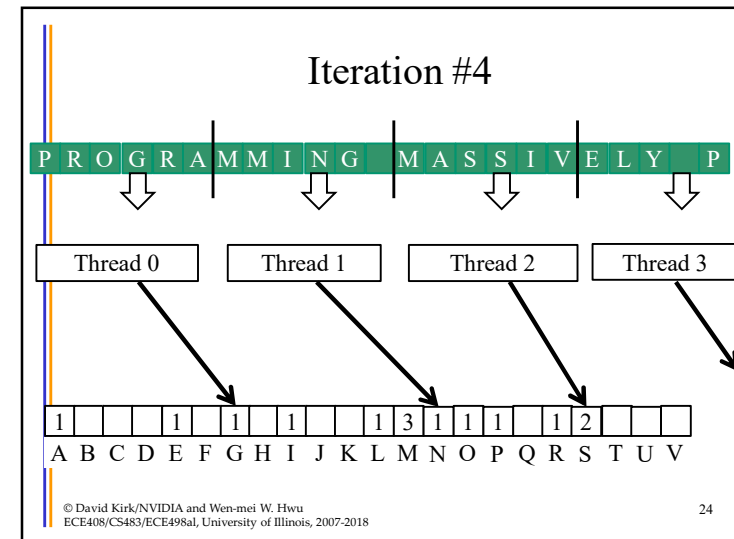
21



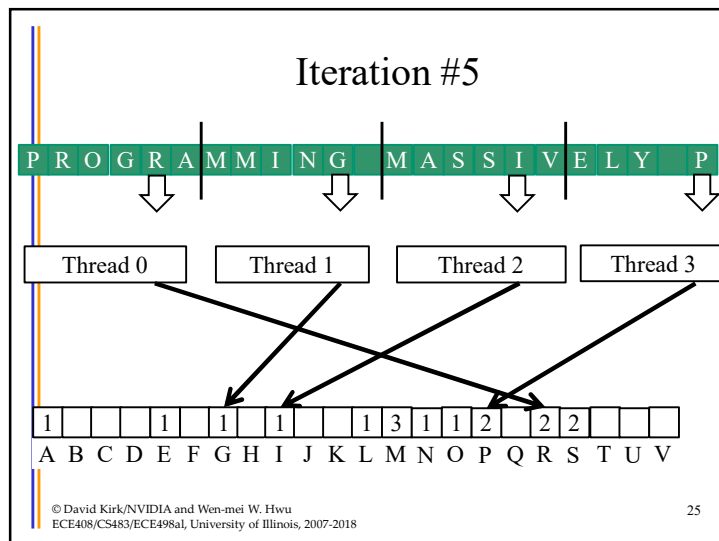
22



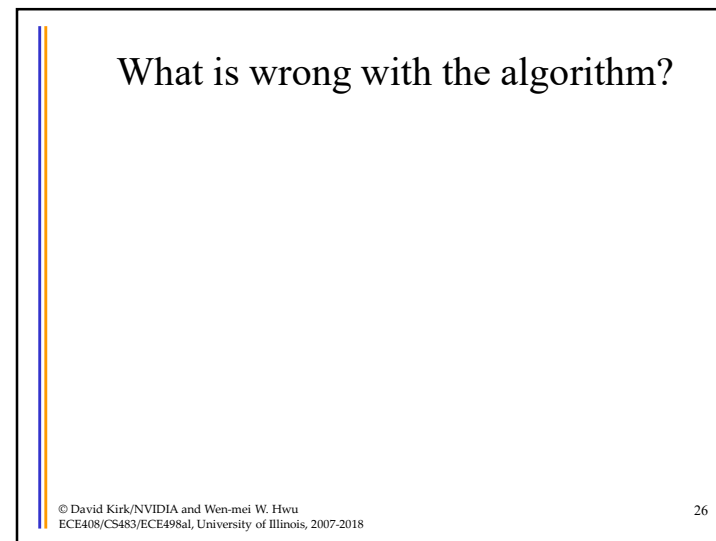
23



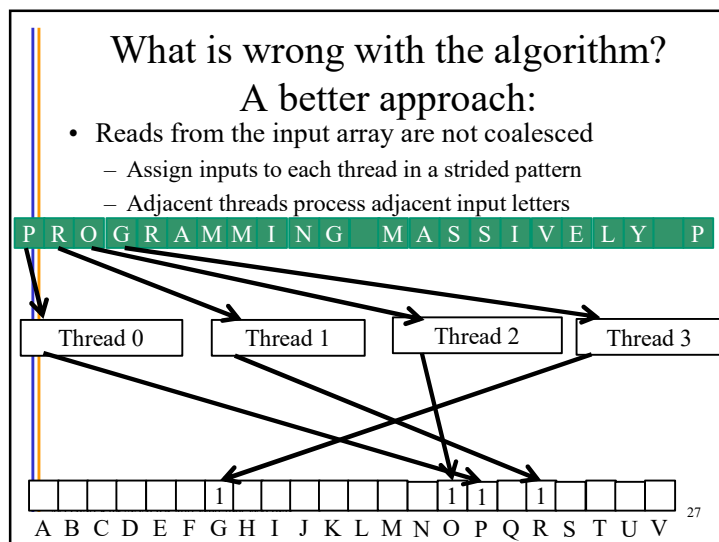
24



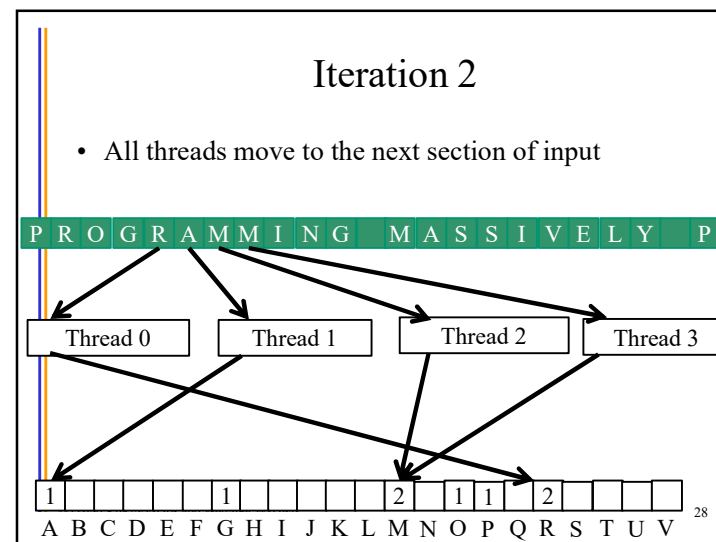
25



26



27



28

A Histogram Kernel

- The kernel receives a pointer to the input buffer
- Each thread process the input in a strided pattern

```
__global__ void histo_kernel(unsigned char *buffer,
                             long size, unsigned int *histo)
{
    int i = threadIdx.x + blockIdx.x * blockDim.x;

    // stride is total number of threads
    int stride = blockDim.x * gridDim.x;
```

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

29

29

More on the Histogram Kernel

```
// All threads in the grid collectively handle
// blockDim.x * gridDim.x consecutive elements
```

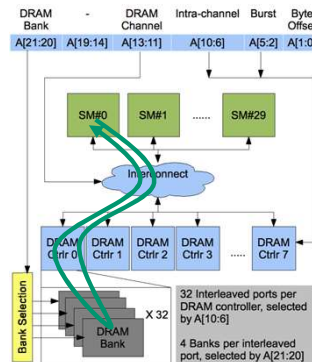
```
while (i < size) {
    atomicAdd( &(histo[buffer[i]]), 1);
    i += stride;
}
```

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

30

30

Atomic Operations on DRAM



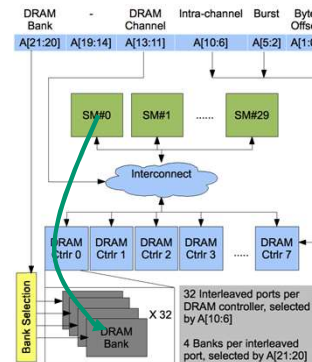
- An atomic operation starts with a read, with a latency of a few hundred cycles

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

31

31

Atomic Operations on DRAM



- An atomic operation starts with a read, with a latency of a few hundred cycles
- The atomic operation ends with a write, with a latency of a few hundred cycles
- During this whole time, no one else can access the location

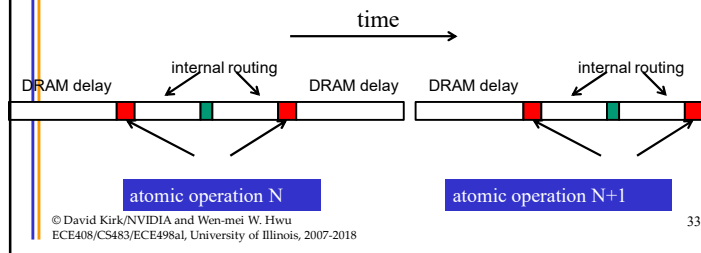
© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498a, University of Illinois, 2007-2018

32

32

Atomic Operations on DRAM

- Each Load-Modify-Store has two full memory access delays
 - All atomic operations on the same variable (RAM location) are serialized



33

Latency determines throughput of atomic operations

- Throughput of an atomic operation is the rate at which the application can execute an atomic operation on a particular location.
- The rate is limited by the total latency of the read-modify-write sequence, typically more than 1000 cycles for global memory (DRAM) locations.
- This means that if many threads attempt to do atomic operation on the same location (contention), the memory bandwidth is reduced to $< 1/1000!$

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

34

34

You may have a similar experience in supermarket checkout

- Some customers realize that they missed an item after they started to check out
- They run to the isle and get the item while the line waits
 - The rate of check is reduced due to the long latency of running to the isle and back.
- Imagine a store where every customer starts the check out before they even fetch any of the items
 - The rate of the checkout will be $1 / (\text{entire shopping time of each customer})$

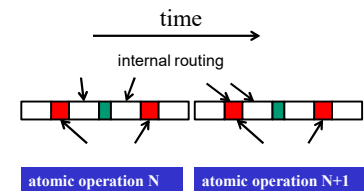
© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498a, University of Illinois, 2007-2018

35

35

Hardware Improvements

- Atomic operations on L2 cache
 - medium latency, but still serialized
 - Global to all blocks
 - “Free improvement” on Global Memory atomics



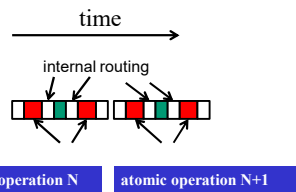
© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498a, University of Illinois, 2007-2018

36

36

Hardware Improvements

- Atomic operations on Shared Memory
 - Very short latency, but still serialized
 - Private to each thread block
 - Need algorithm work by programmers (more later)



© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

37

37

Atomics in Shared Memory Requires Privatization

- Create private copies of the histo[] array for each thread block

```
__global__ void histo_kernel(unsigned char *buffer,
                             long size, unsigned int *histo)
{
    __shared__ unsigned int histo_private[256];
    if (threadIdx.x < 256) histo_private[threadIdx.x] = 0;
    __syncthreads();
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2017

38

38

Build Private Histogram

```
int i = threadIdx.x + blockIdx.x * blockDim.x;
// stride is total number of threads
int stride = blockDim.x * gridDim.x;
while (i < size) {
    atomicAdd( &(private_histo[buffer[i]], 1);
    i += stride;
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2017

39

39

Build Final Histogram

```
// wait for all other threads in the block to finish
__syncthreads();

if (threadIdx.x < 256)
    atomicAdd( &(histo[threadIdx.x]),
               private_histo[threadIdx.x] );
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2017

40

40

More on Privatization

- Privatization is a powerful and frequently used techniques for parallelizing applications
- The operation needs to be associative and commutative
 - Histogram add operation is associative and commutative
- The histogram size needs to be small
 - Fits into shared memory
- What if the histogram is too large to privatize?

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2017

41

41

**ANY MORE QUESTIONS
READ CHAPTER 9**

© David Kirk/NVIDIA and Wen-mei W. Hwu
ECE408/CS483/ECE498aI, University of Illinois, 2007-2018

42

42