



GPU Teaching Kit

Accelerated Computing



Module 7.4 – Parallel Computation Patterns (Histogram)

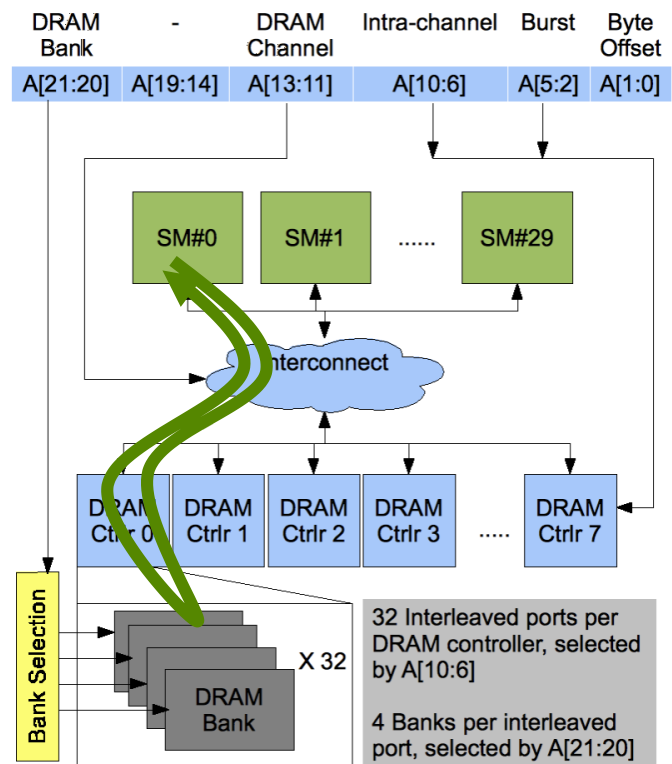
Atomic Operation Performance

Objective

- To learn about the main performance considerations of atomic operations
 - Latency and throughput of atomic operations
 - Atomic operations on global memory
 - Atomic operations on shared L2 cache
 - Atomic operations on shared memory

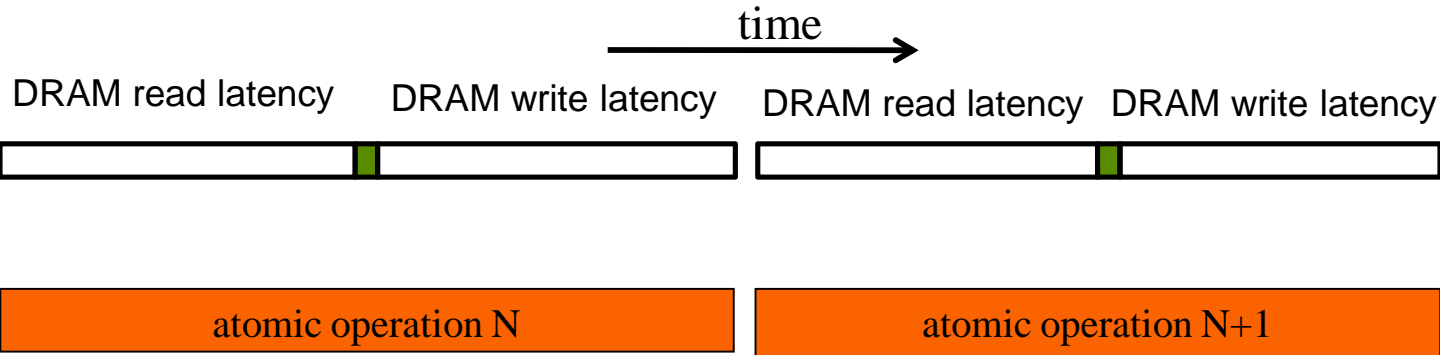
Atomic Operations on Global Memory (DRAM)

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



Atomic Operations on DRAM

- Each Read-Modify-Write has two full memory access delays
 - All atomic operations on the same variable (DRAM location) are serialized



Latency determines throughput

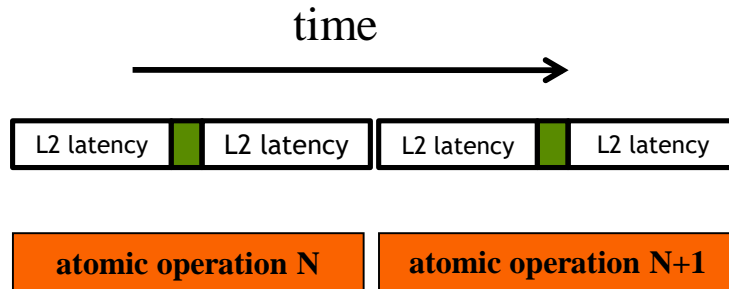
- Throughput of atomic operations on the same DRAM location is the rate at which the application can execute an atomic operation.
- The rate for atomic operation on a particular location is limited by the total latency of the read-modify-write sequence, typically more than 1000 cycles for global memory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

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 aisle and get the item while the line waits
 - The rate of checkout is drastically reduced due to the long latency of running to the aisle 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 $1 / (\text{entire shopping time of each customer})$

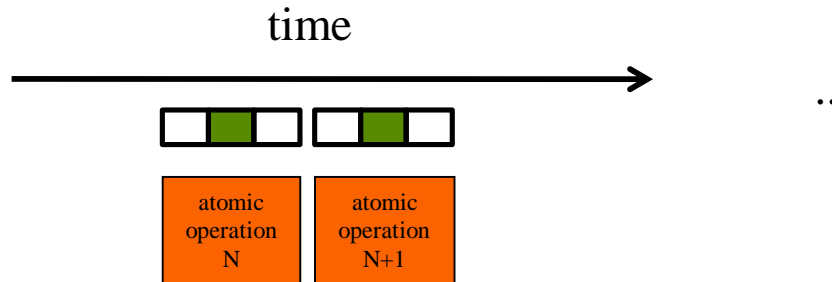
Hardware Improvements

- Atomic operations on Fermi L2 cache
 - Medium latency, about 1/10 of the DRAM latency
 - Shared among all blocks
 - “Free improvement” on Global Memory atomics



Hardware Improvements

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





GPU Teaching Kit

Accelerated Computing



The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under the [Creative Commons Attribution-NonCommercial 4.0 International License](https://creativecommons.org/licenses/by-nc/4.0/).