

# CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Iowa State University

October 21, 2013

# Outline

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Outline

## Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Race conditions

- ▶ Let `int *x` point to global memory. `*x++` happens in 3 steps:
  1. Read the value in `*x` into a register.
  2. Add 1 to the value read in step 1.
  3. Write the result back to `*x`.
- ▶ If we want parallel threads A and B to both increment `*x`, then we want something like:
  1. Thread A reads the value, 7, from `*x`.
  2. Thread A adds 1 to its value, 7, to make 8.
  3. Thread A writes its value, 8, back to `*x`.
  4. Thread B reads the value, 8, from `*x`.
  5. Thread B adds 1 to its value, 8, to make 9.
  6. Thread B writes the value, 9, back to `*x`.

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Race conditions

- ▶ But since the threads are parallel, we might actually get:
  1. Thread A reads the value, 7, from `*x`.
  2. Thread B reads the value, 7, from `*x`.
  3. Thread A adds 1 to its value, 7, to make 8.
  4. Thread A writes its value, 8, back to `*x`.
  5. Thread B adds 1 to its value, 7, to make 8.
  6. Thread B writes the value, 8, back to `*x`.

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Example: race\_condition.cu

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <cuda.h>
4 #include <cuda_runtime.h>
5
6 __global__ void colonel(int *a_d){
7     *a_d += 1;
8 }
9
10 int main(){
11
12     int a = 0, *a_d;
13
14     cudaMalloc((void**) &a_d, sizeof(int));
15     cudaMemcpy(a_d, &a, sizeof(int), cudaMemcpyHostToDevice);
16
17     float   elapsedTime;
18     cudaEvent_t start, stop;
19     cudaEventCreate(&start);
20     cudaEventCreate(&stop);
21     cudaEventRecord( start, 0 );
22
23     colonel <<<1000,1000>>>(a_d);

```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## Example: race\_condition.cu

```

24  cudaEventRecord( stop, 0 );
25  cudaEventSynchronize( stop );
26  cudaEventElapsedTime( &elapsedTime, start, stop );
27  cudaEventDestroy( start );
28  cudaEventDestroy( stop );
29  printf("GPU Time elapsed: %f seconds\n", elapsedTime/1000.0);
30
31
32  cudaMemcpy(&a, a_d, sizeof(int), cudaMemcpyDeviceToHost);
33
34  printf(" a = %d\n", a);
35  cudaFree(a_d);
36
37  }

```

```

1 > nvcc race_condition.cu -o race_condition
2 > ./race_condition
3 GPU Time elapsed: 0.000148 seconds
4 a = 88

```

- ▶ Since we started with  $a$  at 0, we should have gotten  $a = 1000 \cdot 1000 = 1,000,000$ .

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Race conditions

- ▶ **Race condition:** A computational hazard that arises when the results of the program depend on the timing of uncontrollable events, such as the execution order or threads.
  - ▶ Many race conditions are caused by violations of the SIMD paradigm.
  - ▶ Atomic operations and locks are brute force ways to fix race conditions.

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps



# Outline

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Atomics

- ▶ **Atomic operation:** an operation that forces otherwise parallel threads into a bottleneck, executing the operation one at a time.
- ▶ In `colonel()`, replace

```
*a_d += 1;
```

with an atomic function,

```
atomicAdd(a_d, 1);
```

to fix the race condition in `race_condition.cu`.

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## race\_condition\_fixed.cu

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <cuda.h>
4 #include <cuda_runtime.h>
5
6 __global__ void colonel(int *a_d){
7     atomicAdd(a_d, 1);
8 }
9
10 int main(){
11
12     int a = 0, *a_d;
13
14     cudaMalloc((void**) &a_d, sizeof(int));
15     cudaMemcpy(a_d, &a, sizeof(int), cudaMemcpyHostToDevice);
16
17     float   elapsedTime;
18     cudaEvent_t start, stop;
19     cudaEventCreate(&start);
20     cudaEventCreate(&stop);
21     cudaEventRecord( start, 0 );
22
23     colonel <<<1000,1000>>>(a_d);

```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## race\_condition\_fixed.cu

```
24  cudaEventRecord( stop , 0 );
25  cudaEventSynchronize( stop );
26  cudaEventElapsedTime( &elapsedTime , start , stop );
27  cudaEventDestroy( start );
28  cudaEventDestroy( stop );
29  printf("GPU Time elapsed: %f seconds\n" , elapsedTime/1000.0);
30
31
32  cudaMemcpy(&a , a_d , sizeof(int) , cudaMemcpyDeviceToHost);
33
34  printf(" a = %d\n" , a);
35  cudaFree(a_d);
36
37 }
```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## race\_condition\_fixed.cu

```
1 > nvcc race_condition_fixed.cu -arch sm_20 -o  
  race_condition_fixed  
2 > ./race_doncition_fixed  
3 GPU Time elapsed: 0.01485 seconds  
4 a = 1000000
```

- ▶ We got the right answer this time, and execution was slower because we forced the threads to execute the addition sequentially.
- ▶ If you're using builtin atomic functions like `atomicAdd()`, use the `-arch sm_20` flag in compilation.
  - ▶ This is to make sure you're using CUDA compute capability (version) 2.0 or above.

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# CUDA C builtin atomic functions

- ▶ With CUDA compute capability 2.0 or above, you can use:
  - ▶ `atomicAdd()`
  - ▶ `atomicSub()`
  - ▶ `atomicMin()`
  - ▶ `atomicMax()`
  - ▶ `atomicInc()`
  - ▶ `atomicDec()`
  - ▶ `atomicAdd()`
  - ▶ `atomicExch()`
  - ▶ `atomicCAS()`
  - ▶ `atomicAnd()`
  - ▶ `atomicOr()`
  - ▶ `atomicXor()`
- ▶ For documentation, refer to the [CUDA C programming guide](#).

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

`atomicCAS(int *address, int compare, int val)`: needed for locks

1. Read the value, `old`, located at `address`.
2. `*address = (old == compare) ? val : old;`
3. Return `old`.

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## Locks and mutex

- ▶ **Lock:** a mechanism in parallel computing that forces an entire segment of code to be executed atomically.
- ▶ **mutex**
  - ▶ “mutual exclusion”, the principle behind locks.
  - ▶ While a thread is running code inside a lock, it shuts all the other threads out of the lock.

```
1  __global__ void someKernel(void){
2  Lock myLock;
3
4  // some parallel code
5
6  mylock.lock();
7  // some sequential code
8  mylock.unlock();
9
10 // some parallel code
11 }
```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

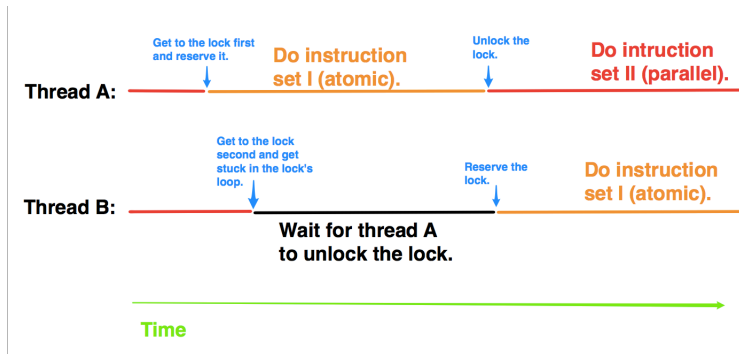
Race conditions

Brute force fixes: atomics, locks, and mutex

Warps



# The concept



CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Lock.h

```

1  struct Lock {
2      int *mutex;
3
4      Lock(){
5          int state = 0;
6
7          cudaMalloc((void**) &mutex, sizeof(int));
8          cudaMemcpy(mutex, &state, sizeof(int), cudaMemcpyHostToDevice);
9      }
10
11     ~Lock(){
12         cudaFree(mutex);
13     }
14
15     __device__ void lock(){
16         while(atomicCAS(mutex, 0, 1) != 0);
17     }
18
19     __device__ void unlock(){
20         atomicExch(mutex, 0);
21     }
22
23 };

```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## A closer look at the lock function

```

15  __device__ void lock(){
16     while(atomicCAS(mutex, 0, 1) != 0);
17 }

```

### ► In pseudocode:

```

1  __device void lock(){
2     repeat{
3         do atomically{
4
5             if(mutex == 0){
6                 mutex = 1;
7                 return_value = 0;
8             }
9
10            else if(mutex == 1){
11                return_value = 1;
12            }
13        } // do atomically
14
15        if(return_value == 0)
16            exit loop;
17
18    } // repeat
19 } // lock

```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## Example: counting the number of blocks

- ▶ Compare the two kernels:

```

5  __global__ void blockCounterUnlocked(int *
      nblocks){
6      if(threadIdx.x == 0){
7          *nblocks = *nblocks + 1;
8      }
9  }

```

```

11 __global__ void blockCounter1(Lock lock , int *
      nblocks){
12     if(threadIdx.x == 0){
13         lock.lock();
14         *nblocks = *nblocks + 1;
15         lock.unlock();
16     }
17 }

```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## blockCounter.cu

```

1 #include "../common/lock.h"
2 #define NBLOCKS_TRUE 512
3 #define NTHREADS_TRUE 512 * 2
4
5 __global__ void blockCounterUnlocked( int *nblocks ){
6     if( threadIdx.x == 0){
7         *nblocks = *nblocks + 1;
8     }
9 }
10
11 __global__ void blockCounter1( Lock lock, int *nblocks ){
12     if( threadIdx.x == 0){
13         lock.lock();
14         *nblocks = *nblocks + 1;
15         lock.unlock();
16     }
17 }
18
19 int main(){
20     int nblocks_host, *nblocks_dev;
21     Lock lock;
22     float elapsedTime;
23     cudaEvent_t start, stop;
24
25     cudaMalloc((void**) &nblocks_dev, sizeof(int));
26
27     //blockCounterUnlocked:
28
29     nblocks_host = 0;
30     cudaMemcpy( nblocks_dev, &nblocks_host, sizeof(int),
                cudaMemcpyHostToDevice );

```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## blockCounter.cu

```

32  cudaEventCreate(&start);
33  cudaEventCreate(&stop);
34  cudaEventRecord( start , 0 );
35
36  blockCounterUnlocked<<<NBLOCKS.TRUE, NTHREADS.TRUE>>>(nblocks_dev);
37
38  cudaEventRecord( stop , 0 );
39  cudaEventSynchronize( stop );
40  cudaEventElapsedTime( &elapsedTime , start , stop );
41
42  cudaEventDestroy( start );
43  cudaEventDestroy( stop );
44
45  cudaMemcpy( &nblocks_host , nblocks_dev , sizeof(int) ,
              cudaMemcpyDeviceToHost );
46  printf(" blockCounterUnlocked <<< %d, %d >>> () counted %d blocks in
         %f ms.\n" ,
         NBLOCKS.TRUE,
         NTHREADS.TRUE,
         nblocks_host ,
         elapsedTime);
50

```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## blockCounter.cu

```

51 //blockCounter1:
52
53 nblocks_host = 0;
54 cudaMemcpy( nblocks_dev , &nblocks_host , sizeof(int) ,
             cudaMemcpyHostToDevice );
55
56 cudaEventCreate(&start);
57 cudaEventCreate(&stop);
58 cudaEventRecord( start , 0 );
59
60 blockCounter1<<<NBLOCKS_TRUE, NTHREADS_TRUE>>>(lock , nblocks_dev);
61
62 cudaEventRecord( stop , 0 );
63 cudaEventSynchronize( stop );
64 cudaEventElapsedTime( &elapsedTime , start , stop );
65
66 cudaEventDestroy( start );
67 cudaEventDestroy( stop );
68
69 cudaMemcpy( &nblocks_host , nblocks_dev , sizeof(int) ,
             cudaMemcpyDeviceToHost );
70 printf(" blockCounter1 <<< %d, %d >>> () counted %d blocks in %f ms
       .\n" ,
       NBLOCKS_TRUE,
       NTHREADS_TRUE,
       nblocks_host ,
       elapsedTime);
75
76 cudaFree( nblocks_dev );
77 }

```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# blockCounter.cu

```
78 > nvcc blockCounter.cu -arch sm_20 -o blockCounter
79 > ./blockCounter
80 blockCounterUnlocked <<< 512, 1024 >>> () counted 47 blocks in
    0.057920 ms.
81 blockCounter1 <<< 512, 1024 >>> () counted 512 blocks in 0.636064 ms.
```

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps



# blockCounter.cu pauses indefinitely with this kernel

```
1  __global__ void blockCounter2(Lock lock , int *
    nblocks){
2  lock.lock();
3  if(threadIdx.x == 0){
4      *nblocks = *nblocks + 1;
5  }
6  lock.unlock();
7  }
```

- ▶ Why? warps.

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Outline

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

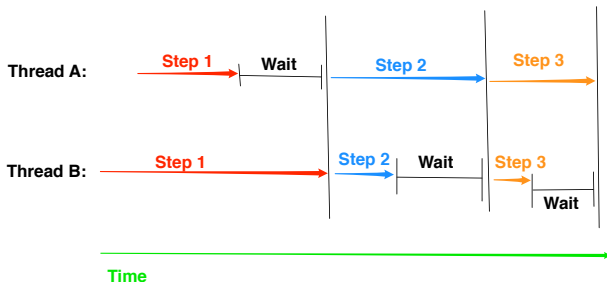
Brute force fixes: atomics, locks, and mutex

Warps

# Warps

- ▶ Warp: a group of 32 threads in the same block that execute in lockstep.
  - ▶ That is, they synchronize after every step (as if `__syncthreads()` is called as often as possible).
  - ▶ All blocks are partitioned into warps.

## Threads in the same warp:



CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Warps

CUDA C: race conditions, atomics, locks, mutex, and warps

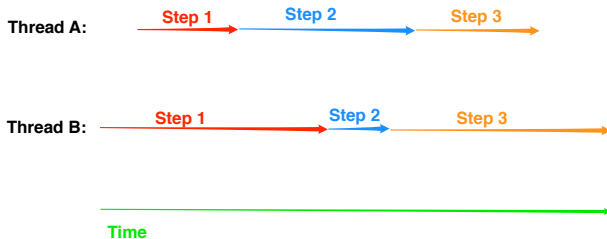
Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## Threads in different warps:



# Warps and locks

CUDA C: race conditions, atomics, locks, mutex, and warps

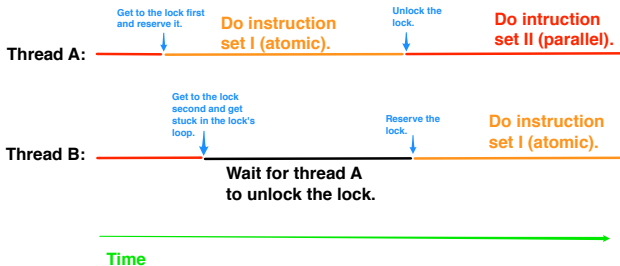
Will Landau

Race conditions

Brute force fixes:  
atomics, locks, and mutex

Warps

## Threads in different warps:



# Warps and locks

CUDA C: race conditions, atomics, locks, mutex, and warps

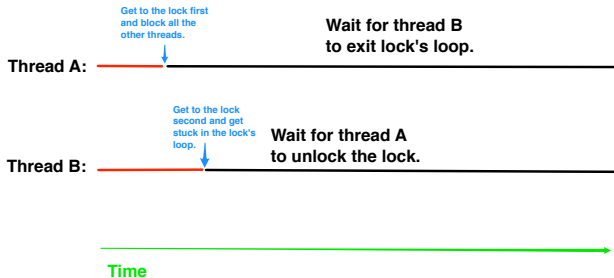
Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

## Threads in the same warp:



# Outline

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps

# Resources

▶ Texts:

1. J. Sanders and E. Kandrot. *CUDA by Example*. Addison-Wesley, 2010.
2. D. Kirk, W.H. Wen-mei, and W. Hwu. *Programming massively parallel processors: a hands-on approach*. Morgan Kaufmann, 2010.

▶ Code from today:

- ▶ [race\\_condition.cu](#)
- ▶ [race\\_condition\\_fixed.cu](#)
- ▶ [blockCounter.cu](#)

▶ Dot product with atomic operations:

- ▶ [dot\\_product\\_atomic\\_builtin.cu](#)
- ▶ [dot\\_product\\_atomic\\_lock.cu](#)

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps



That's all for today.

- ▶ Series materials are available at <http://will-landau.com/gpu>.

CUDA C: race conditions, atomics, locks, mutex, and warps

Will Landau

Race conditions

Brute force fixes: atomics, locks, and mutex

Warps