

Shared memory.

Within a thread block - sharing

Programmable L1 cache / scratchpad memory

Useful for repeated small data / synchronization.

-- shared -- float a[N]; // inside kernel

-- shared -- unsigned s;

a[id] = id;

if (id == 0) s = 1;

what is this id ? - threadIdx only

no block ID ... as it is within a thread block.

5/2/2020.

Classwork.

You are given a 1024×1024 integer matrix M . Each row is assigned to a thread block. Each thread is assigned a matrix element $M[i][j]$.

It changes $M[i][j]$ to $M[i][j+1]$ (where possible). Exploit shared memory to do this in $O(1)$.

- void fun(~~int i, int j~~);

shared ~~matrix[N][N]~~, it creates a lot of contention for shared memory between threads in the same block.

Change ~~int i, int j~~ to blockIdx.x:

int colNumber = threadIdx.x;

if (colNumber < N) {

 a[~~rowNumber * N + i~~] =

 matrix[~~rowNumber * N + i + 1~~][colNumber];

+ ~~matrix[~~rowNumber * N + i + 1~~][colNumber]~~

→ It would be better if we have 1024 per thread block in shared memory

in shared memory

~~int rowNumber = blockIdx.y * 1024 + threadIdx.y;~~

~~int colNumber = threadIdx.x;~~

~~if (colNumber < 1024) {~~

~~a[~~rowNumber * 1024 + i~~] =~~

~~matrix[~~rowNumber * 1024 + i + 1~~][colNumber];~~

~~+ matrix[~~rowNumber * 1024 + i + 1~~][colNumber]~~

~~}~~

int * M2

91 we want 4 KB per. row

4 KB per. row

global void thread (int * M, int N){

shared -- int mrow [1024];

A copy from M to mrow using blockedIdx.x

~~for (int i = 0; i < N; i++) {~~ threadIdx.x * blockDim.x + i] = M [blockIdx.x * blockDim.x + i];

3

if (blockIdx.x < 1022) {

if (threadIdx.x * blockDim.x + threadIdx.x < N) {

cout << mrow [~~threadIdx.x~~ * 1024] +

mrow [threadIdx.x + 1];

3

3

* Need to ensure that all the copying from M to mrow

is done before the compute step starts and all the

* for all the threads within a warp, no interleaving

0 - 30 ✓ 31 ← can be a problem

next warp 0th thread is done / not??

Precise dependency: last thread of 1 warp waits for first
thread of next warp.

Overapprox: All writes to mrow finish before the
compute part starts.

⇒ Barrier, after copy step, before compute step.

--syncthreads(); // ensures all the threads within a
thread block reach that point before
proceeding any further.

Within a warp, implicit barrier.

```
if (threadIdx.x == 0) s = 0; } within a warp,  
if (threadIdx.x == 1) s+=1; } one after another  
if (threadIdx.x == 100) s+=2; → another warp.  
if (threadIdx.x == 0) printf("s=%d\n", s);
```

6/2/2020.

CW. What is the o/p?

```
#include <stdio.h>  
#include <cuda.h>  
#define BLOCKSIZE 26  
__global__ void dkernel()  
{  
    __shared__ char str[BLOCKSIZE+1];  
    str[threadIdx.x] = 'A' + (threadIdx.x + blockIdx.x) / .BLOCKSIZE;  
    if (threadIdx.x == 0) {  
        str[BLOCKSIZE] = '\0';  
    }  
    if (threadIdx.x == 0) {  
        printf("%d: %s\n", blockIdx.x, str);  
    }  
}
```

3 $\xrightarrow{\text{32 so, within a warp, so it is executed in SIMD fashion}}$
main() {
 dkernel <<< 10, BLOCKSIZE >>>();
 cudaDeviceSynchronize();
}

blockIdx.x = 0 $\xrightarrow{\text{idx}} \text{str}[0] = 'A'$
 $\xrightarrow{\text{idx}} \text{str}[1] = 'A' + 1 = 'B'$
 $\xrightarrow{\text{idx}} \text{str}[25] = 'A' + 25 = 'Z'$
 $\xrightarrow{\text{idx}} \text{str}[26] = '\0'$

6/2/2020. 88

A B C D ... Z

Block size = 32
 if i < 32 str(B) = 'A' + ① => B
 So 32
 C
 D
 E
 F
 G
 H
 I
 J
 K
 L
 M
 N
 O
 P
 Q
 R
 S
 T
 U
 V
 W
 X
 Y
 Z

Block size = 32
 So 32

: 9 - JKLMHI

Order of 32 → can be any order

Q) what is the potential bug in this code?

→ If $\text{BLOCKSIZE} > 32$, then the ^{multiple} warps may run in asynchronous manner --

Put `syncthreads();` b/w the 2 if's ?

both are okay.

(2) before first if ?

: 2 if's are going to be within a thread, so req.

Where can data race happen?

CW: convert the following c codes to optimized CUDA

~~for (i=0; i<N; i++) {~~
~~p += a[i];~~
~~q = g * a[i];~~
~~r = q / a[i];~~
~~}~~

~~-- global -- void kernel (int k, int N,~~
~~int *p, int *g, int *r) {~~
~~p[0] = 0;~~
~~for (int i = 1; i < N; i++) {~~
~~p[i] = p[i-1] + a[i];~~
~~q = g[i] * a[i];~~
~~r[i] = q / a[i];~~
~~}~~

P | 1st option - sequential instruction
 P | 2nd option - atomic instruction
 P | atomicAdd (&p, a[i]) is
 P | can be slower than
 P | for loop on GPU.
 P | 3rd option - divide & conquer / reduction.
 P | likely to be faster
 P | than atomic add

Q) Do we need branch b/w

p, q & r assignment?

→ within a thread, it is ok

Warp within a thread block / across blocks

⇒ atomic instruction & divide & conquer are ok

or same fundas as P.

~~r ← final q / a[i] final]. (if only last one~~
~~is read)~~

④ Local memory access and coalescence

⑤ Local memory access iteration affects
local memory coalescence

```
for (i=0; i<N; ++i) {  
    for (j=0; j<M; ++j) {  
        p = a[i*M + j];  
        q = g * a[i*M + j];  
        r = q / a[i*M + j];  
    }  
}
```

10/21/2020.

LI v/s shared

```
cudaFuncSetCacheConfig( dkernel, cudaMemcpyPrefetchMode);  
" " " " LI  
" " " " Shared
```

Dynamic shared memory

When the amount of shared memory reqd is unknown at compile time,
dynamic shared memory can be used.

```
-- dkernel {  
    dkernel << --, --, (L) >> c);
```

out of dynamic shared memory we require

Similar to main (int argc, char* argv[], char* env[])

↑
Environ vars.

most of the time, we would not require dynamic shared memory.
Only when the amount of work done by a thread depends on the
input, we would need dynamic shared memory.

Q) How can we specify multiple shared variables? ... $s[]$
... $s1[]$

Bank conflicts.

- Shared memory is organized into banks.
- Accesses to a bank are sequential.
↑ queue of requests to one bank.

If all the threads ~~access~~ require data from the same bank,
then the requests will be serviced sequentially.
~~4 bytes = 1 word~~ ⇒ characters will be stored 4 at a time; bank conflict
Consecutive words are stored in adjacent banks → Useful for
Coalesced access.

integers/floats ⇒ one in each bank.
we cannot change the way banks store data. But what we
can change is our access pattern.

EXCEPTION: Warp accesses to the same word are not
sequentialized.

Warp accesses to different words of same bank
⇒ ^{bank} conflict

```
-- global -- void bankNoConflict () {  
    -- shared -- unsigned s[1024];  
    s[(*threadIdx.x)] = threadIdx.x;  
    // stride of 1  
    // normal access.
```

```
-- global void bankConflict() {  
    -- shared unsigned s[1024],  
    s[32 * ThreadIdx.x] = ThreadIdx.x; // stride of 32
```

§

Suppose there are 32 banks — **DIAGRAM**

Accesses to diff words in same bank → Bank Conflict

* Recall False sharing.

Multicore programming ←

write to an array element. — spatial locality is a liability here
invalidate again & again.

Hw: create programs that create bank conflicts
measure time taken

Texture Memory

- Fast read only memory cache
- Optimized for 2d spatial access.
- texture <float 2, cudaMemcpyMode ElementType> tex
#dim ↑
in main(): cudaBindTextureToArray (tex, cuArray, ...);
in kernel: -tex2D (tex, ...);
 ↑
 cuda malloc'd array

→ Example from CUDA SDK

Cannot write to texture inside ~~a kernel~~ a kernel.
But can do multiple dimensions.

constant memory.

- read only memory
- 64 KB per SM
- -- constant -- assigned meta
- main(): cudaMemcpyToSymbol(meta, &meta, sizef);
- kernel: data[threadIdx.x] = meta[0];

compute capability

Version no. major.minor reg 3.5

- features supported by GPU hardware.
- used by application at runtime.

- arch = sm_35. Compatibility. -

~~source code defined in device code~~

Backwards
Compatibility ↗

- 1 - Tesla
 - 2 - Fermi
 - 3 - Kepler
- and so on.

Cuda version ⇒ Software part

CUDA-ARCH ← device code

≥ 350 for simulated thing. → atoms on double
simulated using float.

CW

Write CUDA code for the foll. functionality; use shared mem
Assume foll. data type, filled with some values.

struct Point { int x, y; } arr[N];

Each thread operates on 4 elements.

Find the average AVG of the 'x' values.

If a thread sees 'y' value above the average, it
replaces all 4 y values to AVG.

Else, it adds y values to a global sum.

Host prints # elements set to AVG.

```
# - global - void k() {  
    extern __shared__ int x[];  
    extern __shared__ int y[];  
    ...  
    <<< x, y, 2000 >>>();  
}
```

11-2-2020.

SYNCHRONIZATION.

- Mutual Exclusion No two processes are accessing the same shared resource at the same time.
 - * Critical section - arbitrary piece of Gde, such that, if all threads follow the same discipline / protocol, for dealing with the critical section, they'll all see a semantically consistent value, a sensitive piece of Gde. Even if one doesn't follow, we'll be in trouble.
- Atomic section - Effect of access is, that section is fully executed, or not executed at all.
- Deadlock:
 - * ~~deadlock won't~~
 - * Parallel / concurrent execution.
 - * cyclic dependency.
 - * Hold and wait.
 - * No preemption.
 - * At least 2 resources (shared) reqd at same time by more than 1 thread. But each thread has some reason & needs / waits for some others. locks are usually blocking

b) How to avoid deadlocks?

- ① Deadlock prevention → Single lock for all resources, b/c - Performance will suffer a bit. ∵ not all resources may be used at a time
- ② Allow preemption. → removes one of the necessary conditions for deadlock
- ③ Enforce an order ^(total order) on the acquisition of resources.
 - breaks the cyclic dependency.
 - * avoids deadlock but starvation.

How to avoid starvation? Time-based priorities assigned to the locks/resources.
↓
e.g. priority ↑ with time, so a process waiting for long time gets more priority.

(Process with lowest priority would starve otherwise--)

Linux → CFS: Completely Fair Scheduling

multi-level feedback queues).

* Another way is to prove that deadlocks can happen, they are very rare. Many practical systems follow this ^{though}.

→ T_1 locks R_1 T_2 locks R_2

∴ it doesn't have R_2 , ∴ it doesn't have R_1 ,

T_1 releases R_1 & T_2 releases R_2 &

goes away goes away.

This keeps happening again & again

⇒ Not deadlock, but ~~deadlock~~ live lock (Pehle aap)

Not blocked altogether, but not making progress either.

Data Race. 4 necessary conditions:

1. multiple threads

Data races are usually not good news.

2. common memory location

But for implementation of locks, data races are useful.

3. At least one write

4. concurrent execution

b) How to get rid of data race?

→ avoid one of the 4 necessary conditions.

1. Execute sequentially

2. Privatization / data replication (fol. by some kind of merging)

3. Separating reads + writes by a barrier.

4. Mutual exclusion.

classwork:

T1

```
1 flag = 1;  
2 while (flag) {  
3   :  
4   S1;
```

T2

```
4 while (!flag) {  
  :  
  g  
5 S2;  
6 flag = 0;
```

1 2 stack

~~-----~~ ~~s1~~

1 4 5 6 2 3 S2 S1

1 4 2 3 5 6 S1 S2

4 1 5 6 2 3 S2 S1

4 1 2 5 3 6 S2 S1

4 5 6 1 2 3 S2 S1

4 stack.

12/2/2020.

Class work

Given roll numbers & marks of 80 students in GPU programming, assign grades.

S = 90, A = 80, B = 70, C = 60, D = 50, E = 40, and U.

use input arrays and output arrays.

→ Obtain: one thread for each student, within one thread block.
-- global -- void kernel (int *marks, char *grades) {
 int i; // our roll no
 has the format like
 CS16B011,
 8 chars bounded
 int uid = (blockDim.x * blockIdx.x) + threadIdx.x;
 if (uid < N) {
 if (marks[uid] >= 90) {
 grades[uid] = 'S';
 }
 else if (marks[uid] >= 80 && marks[uid] < 90) {
 grades[uid] = 'A';
 }
 else {
 grades[uid] = 'U';
 }
 }
}

No real synchronization reqd

Suppose we also want the total counts of no. of students with each grade → histGrades [7]

#S, #A, #B, #C, #D, #E, #U

Now, we need to maintain counts of each grade.
There can be 2 threads trying to access same counter. \rightarrow synchronization
reqd.

```
-- global -- void __kernel (int *marks, int N) {
    -- shared -- char grades [4];
    -- shared -- int histGrades [7];

    int uid = -1;
    if (uid < N) {
        if (marks[uid] > 90) {
            Grade[uid] = 'S';
            histGrades[0]++;
        }
        atomic_inc(&histGrades[0], INTMAX);
    }
}
```

Guarantees that, if 2 threads are trying to access same time, they will go one after another, resulting in correct results

3
 Σ
=

Q) Do we need ~~or~~ ^{or rate} syncThreads() ? \rightarrow No. It doesn't matter in which order the diff threads are computing.
(in this case)

Ex: If we want the max of all elements in histGrades, we'll have to wait for all elements of histGrades to be computed, then only we ~~can~~ do atomicMax.

So, -syncThreads() would be reqd here.

^{blocks of}
Statements are separated by barriers into "phases"

We need to check for data race only within a phase,
not across phases.

Shortest Path in a Graph.

^{of India}

Given an input graph, you want to compute the shortest path from Nagpur to every other city. Assume you are given a C++ graph library and the associated routines. Each thread operates on a node and updates distances of the neighbours (Bellman Ford).
↓

nicely parallel.

contrary to Dijkstra's

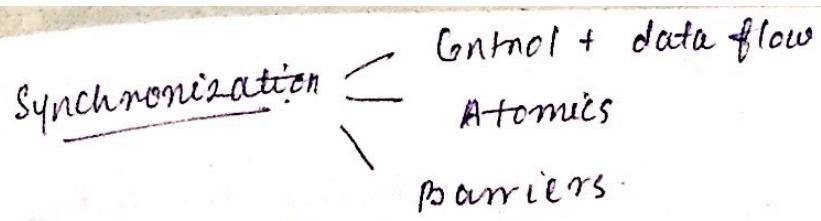
```
-- global -- void dsssp (Graph g, unsigned *dist),  
  (inherently sequential,  
  uses priority queue)  
  unsigned id = ...;  // one thread - one node  
  for each n in g.all_neighbours (id) {  
    unsigned altdist = dist [id] + weight (id, n);  
    if (altdist < dist [n]) {  
      dist [n] = altdist; // incorrect  
      }  
    }
```

3

3

Kernel should be called ^a number of times from CPU
→ max $(n-1)$ times, where n = no. of nodes

- ④ Persistent kernel  if some temp. computation is stored in shared memory
 -  great load imbalance → e.g. in a graph/social network



CW: Implement mutual exclusion for two threads - using only control & data flow methods

| | |
|----------------|----------------|
| T ₁ | T ₂ |
| { | { |
| S ₁ | S ₂ |

S₁, S₂ both are writing to same shared variable.

How to prevent this from happening?

bool flag = false; // Shared

T₁

≡

while(flag);

S₁;

flag = true;

≡

≡

S₂;

flag = false;

≡

Issues

- // Performance drop - spinning on while - unnecessarily wastes CPU cycles
- // Codes for diff threads are different.
- // If we launch less number of threads (in this case), then it gets blocked.

14/2/2020

Debugging and Profiling.

Debugging.

- * Non-determinism due to thread-scheduling - diff outputs possible
- * Depending on the application semantics, intermediate steps may also need to be deterministic.
- * No. of correct intermediate values - may be large.
- * cuda-gdb - for debugging CUDA programs on real hardware.
- * There are GPU simulators also--
- * cuda-gdb allows break points, single step, read/write memory contents.

⑥ cuda-gdb → Design decision → ~~restricts parallelism~~ do a single thread block.

prev errors - probably overflow/etc
may still continue.

in main():

If there are multiple errors,

→ we'll get info about last one only.

```
cudaError_t err = cudaGetLastError();
```

```
printf("error = %d, %s, %s\n", err, cudaGetErrorName(err),  
       cudaGetErrorString(err));
```

e.g. error = 77, cudaErrorIllegalAddress, an illegal memory access was encountered.

cuda-gdb.

```
nvcc -g -G file.cu
```

↑ ↗
adds symbol info into a.out
a.out
(CPU code)

adds symbol info into a.out
(GPU Gde)-
kernels.

// also, disables optimizations.
∴ we would be single-stepping through the source code, so we should be able to see deterministic execution for our source code.

cuda-gdb a.out // loads a.out symbol info and stops

> // prompt

> run // to execute the program

* may have to stop windows manager - GUI part -- graphics application

* due to lots of threads, cuda-gdb works with a focus
(current thread).

↓
probably doesn't perform
SIMD execution -- confirm this --

(cuda-gdb) run

CUDA exception: Device illegal address starts with 1 rather than 0.

[Switching focus to CUDA kernel 0, grid 1, block (1,0,0),
thread (0,0,0), device 0, sm 13, warp 0, lane 0]

↑
thread neither warp

(err: ~~at~~ shows printf(...); ← prev stmt is executed, error occurs
& the cursor is here now.

Kernel << (1,2,3), (4,5,6) >> ();

↓ ↓
3 6

c-semantics

(only 1 in this case)

(cuda-gdb) info cuda kernels ← information about kernels

(cuda-gdb) info threads ← information about each thread.
→ shows only 3 threads, but we launched 20. (2x10)

includes both CPU & GPU threads

2 1 → erroneous one

(cuda-gdb) info cuda threads ← information about GPU
threads only

(cuda-gdb) cuda kernel block thread ← current focus

(cuda-gdb) cuda block 1 thread 0 ← switches focus to another block, thread.

Breakpoints:

or any function name / line #
(cuda-gdb) break main file.cu:223

(cuda-gdb) set cuda break-on-launch application

Conditional breakpoint

break file.cu:29 if threadIdx == 0

Single step: s, step.

Info about SMS: info cuda sms.

Info about warp: info cuda lanes. → shows that warp execution is ~~actually~~ in SIMD fashion for this execution
only 10 threads in a block

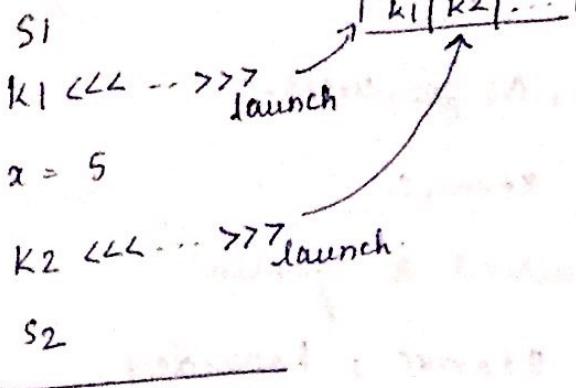
Profiling → next week

2 kernels → by default same stream (queue type)

⇒ no need to have cudaDeviceSynchronize() b/w ² GPU kernels.

but is a good practice to have b/w GPU & CPU instr

11/2/2020.



S_2 can be executing in parallel with K_1 or K_2 .

only the launches happen seq

Synchronization

- control + data flow (initially, flag: false)

```
while (!flag) {  
    S2;  
    flag = true;  
    S1;  
}
```

mutual Exclusion b/w 2 threads.

V1 lock(), unlock()

Methods should be same.

Say, tids $\Rightarrow 0, 1$:

lock()

1 me = tid

2 other = 1 - me;

3 flag [me] = true;

4 while (flag [other]);

unlock()

11 flag [tid] = false;

T1 - L1 L2 L3

T2 - L1 L2 L3

Q1) can the threads wait indefinitely? Yes - both wait-for each other

Q2) Should flags be initialized to false? \rightarrow Yes. (otherwise, if 1st doesn't execute at all, other also can't go into CS)

Q3) Does the code ensure mutual exclusion? \rightarrow Yes.

deadlock (lock \rightarrow threads release, again get lock, again have to release \rightarrow)

both wait-for each other
(\neq 3TAA)
(otherwise, if 1st doesn't execute at all, other also can't go into CS,
Other cannot go into CS.)

- * Mutual exclusion guaranteed
- * May lead to deadlock
- * If one thread runs before the other, all goes well.

④ Primitive assignments must be atomic

Flag [me] = true / false // Inherent assignments

Rest assured that these will be atomic, provided by hardware:

keyword in C → disables optimizations for that variable.
Compiler
primitive → to compiler to the compiler

V2 volatile int victim;

lock()

me = tid;

victim = me;

while (victim == me);

unlock()

;

can remove the loop or CS :: it seems like it'll never happen

Also. if var is in var's private cache, the var stays shared in memory - consistency issues avoided

Cooperative protocol - threads are not selfish

Give importance to "other".

individual threads have to

involve coordination b/w threads ⇒ spend more energy

① If only single thread runs, leads to deadlock

② Does ~~it~~ ensure mutual exclusion? Yes.

Tid 0 → ~~lock~~ gets victim = 0

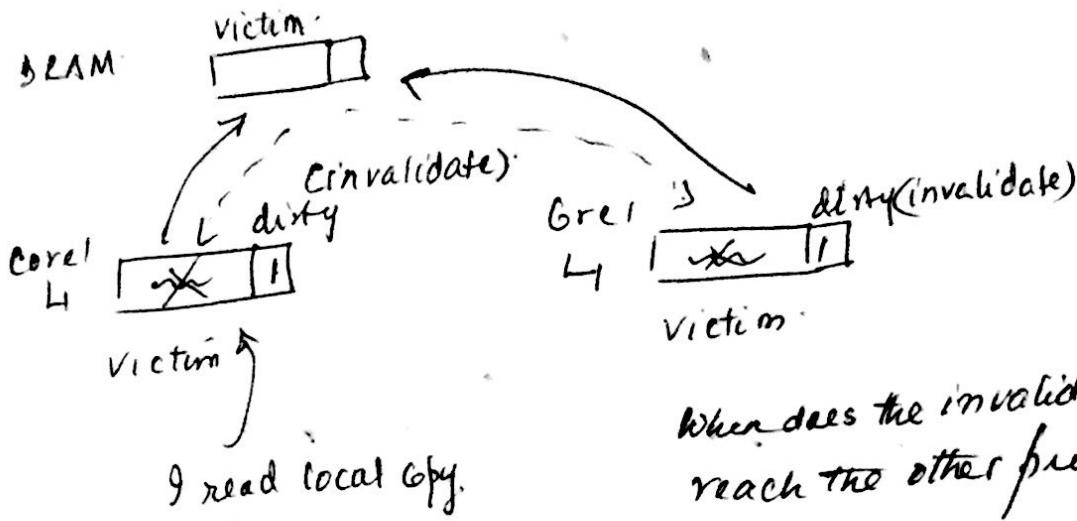
while (victim == 0); ← How this can get in

Now, Tid 1 → ~~lock~~ gets victim = 1) Tid 1 → ~~unlock~~ while loop

Both threads see same value for victim. ← no mutual exclusion

Mutex violated only if Tid 1 sees victim = 0 and Tid 0 sees victim = 1

works well if threads repeatedly take locks



When does the invalidation reach the other private L1?

By marking volatile, it ensures on CPU that:

- a) writing → all the way back to shared mem
- b) reading → all the way from DRAM/shared mem

On GPU:

Reading → all the way from DRAM ✓
 But writing → not all the way from upto shared mem DRAM.

"memory fence" - whatever mem ops are done for a thread are "pushed" to memory

Need to explicitly call memory fence in GPU.

threadfence(); // after every write to shared var

read ← mark var as volatile

* ALSO, | data race is possible |

However, when you implement a synchronization protocol such as lock C, you need to exploit data race.

(v1) → worked when threads were "away" from each other

(v2) → " " " " " close" to each other

(v3) → Combine both.

Volatile bool flag[2]; π must be initialized to false. → crucial to avoid a deadlock if other thread isn't there } to ensure progress, to prevent starvation.

lock()

me = tid;

other = 1 - me;

flag[me] = true;

Victim = me;

while (flag[other] && victim == me);

unlock()

flag[tid] = false;

Peterson's protocol

Q) Does it ensure mutex? ENSURES MUTUAL EXCLUSION.

Q) Can deadlocks occur? No.

Q) Can starvation occur? No.

Q) Is the algorithm "fair"? Yes... check ✓

[Mutex, Progress, Bounded wait]

L) Both threads get opportunity to enter CS.

18/2/2020.

Peterson's lock for N threads.

while(flag[other] && victim == me)



need to change to 'or' of all the flags . . .

Bakery algorithm.

- devised by Lamport
- works with N threads
- maintain FCFS using ever-increasing numbers
initially.

bool flag[N]; // false } should be
int label[N]; // 0 volatile.
lock(L).



Assumption → everything happens
in main mem; if caching is involved,
problem may arise
CHECK: writing to label, reading
from label → ??

me = tid;
flag[me] = true; ~~not atomic.~~ ^{within a warp, will all threads}
label[me] = 1 + max(label); ^{get same label? to break}
white ($\exists k \neq me : flag[k] \& (label[k], k) < (label[me], me)$)
inherent loop - costly operation.

unlock()

flag[tid] = false;

b) Does fairness \Rightarrow no starvation?
CHECK.

- ① Mutual Exclusion \rightarrow preserved ✓
- ② Deadlock \rightarrow No; Strictly ' $<$ ' condition ensures total order
 \Rightarrow No cyclic dependency could arise
- ③ Starvation \rightarrow No; flag[k] would be false for uninterested threads ✓
- ④ Fairness \rightarrow FCFS will ensure fairness..

Baker's Algorithm, as we saw, is suited for CPU.

What about GPU?

Q) Does warp-based execution pose any trouble?

Blocking --- ?

As long as you don't use atomic instr, the same code works on GPU. But performance is slow.

On GPUs, locks are usually prohibited.

* No. of threads Enormous High spinning cost @ large scale.
Each thread will have to wait for long time.

Avoid locks as much as possible.

If really unavoidable & the conflicts ~~should be~~ ^{are} sparse,
then could consider using locks.

Locks can be implemented using atomics.

Atomics.

- primitive operations whose effects are visible either none or fully (never partially).
- needs hardware support
- e.g. atomic CAS, atomic Min, atomicAdd etc

Compare & Swap.

1 mem loc involved - being read & written

+ works with both global mem & shared mem ^{across threads} ^{within} _{thread block}

```
-- global -- void dkernel (int *x) {
    ++x[0];
}
```

Q) what is the value of ~~x[0]~~ at the end of this?

$x = 1 \text{ or } 2$

Thread 0

load $x[0]$, R1

inc R1

store R1, $x[0]$.

Thread 1

load $x[0]$, R2

inc R2

store R2, $x[0]$.

interleaving \Leftarrow could lead to 1/2

a) If $x[0]$ itself req. multiple instructions --- problems.

→ Ensure all-or-none behaviour.

→ Ensure all-or-none behaviour. // deterministic outcome (2)
e.g. atomicInc (& $x[0]$, ...);

② $dkernel \ll< k1, k2 \gg>$ would ensure $x[0]$ to be incremented by exactly $k1 + k2$, irrespective of the thread execution order.

$x: \cancel{x} \cancel{x} \cancel{x}$
1
2
3

PTD No
↑

check
atomicInc --;

8) does atomicInc (& $x[0]$, k);

get split
into multiple
instn - non
atomically?

print(); → all threads within a warp see single final value of x
Across warps → Build be diff.

Computation of $x[k]$ mem address - not part of atomic instruction.

Actual load, add, store is atomic

Recall. dsssp - atomic Min (tdist[n], alldist);

lost update problem solved by atomic min ✓

Every thread has to follow atomic Min ✓ otherwise chaos.

atomic CAS.

oldval = atomicCAS (&var, x, y);

$z = \text{atomicCAS}$
(&var, x, y)

Compare var, x.

If var == x, change var to y; Return x.

otherwise don't do anything; Return x

Load var, R_1
If $R_1 == x$
Store y, var
Return R_1

CAS generic atomic instr

usually available in Hardware

Others like min etc built using CAS

19-2-2020.

Typical uses cases of CAS:

- Locks: critical section processing

- Single: Only one arbitrary thread executes the block.

(Parallel \rightarrow Seq \rightarrow Parallel--)

done by a single

- Other atomic variants: - atomic min etc can be implemented using atomic CAS

Ex. Implement lock with atomic CAS.

$\text{curr} = 0$; // some initialization.

lock()

$\text{mycurr} = \text{atomicCAS}(\&\text{curr}, 0, \text{tid});$ $\text{mycurr} = \text{atomicCAS}(\&\text{curr}, 0, \text{tid}, 0)$

If $\text{curr} = 0$, then curr gets replaced by tid, 0 is returned
⇒ currently tid has acquired lock
so can do "rdm" etc.

If $\text{curr} \neq 0$, then whatever that value is, is returned

lock()

do {

$\text{mycurr} = \text{atomicCAS}(\&\text{curr}, 0, \text{tid}),$ // try trying to
acquire lock.

If $\text{curr} (\text{mycurr}) = 0$

// comes here only if I am able to change curr from 0 to tid.

CS

≡

unlock()

$\text{mycurr} = \text{atomicCAS}(\&\text{curr}, \text{tid}, 0),$ // replace curr value
with 0, so that
some other thread
can acquire

$\text{curr} = 0$ is program

CW: Implement single with atomic CAS.

s1

single {

cs

}

s2

All threads execute s1, s2. Exactly one of the threads executes cs. While cs is being executed, other threads could be executing either s1 or s2.

lockvar = 0; // initial value

s1 // if threads can't get lockvar, go to s2

if (atomicCAS(&lockvar, 0, 1) == 1) then { // doesn't block
other threads from going ahead to

cs

3

2

s2

Also, since we are not unlocking, no other thread can go into the if

If single is inside for loop.

Meet the lockvar only after all the threads have finished s2. // barrier reqd

Remember:

Atomic instructions do not behave like a barrier.

Multiple 'singles' → multiple lock variables.

Q: What's the output? - not deterministic.

```
#include <stdio.h>
```

```
#include <cuda.h>
```

```
__global__ void k1(int *gg) {
```

```
    int old = atomicCAS(gg, 0, threadIdx.x + 1);
```

```
    if (old == 0) { → executed by exactly one thread.
```

```
        printf("Thread %d succeeded 1.\n", threadIdx.x);
```

```
}
```

```
    old = atomicCAS(gg, 0, threadIdx.x + 1);
```

```
    if (old == 0) { → will not be executed by any thread
```

```
        printf("Thread %d succeeded 2.\n", threadIdx.x);
```

```
}
```

```
    old = atomicCAS(gg, threadIdx.x, -1);
```

```
    if (old == threadIdx.x) {
```

```
        printf("Thread %d succeeded 3.\n", threadIdx.x);
```

```
}
```

```
}
```

```
int main() {
```

```
    int *gg;
```

```
    cudaMalloc(&gg, sizeof(int));
```

```
    cudaMemset(&gg, 0, sizeof(int));
```

```
    k1 << 2, 32 >> (gg); → 2 warps.
```

```
    cudaDeviceSynchronize();
```

```
    return 0;
```

```
}
```

Even of warp.
Threads are
executing an
atomic instruction.

only one of the
threads is going to
do it at a time.