

CUDA Memory Model

Or how “our strength grows out of our weakness”

by Georgy Evtushenko @g_evtushenko

on January 5, 2021

» Motivation

Performance A significant slowdown is caused by excess memory ordering

Correctness One source of non-deterministic concurrency bugs is a misunderstanding of the architecture's shared memory consistency model i.e. what values can be read from shared memory when issued concurrently with other reads and writes.¹

¹Daniel J. Sorin, Mark D. Hill, and David A. Wood. *A Primer on Memory Consistency and Cache Coherence*. 1st. 2011. ISBN: 1608455645.

» Definitions

Read operation All variants of **ld** instruction and **atom** instruction.

Write operation All variants of **st** instruction and **atom** instruction.

Memory operation A **read** or **write** operation.

Memory model Set of rules defining how **memory operations** on shared memory from multiple threads are processed. Tradeoff between ease of programmability and efficiency.

Memory order Total order of **memory operations**.

» Sequential consistency

“... the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program.”²

Thread 1

```
1  x = 1
2  r1 = y
```

Thread 2

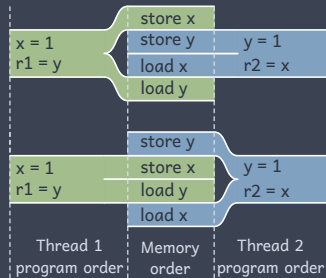
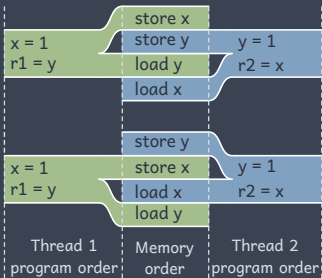
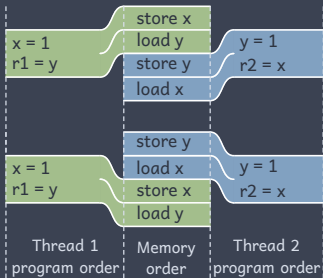
```
1  y = 1
2  r2 = x
```

²L. Lamport. “How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs”. In: *IEEE Transactions on Computers* C-28 (1979), pp. 690–691.

» Sequential consistency

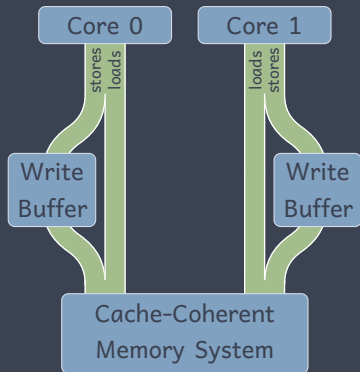
Memory order

$$A <_p B \Rightarrow A <_m B$$



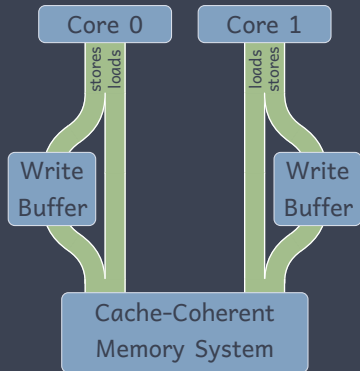
» Total Store Order

TSO



TSO/x86 memory model allows some non-SC executions by holding stores in write buffers:

- * A store enters the write buffer immediately
- * A store exits the write buffer when cache in a read-write coherence state



TSO/x86 memory model allows some non-SC executions by holding stores in write buffers:

- * A store enters the write buffer immediately
- * A store exits the write buffer when cache in a read-write coherence state

Thread 1

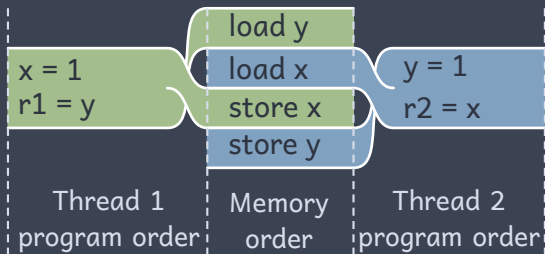
```
1 x = 1;
2 r1 = y;
```

Thread 2

```
1 y = 1;
2 r2 = x;
```

» Total Store Order

TSO



- * $L(a) <_p L(b) \Rightarrow L(a) <_m L(b)$
- * $L(a) <_p S(b) \Rightarrow L(a) <_m S(b)$
- * $S(a) <_p S(b) \Rightarrow S(a) <_m S(b)$
- * $S(a) <_p L(b) \not\Rightarrow S(a) <_m L(b)$
- * Therefore it's possible to have both $r_1 = 0$ and $r_2 = 0$

» Program order

Reordering

C++

```
1 x = y + 1;  
2 y = 0;
```

```
1 g++ -O3 -S -masm=intel ctmo.cpp
```

ASM

```
1 mov eax, DWORD PTR Y[rip]  
2 mov DWORD PTR Y[rip], 0  
3 add eax, 1  
4 mov DWORD PTR X[rip], eax
```

» Program order

Compiler barrier

C++

```
1 x = y + 1;  
2 asm volatile ("" ::: "memory");  
3 y = 0;  
  
1 g++ -O3 -S -masm=intel ctmo.cpp
```

ASM

```
1 mov eax, DWORD PTR Y[rip]  
2 add eax, 1  
3 mov DWORD PTR X[rip], eax  
4 mov DWORD PTR Y[rip], 0
```

» Program order

Compiler barrier

Thread 1

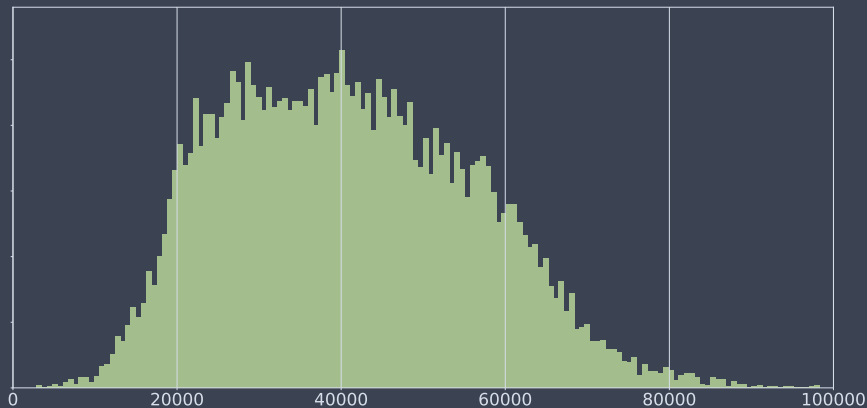
```
1 void thread_1_body ()
2 {
3     x = 1;
4     asm volatile (
5         "" ::: "memory");
6     r1 = y;
7 }
```

Thread 2

```
1 void thread_2_body ()
2 {
3     y = 1;
4     asm volatile (
5         "" ::: "memory");
6     r2 = x;
7 }
```

» Load after Store violations count

i7-7700k



» Program order

Compiler barrier

Thread 1

```
1 void thread_1_body ()
2 {
3     x = 1;
4     asm volatile (
5         "mfence" ::: "memory");
6     r1 = y;
7 }
```

Thread 2

```
1 void thread_2_body ()
2 {
3     y = 1;
4     asm volatile (
5         "mfence" ::: "memory");
6     r2 = x;
7 }
```

Weak memory models preserves only the orders that programmers require.

memory operations can be reordered unless there is a **fence** between them.

- * $L(a) <_p \text{FENCE} \Rightarrow L(a) <_m \text{FENCE}$
- * $S(a) <_p \text{FENCE} \Rightarrow S(a) <_m \text{FENCE}$
- * $\text{FENCE} <_p L(a) \Rightarrow \text{FENCE} <_m L(a)$
- * $\text{FENCE} <_p S(a) \Rightarrow \text{FENCE} <_m S(a)$

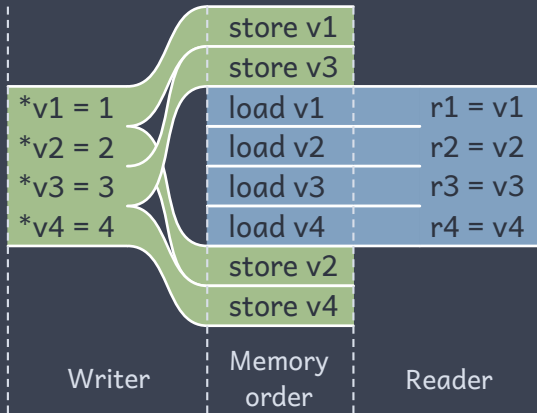
» Operations interleaving

```
1  __device__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1;  
8      *v2 = 2;  
9      *v3 = 3;  
10     *v4 = 4;  
11 }
```

```
1  __device__ void reader (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4,  
6      int *result)  
7  {  
8      int r1 = *v1; int r2 = *v2;  
9      int r3 = *v3; int r4 = *v4;  
10  
11     if (    r1 == 1 && r3 == 3  
12         && r2 == 0 && r4 == 0)  
13         *result = 0;  
14 }
```

» Operations interleaving

TSO



» Operations interleaving

Was interleaved up to **671** times out of 1000 runs on RTX 3090.

Writer		Reader	
1	IMAD.MOV.U32 R0, RZ, RZ, 0x1;	1	LDG.E.STRONG.SYS R0, [UR4];
2	IMAD.MOV.U32 R2, RZ, RZ, 0x2;	2	LDG.E.STRONG.SYS R2, [UR6];
3	IMAD.MOV.U32 R3, RZ, RZ, 0x3;	3	LDG.E.STRONG.SYS R3, [UR8];
4	IMAD.MOV.U32 R4, RZ, RZ, 0x4;	4	LDG.E.STRONG.SYS R4, [UR10];
5	STG.E.STRONG.SYS [UR4], R0;		
6	STG.E.STRONG.SYS [UR6], R2;		
7	STG.E.STRONG.SYS [UR8], R3;		
8	STG.E.STRONG.SYS [UR10], R4;		

» Operations interleaving

```
1  cudaMemset (x, 0, 4 * sizeof (int));  
2  cudaMemset (y, 0, 2 * sizeof (int));  
3  
4  if (single_segment)  
5      kernel<<<1024, 32>>> (x + 0, x + 1, x + 2, x + 3, r);  
6  else  
7      kernel<<<1024, 32>>> (x + 0, y + 0, x + 1, y + 1, r);
```

» Write Buffer

```
1  __host__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {      // <=  
7      *v1 = 1;  
8      *v2 = 2;  
9      *v3 = 3;  
10     *v4 = 4;  
11 }
```

Entry: single word

addr	value

» Write Buffer

```
1  __host__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1; // <=  
8      *v2 = 2;  
9      *v3 = 3;  
10     *v4 = 4;  
11 }
```

Entry: single word

addr	value
200	1

» Write Buffer

```
1  __host__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1;  
8      *v2 = 2; // <=  
9      *v3 = 3;  
10     *v4 = 4;  
11 }
```

Entry: single word

addr	value
200	1
400	2

» Write Buffer

```
1  __host__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1;  
8      *v2 = 2;  
9      *v3 = 3; // <=  
10     *v4 = 4;  
11 }
```

Entry: single word

addr	value
200	1
400	2
204	3

» Write Buffer

```
1  __host__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1;  
8      *v2 = 2;  
9      *v3 = 3;  
10     *v4 = 4; // <=  
11 }
```

Entry: single word

addr	value
200	1
400	2
204	3
404	4

» Merging Write Buffer

```
1  __device__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {      // <=  
7      *v1 = 1;  
8      *v2 = 2;  
9      *v3 = 3;  
10     *v4 = 4;  
11 }
```

Entry: cache block

addr	value	value	value	value

» Merging Write Buffer

```
1  __device__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1; // <=  
8      *v2 = 2;  
9      *v3 = 3;  
10     *v4 = 4;  
11 }
```

Entry: cache block

addr	value	value	value	value
200	1			

» Merging Write Buffer

```
1  __device__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1;  
8      *v2 = 2; // <=  
9      *v3 = 3;  
10     *v4 = 4;  
11 }
```

Entry: cache block

addr	value	value	value	value
200	1			
400	2			

» Merging Write Buffer

```
1  __device__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1;  
8      *v2 = 2;  
9      *v3 = 3; // <=  
10     *v4 = 4;  
11 }
```

Entry: cache block

addr	value	value	value	value
200	1	3		
400	2			

» Merging Write Buffer

```
1  __device__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1;  
8      *v2 = 2;  
9      *v3 = 3;  
10     *v4 = 4; // <=  
11 }
```

Entry: cache block

addr	value	value	value	value
200	1	3		
400	2	4		

» Operations interleaving

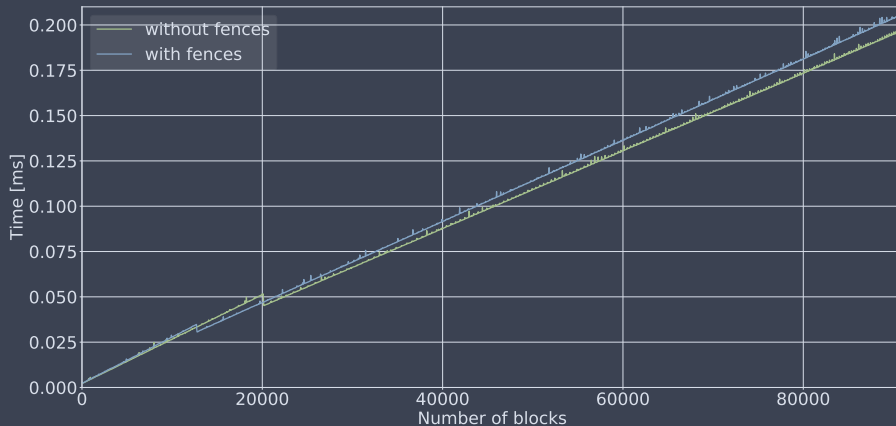
fixed version

```
1  __device__ void writer (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4)  
6  {  
7      *v1 = 1;  
8      *v2 = 2;  
9      __threadfence ();  
10     *v3 = 3;  
11     *v4 = 4;  
12 }
```

```
1  __device__ void reader (  
2      volatile int * v1,  
3      volatile int * v2,  
4      volatile int * v3,  
5      volatile int * v4,  
6      int *result)  
7  {  
8      int r1 = *v1; int r2 = *v2;  
9      __threadfence ();  
10     int r3 = *v3; int r4 = *v4;  
11  
12     if (    r1 == 1 && r3 == 3  
13         && r2 == 0 && r4 == 0)  
14         *result = 0;  
15 }
```

» Operations interleaving

RTX 3090



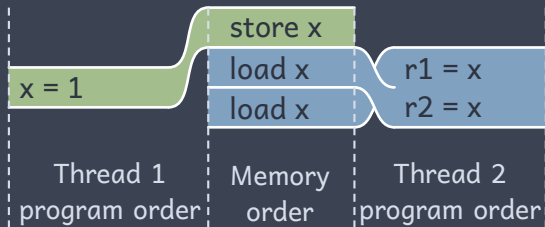
» Read-Read Coherence

coRR

```
1  __global__ void coherence_of_read_read (/*...*/) {
2  __shared__ int cache[BLOCK_SIZE];
3  /* .. Initialize shared memory .. */
4  __syncthreads ();
5  if (tid == wthread) {
6      cache[tid] = 1;
7  }
8  else {
9      const int r1 = cache[reordering_1[tid]];
10     const int r2 = cache[reordering_2[tid]];
11     if (r1 == 1 && r2 == 0)
12         if (tid == rthread)
13             *violated = 1;
14 }
15 }
```

» Read-Read Coherence

coRR



Violated **157** times out of 100'000 runs on GTX 560 (Fermi). Fixed in Maxwell.

Thread 1

```
1 MOV32I R0, 0x1;  
2 STS [R7], R0;
```

Thread 2

```
1 LDS R0, [R0];  
2 LDS R4, [R4]; /* [R4] = [R0] */  
3 ISETP.EQ.AND P0, PT, R0, 0x1, PT;  
4 ISETP.EQ.AND P0, PT, R4, RZ, P0;
```

» Message Passing

```
1 __global__ void kernel (  
2     int n, int * flag, int * data, int *result) {
```

Writer

```
1 for (int i = 0; i < n; i++)  
2     data[i] = i + 1;  
3  
4 *flag = 1;
```

Reader

```
1 while (*flag == 0);  
2  
3 for (int i = 0; i < n; i++)  
4     if (data[i] == 0)  
5         *result = 0;
```

» Message Passing

```
1  __global__ void kernel (  
2      int n, int * flag, int * data, int *result) {
```

Writer

```
1  for (int i = 0; i < n; i++)  
2      data[i] = i + 1;  
3  
4  *flag = 1;
```

SASS

```
1  IADD R3, R0, 0x1;  
2  IADD R4, R0, 0x2;  
3  STG.E [R1], R3  
4  IADD R3, R0, 0x3;  
5  STG.E [R2], R4;  
6  ; /* ... */  
7  MOV32I R5, 0x1;  
8  STG.E [flag], R5;
```

» Message Passing

```
1 __global__ void kernel (  
2     int n, int * flag, int * data, int *result) {
```

Reader

```
1 while (*flag == 0);  
2  
3 for (int i = 0; i < n; i++)  
4     if (data[i] == 0)  
5         *result = 0;
```

SASS

```
1 LDG.E R1, [R0];  
2 ISETP.NE.AND P0, PT, R2, RZ, PT;  
3 @!P0 BRA .L1;  
4 LDG.E R3, [R2];  
5 LDG.E R4, [R2+0x4];  
6 LDG.E R5, [R2+0x8];  
7 LDG.E R6, [R2+0xc];  
8 ; /* ... */  
9 EXIT;  
10 .L1:  
11 BRA .L1;
```

» Message Passing

first fix

```
1  __global__ void kernel (  
2      int n, volatile int * flag, int * data, int *result) {
```

Reader

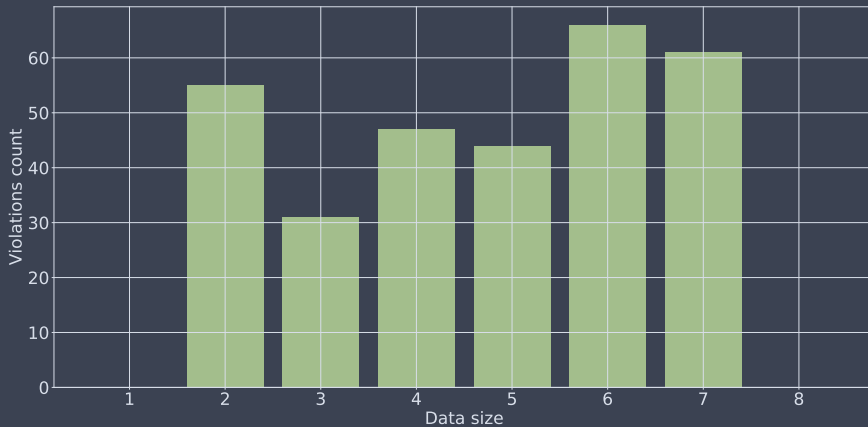
```
1  while (*flag == 0);  
2  
3  for (int i = 0; i < n; i++)  
4      if (data[i] == 0)  
5          *result = 0;
```

SASS

```
1  .L1:  
2  LDG.E.STRONG.SYS R1, [R0];  
3  ISETP.NE.AND P0, PT, R2, RZ, PT;  
4  @!P0 BRA .L1;  
5  LDG.E R3, [R2];  
6  LDG.E R4, [R2+0x4];  
7  LDG.E R5, [R2+0x8];  
8  LDG.E R6, [R2+0xc];  
9  ; /* ... */  
10 EXIT;
```

» Stale data found

RTX 2080



» Message Passing

final fix

```
1 __global__ void kernel (  
2     int n, volatile int * flag, int * data, int *result) {
```

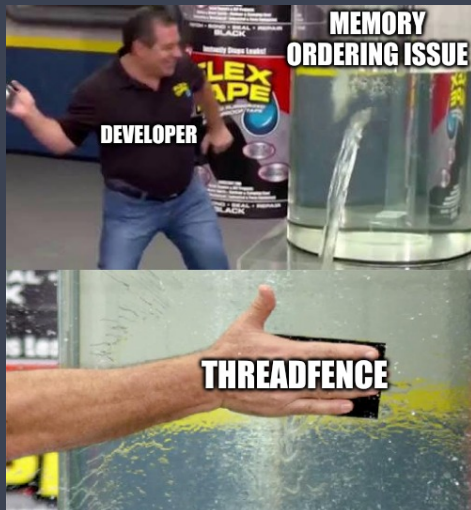
Writer

```
1 for (int i = 0; i < n; i++)  
2     data[i] = i + 1;  
3 __threadfence ();  
4  
5 *flag = 1;
```

Reader

```
1 while (*flag == 0);  
2 __threadfence ();  
3  
4 for (int i = 0; i < n; i++)  
5     if (data[i] == 0)  
6         *result = 0;
```

» Message Passing



» Message Passing

Controversial version

```
1  __device__ int ldg (const int * p)
2  {
3      int out;
4      asm volatile("ld.global.cg.s32 %0, [%1];" : "=r"(out) : "l"(p));
5      return out;
6  }
```

Writer

```
1  for (int i = 0; i < n; i++)
2      data[i] = i + 1;
3  __threadfence ();
4
5  *flag = 1;
```

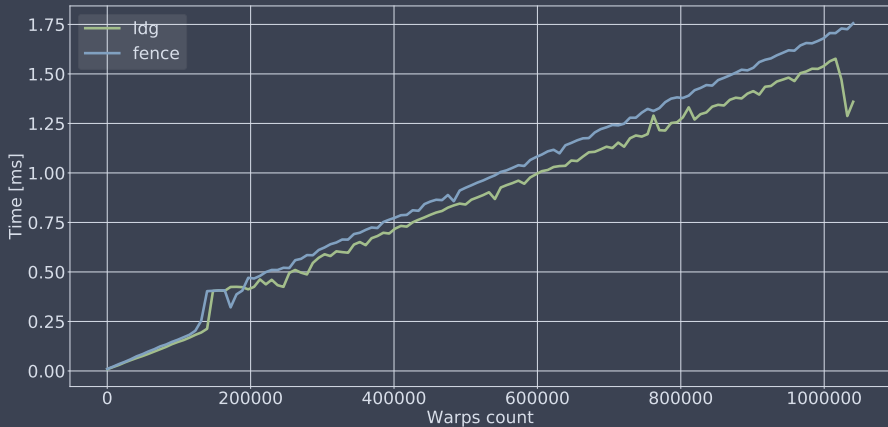
Reader

```
1  while (*flag == 0);
2  // __threadfence ();
3
4  for (int i = 0; i < n; i++)
5      if (ldg (data + i) == 0)
6          *result = 0;
```

» Message Passing

n=1

Average speedup is about 8.4%.



» Message Passing

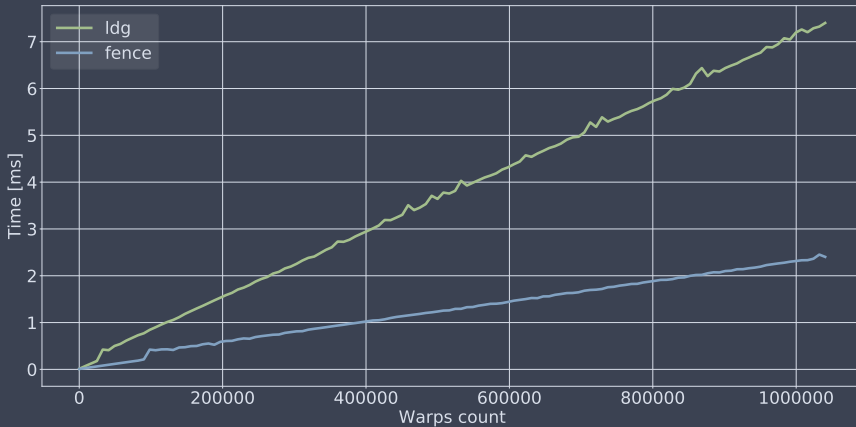
“It’s possible to break the dependency order between the flag load and the later, dependent ld.cg load, that makes your test pass. Because it’s possible to write a program that breaks this order, it’s not possible for us to write a memory model that says the order is enforced.”³

³Olivier Giroux. personal communication.

» Message Passing

n=9

Average slowdown is about 208.3%.



» Message Passing

The fence can be eliminated if the status flag and the corresponding value being updated can be combined into a single architectural word.⁴

⁴M. Garland D. Merrill. “Single-pass Parallel Prefix Scan with Decoupled Look-back”. In: *NVIDIA Technical Report NVR-2016-002* (2016).

» Message Passing

Single word

```
1 union data_flag
2 {
3     struct { int32_t data; int32_t flag; } fields;
4     int64_t vec;
5 };
```

Writer

```
1 data_flag tmp;
2 tmp.fields.data = 1;
3 tmp.fields.flag = 1;
4
5 df->vec = tmp.vec;
```

Reader

```
1 while (!tmp.fields.flag)
2     tmp.vec = df->vec;
3
4 if (tmp.fields.data == 0)
5     *result = 0;
```

» Message Passing

Single word



» Exclusive scan

sequential version

Scan produces an output sequence where each element is computed to be the reduction of the elements occurring earlier in the input sequence.

```
1  int sum = 0;
2  for (int i = 0; i < n; i++)
3      {
4          out[i] = sum;
5          sum += in[i];
6      }
```

in	1	0	1	1	0	1
out	0	1	1	2	3	3

» Exclusive scan

inner scan

in	1	0	1	1	0	1
bs	0	1	0	1	0	0

» Exclusive scan

blocks' results

in	1	0	1	1	0	1
bs	0	1	0	1	0	0
		1		2		1

» Exclusive scan

inclusive scan of blocks' results

in	1	0	1	1	0	1
bs	0	1	0	1	0	0
		1		2		1
		1		3		4

» Exclusive scan

final result

in	1	0	1	1	0	1
bs	0	1	0	1	0	0
		1		2		1
		1		3		4
out	0	1	1	2	3	3

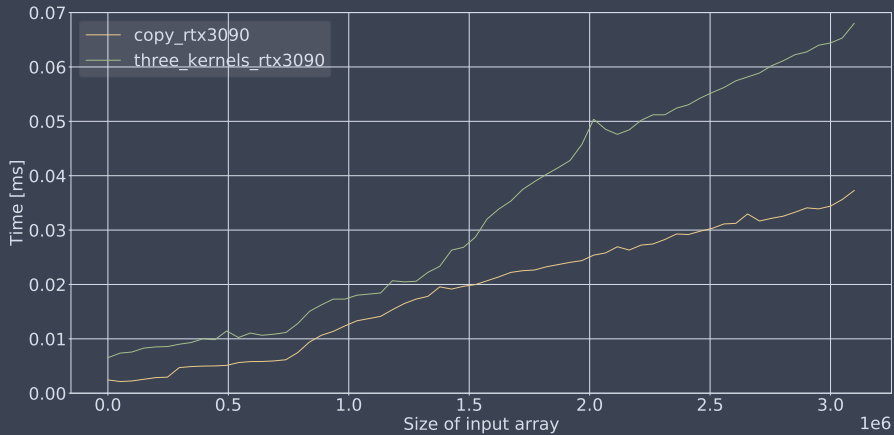
» Exclusive scan

three-kernels version

```
1 hierarchical_partial_scan<<<blocks, hierarchical_block_size>>> (  
2     in, out, helper_buffer);  
3  
4 hierarchical_blocks_results_scan<<<1, 1024>>> (  
5     blocks, helper_buffer, helper_buffer);  
6  
7 hierarchical_final_adjust<<<blocks, hierarchical_block_size>>> (  
8     out, helper_buffer);
```

» Exclusive scan

three-kernels scan performance



» Exclusive scan

sequential version

```
1 // Load
2 data_type thread_data[thread_data_size];
3
4 BlockLoad (temp_storage).Load (in + block_offset, thread_data);
5 __syncthreads ();
6
7 // This thread block aggregate
8 data_type block_reduce =
9     BlockReduce (temp_storage).Sum (thread_data);
10
11 /// ...
```

» Exclusive scan

sequential version

```
1  if (threadIdx.x == 0) {
2      if (blockIdx.x == 0) {
3          commit (block_statuses + blockIdx.x, block_reduce);
4      }
5      else {
6          do {
7              prev_block_status = load (block_statuses + blockIdx.x - 1);
8              } while (prev_block_status.flag == 0);
9
10         thread_data[0] += prev_block_status.data;
11         commit (
12             block_statuses + blockIdx.x,
13             prev_block_status.data + block_reduce);
14     }
15 }
16 __syncthreads ();
```

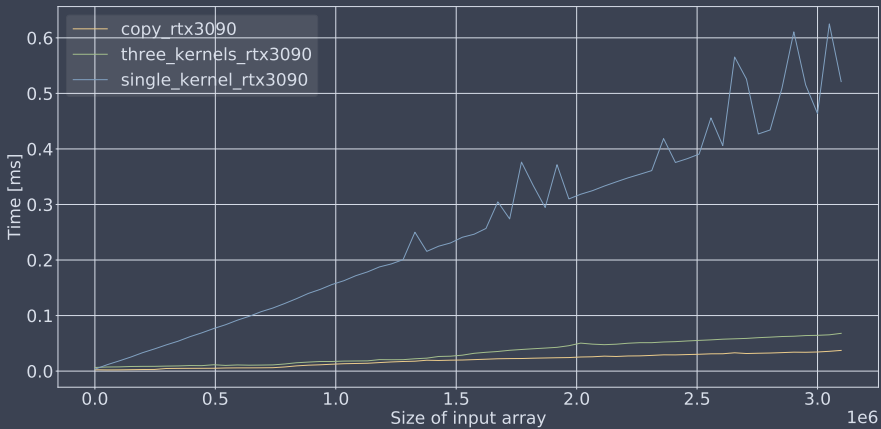

» Exclusive scan

sequential version

```
1 BlockScan (temp_storage).ExclusiveSum (thread_data, thread_data);
2 __syncthreads ();
3
4 // Store final result
5 BlockStore (temp_storage).Store (out + block_offset, thread_data);
```

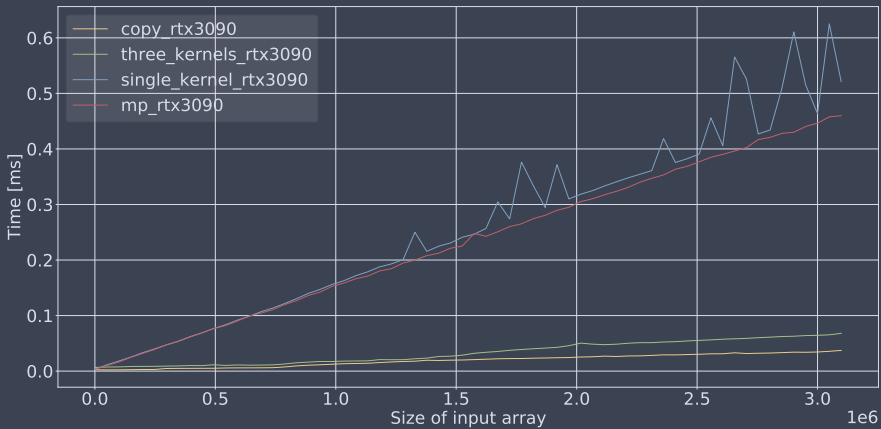
» Exclusive scan

single scan performance

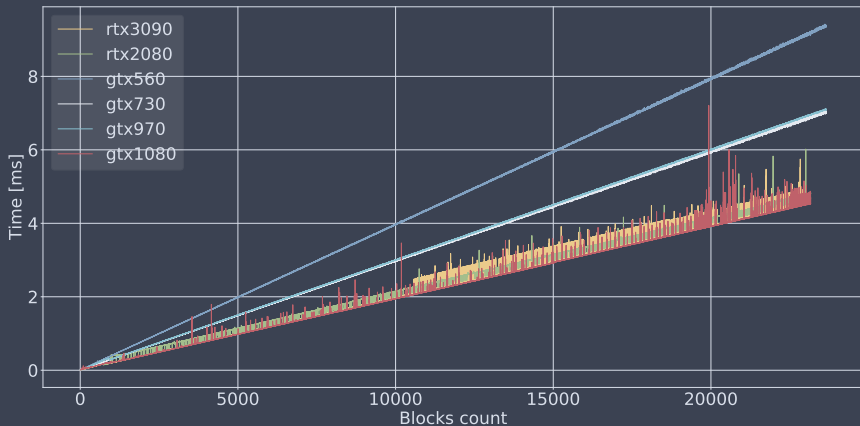


» Exclusive scan

single scan performance



» Message Passing



» Exclusive scan

decoupled look-back

```
1 if (blockIdx.x == 0)
2   if (threadIdx.x == 0)
3     store (block_statuses + blockIdx.x, 2, block_aggregate);
```

» Exclusive scan

decoupled look-back

```
1 if (blockIdx.x != 0)
2 {
3     if (threadIdx.x == 0)
4         store (block_statuses + blockIdx.x, 1, block_aggregate);
5
6     if (threadIdx.x / WARP_SIZE == 0)
7     {
8         // ...
```

» Exclusive scan

decoupled look-back

```
1 int predecessor_idx = blockIdx.x - WARP_SIZE + threadIdx.x;
2 data_type thread_sum = data_type ();
3
4 while (true)
5 {
6     prev_block_status = load (block_statuses + predecessor_idx);
7
8     while (prev_block_status.flag == 0)
9         prev_block_status = load (block_statuses + predecessor_idx);
```

» Exclusive scan

decoupled look-back

```
1  if (__all_sync (WARP_MASK, prev_block_status.flag < 2)) {
2      thread_sum += prev_block_status.data;
3      predecessor_idx -= WARP_SIZE;
4      continue;
5  }
6  else {
7      const int final_mask =
8          __ballot_sync (WARP_MASK, prev_block_status.flag > 1);
9
10     const int rightmost_thread = WARP_SIZE - 1 - __clz (final_mask);
11
12     thread_sum += threadIdx.x % WARP_SIZE < rightmost_thread
13                 ? 0
14                 : prev_block_status.data;
15     break;
16 }
```


» Exclusive scan

decoupled look-back

```
1 data_type aggregate =  
2   WarpReduce (temp_storage.warp_reduce).Sum (thread_sum);  
3  
4 if (threadIdx.x == 0) {  
5   store (block_statuses + blockIdx.x, 2, aggregate + block_aggregate);  
6   prev_aggregate = aggregate;  
7 }
```

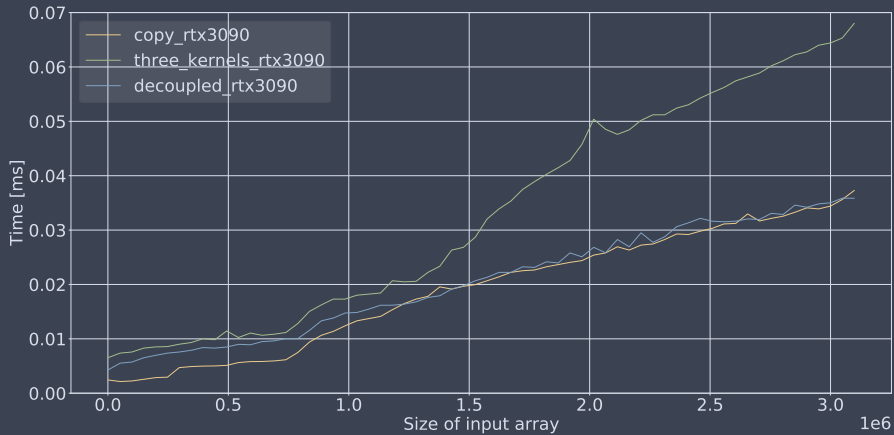
» Exclusive scan

decoupled look-back

```
1      // ...
2  }
3  }
4  __syncthreads ();
5
6  const data_type reg_prev_aggregate = prev_aggregate;
7  for (int i = 0; i < thread_data_size; i++)
8      thread_data[i] += reg_prev_aggregate;
9
10 BlockStore (temp_storage).Store (out + block_offset, thread_data);
```

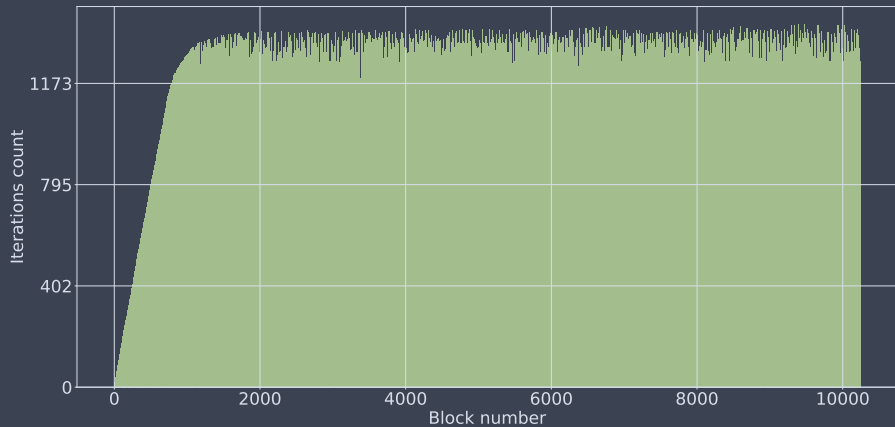
» Exclusive scan

decoupled look-back



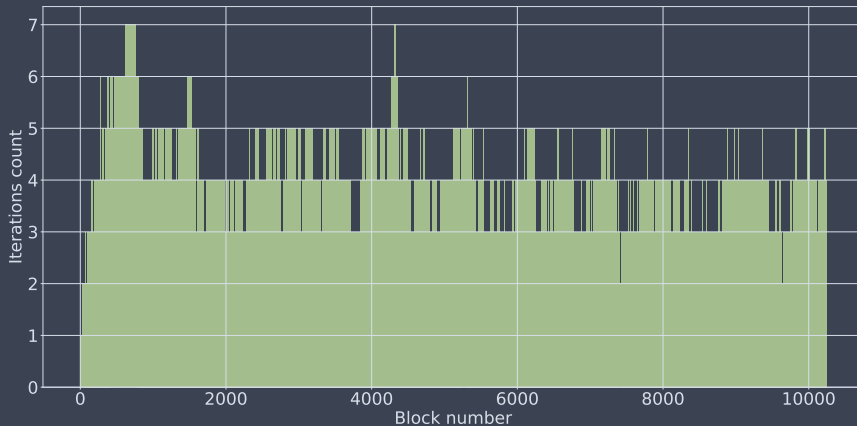
» Exclusive scan

single-kernel



» Exclusive scan

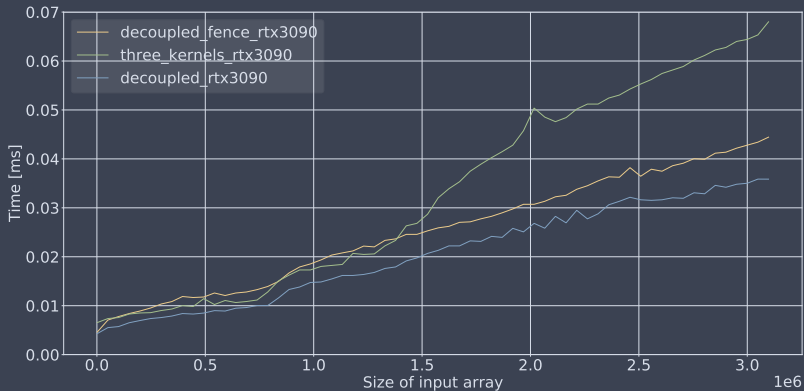
decoupled look-back



» Exclusive scan

decoupled look-back with fence

Fence version: 23% performance degradation



» Release and Acquire Patterns

fence

Release pattern on a location M consists of fence followed by a strong write on M in program order.

```
1 __threadfence ();  
2 *flag = 1; // volatile int *flag
```

Acquire pattern on a location M consists of strong read on M followed by a fence in program order.

```
1 int reg = *flag; // volatile int *flag  
2 __threadfence ();
```

» Definitions

Read operation All variants of **ld** instruction and **atom** instruction.

Write operation All variants of **st** instruction and **atom** instruction.

Memory operation A **read** or **write** operation.

» Definitions

Read operation All variants of `ld` instruction and `atom` instruction.

Write operation All variants of `st` instruction and `atom` instruction.

Memory operation A `read` or `write` operation.

Strong operation A `fence` operation, or a memory operation with a `.relaxed`, `.acquire`, `.release`, `.acq_rel` or `.volatile` qualifier.

» Release and Acquire Patterns

acquire/release operations

Release pattern on a location M consists of a release operation on M.

```
1 __device__ void store_release (int * flag, int val)
2 {
3     asm volatile (
4         "st.release.gpu.b32 [%0], %1;"
5         ":: "l"(flag), "r"(val) : "memory");
6 }
```

Acquire pattern on a location M consists of an acquire operation on M.

```
1 __device__ int load_acquire (int * flag)
2 {
3     int reg;
4     asm volatile (
5         "ld.acquire.gpu.b32 %0, [%1];"
6         : "=r"(reg) : "l"(flag) : "memory");
7     return reg;
8 }
```

» Release and Acquire Patterns

acquire/release operations

st.release.gpu

```
1 MEMBAR.ALL.GPU
2 ST.E.STRONG.GPU [UR4], R0
```

ld.acquire.gpu

```
1 LD.E.STRONG.GPU R0, [UR4]
2 CCTL.IVALL
```

» Release and Acquire Patterns

libcu++

```
1  __device__ void old_writer (  
2      int n, volatile int * flag, int * data, int *result) {  
3      for (int i = 0; i < n; i++)  
4          data[i] = i + 1;  
5  
6      __threadfence ();  
7      *flag = 1;  
8  }
```

```
1  __device__ void new_writer (  
2      int n, atomic<int> &flag, int * data, int *result) {  
3      for (int i = 0; i < n; i++)  
4          data[i] = i + 1;  
5  
6      flag.store (1, memory_order_release);  
7  }
```

» Reduce

sequential version

Reduce produces a single aggregate from a list of input elements.

```
1 int sum = 0;
2 for (int i = 0; i < n; i++)
3     sum += in[i];
4 *out = sum;
```



» Device reduction

```
1 reduce<<<blocks, block_size>>> (in, block_results);  
2 reduce_block_results<<<1, block_size>>> (  
3     blocks, block_results, block_results);
```

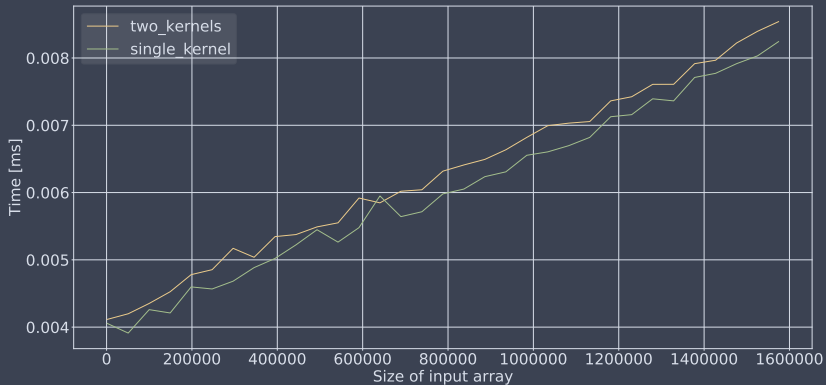
» Device reduction

```
1 // .. block-wide reduce ..
2
3 if (threadIdx.x == 0)
4 {
5     block_results[blockIdx.x] = block_result;
6     __threadfence ();
7
8     const int prev_count = atomicAdd (count, 1);
9     temp_storage.need_to_perform_final_reduce
10     = prev_count == gridDim.x - 1;
11 }
```

» Device reduction

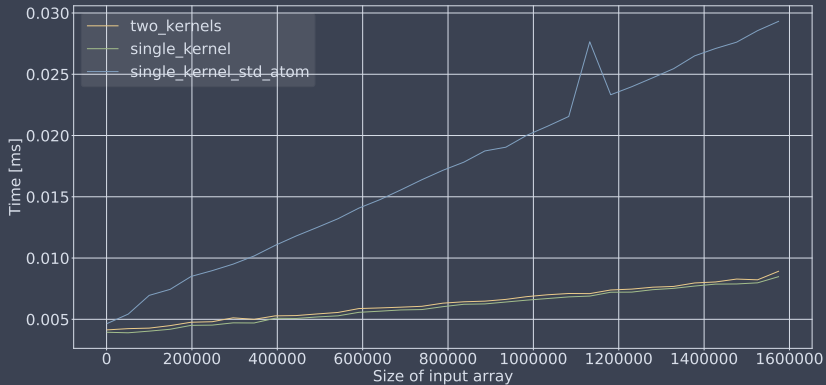
performance

Single-kernel version: $\approx 4\%$ performance improvement



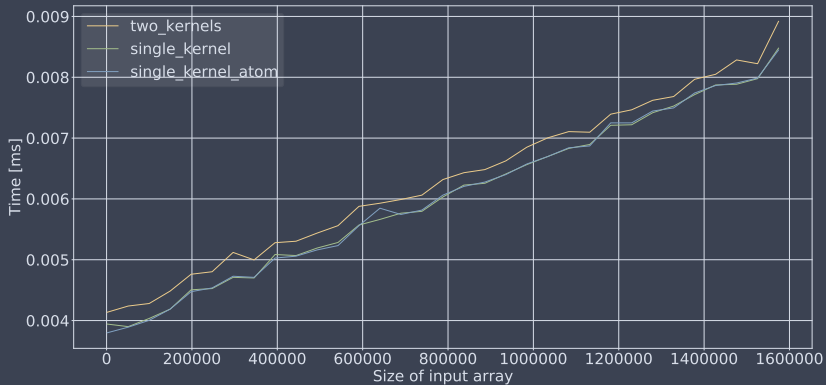
» Device reduction

```
1 // .. block-wide reduce ..
2
3 if (threadIdx.x == 0)
4 {
5     block_results[blockIdx.x] = block_result;
6
7     const int prev_count =
8         count->fetch_add (1, cuda::std::memory_order_release);
9
10    temp_storage.need_to_perform_final_reduce
11        = prev_count == gridDim.x - 1;
12 }
```



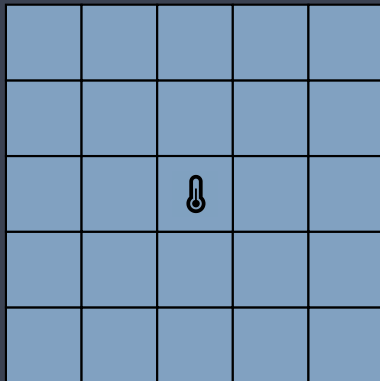
» Device reduction

```
1 // CUDA C++, __host__ __device__.
2 // Strictly conforming to the C++ Standard.
3 #include <cuda/std/atomic>
4 cuda::std::atomic<int> x;
5
6 // CUDA C++, __host__ __device__.
7 // Conforming extensions to the C++ Standard.
8 #include <cuda/atomic>
9 cuda::atomic<int, cuda::thread_scope_block> x;
```



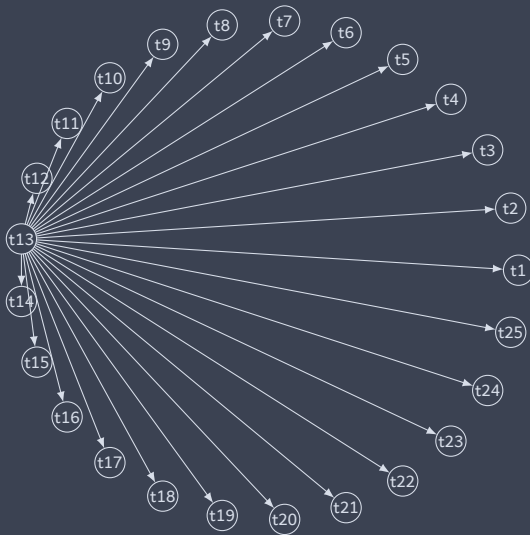
» Broadcast

illustration



» Broadcast

illustration



» Broadcast

basic version

```
1 unsigned int special_value = (unsigned int) -1;
2 unsigned int last_it = special_value;
3 int stride = blockDim.x * gridDim.x;
4
5
6 for (unsigned it = 0; it < last_it; it++) {
7     int i = threadIdx.x + blockIdx.x * blockDim.x;
8     for (; i < n; i += stride) {
9         data[i] += value;
10        if (data[i] > threshold)
11            stop = true;
12    }
13
14    if (stop)
15        last_it = it;
16 }
```

» Broadcast

volatile global

sensor owner

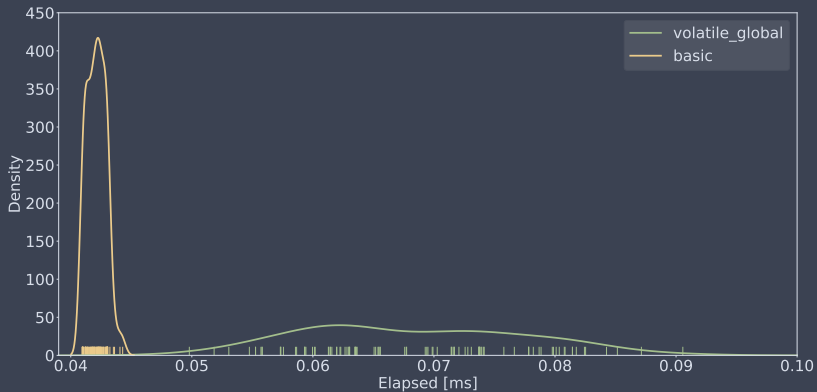
```
1 if (data[i] > threshold)
2   stop = true;
3
4 store (flag, stop, it + 1);
```

other threads

```
1 if (special_value == last_it) {
2   status_word it_state
3   = load (flag);
4
5   while (it_state.data <= it)
6     it_state = load (flag);
7
8   if (it_state.flag)
9     last_it = it_state.data;
10 }
```


» Broadcast

volatile global performance



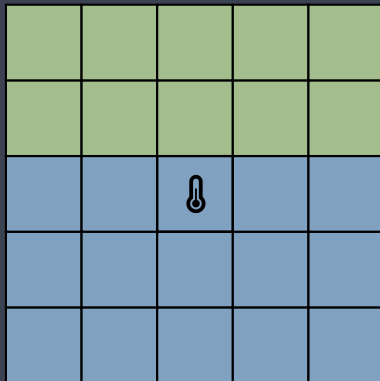
» Broadcast

volatile global block

```
1  if (last_it == special_value) {
2      if (threadIdx.x == 0) {
3          status_word it_state = load (flag);
4
5          while (it_state.data <= it)
6              it_state = load (flag);
7
8          if (it_state.flag)
9              block_last_it = it_state.data;
10     }
11     __syncthreads ();
12
13     last_it = block_last_it;
14 }
```

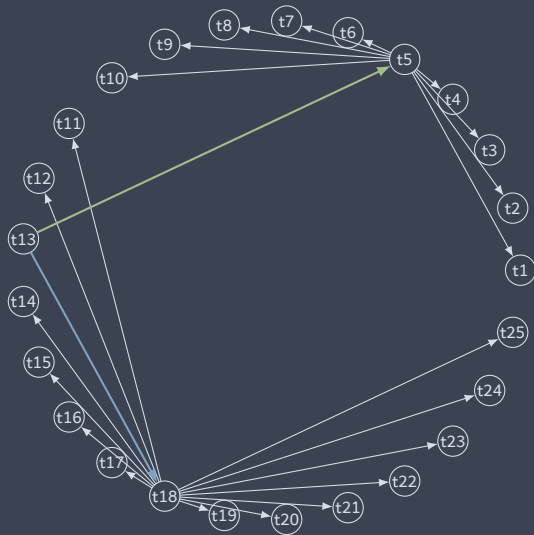
» Broadcast

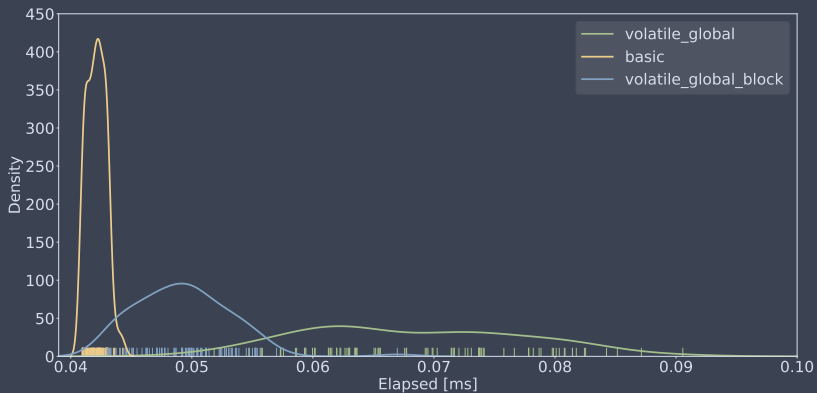
illustration



» Broadcast

illustration





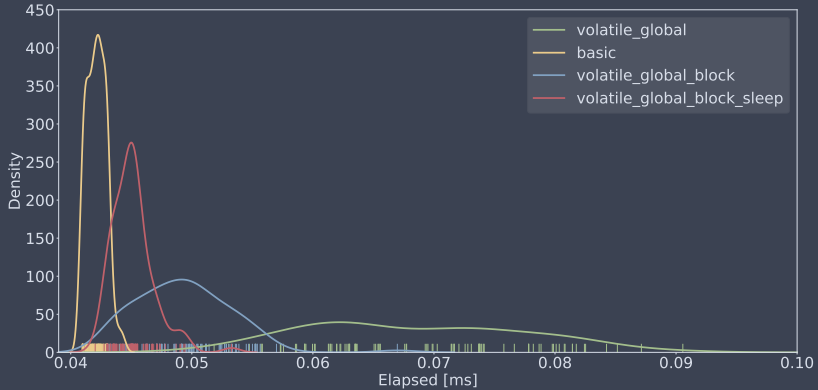
» Broadcast

volatile global block sleep

```
1  if (last_it == special_value) {
2      if (threadIdx.x == 0) {
3          status_word it_state = load (flag);
4
5          while (it_state.data <= it) {
6              __nanosleep (250); //< new line
7              it_state = load (flag);
8          }
9
10         if (it_state.flag)
11             block_last_it = it_state.data;
12     }
13     __syncthreads ();
14
15     last_it = block_last_it;
16 }
```

» Broadcast

volatile global block sleep performance



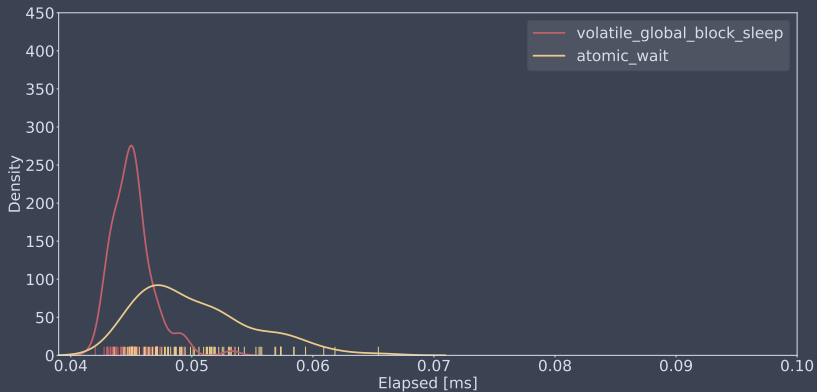
» Broadcast

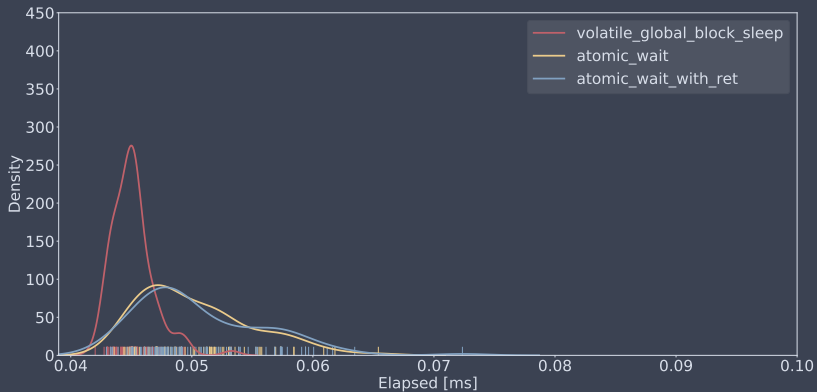
atomic (sensor owner)

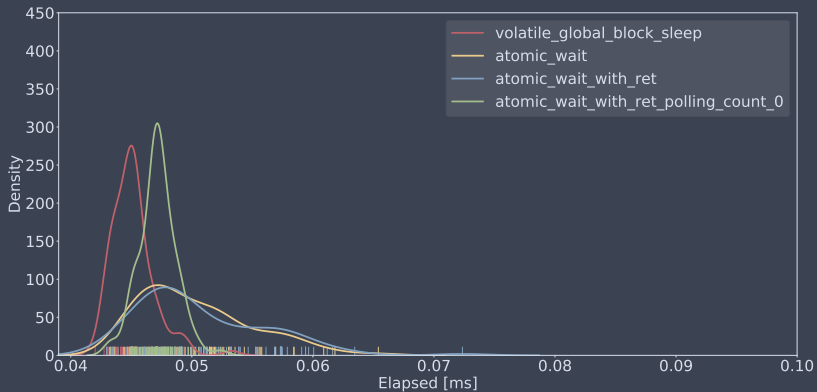
```
1 if (data[i] > threshold)
2     stop = true;
3
4 flag->store (
5     status_word {stop, it + 1},
6     cuda::memory_order_relaxed);
7 flag->notify_all ();
```


other threads

```
1  if (threadIdx.x == 0) {
2      status_word it_state
3      = flag->load (cuda::memory_order_relaxed);
4
5      if (it_state.data <= it) {
6          flag->wait (it_state, cuda::memory_order_relaxed);
7          it_state = flag->load (cuda::memory_order_relaxed);
8      }
9
10     if (it_state.flag)
11         block_last_it = it_state.data;
12 }
13 __syncthreads ();
14
15 last_it = block_last_it;
```



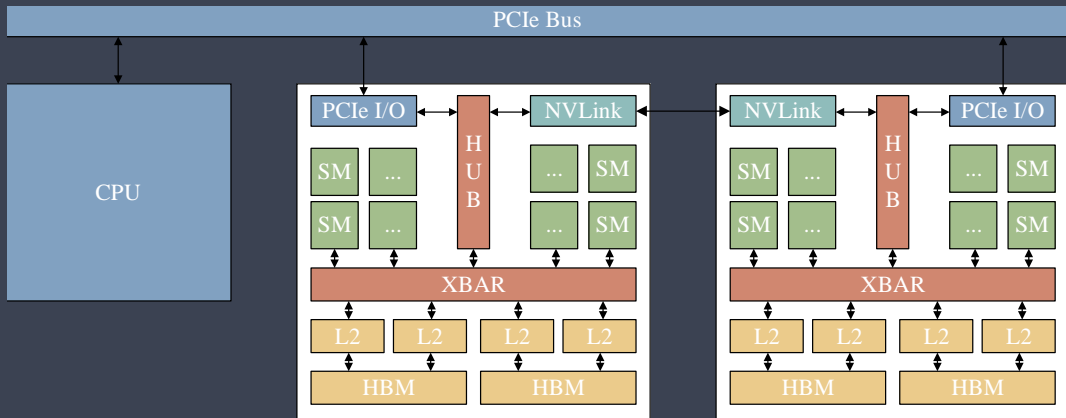






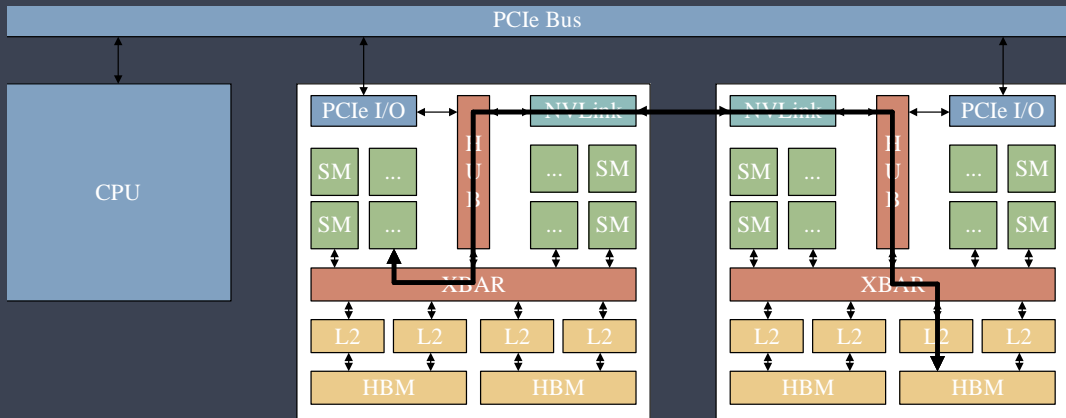
» Broadcast

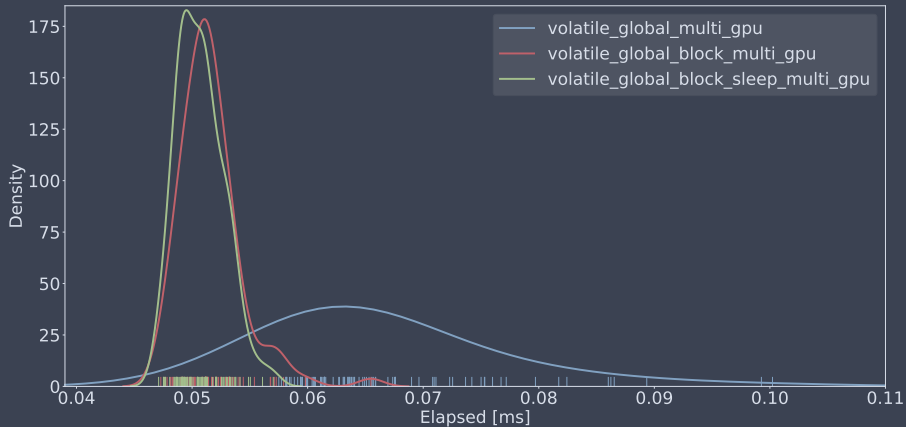
multi-GPU



» Broadcast


multi-GPU





Recommended materials

- * CppCon 2018: Olivier Giroux “High-Radix Concurrent C++.”
- * CppCon 2019: Olivier Giroux “The One-Decade Task: Putting `std::atomic` in CUDA.”
- * Cpp Toronto 2020: Bryce Adelstein Lelbach “The CUDA C++ Standard Library by.”
- * Sorin, Daniel J. and Hill, Mark D. and Wood, David A. “A Primer on Memory Consistency and Cache Coherence”

Source codes and presentation will be published soon, check @g_evtushenko