CS 758: Advanced Topics in Computer Architecture

Lecture #8: The GPU Memory System (Part 2)

Professor Matthew D. Sinclair

Backup of these slides were developed by Tim Rogers at the Purdue University.

Slides enhanced by Matt Sinclair

In-Class Activity

- With a partner, work on the 3 problems for 5 minutes
 - Problem 1: GPU MCM
 - Problem 2: GPU MCM with scopes
 - Problem 3: bugs with writing synchronization code on GPUs
- Then we'll come together as a group and discuss

CPU Coherence: MESI



- Write miss: Get ownership, invalidate all sharers
- Read miss: Update sharers list
- Synchronization points are cheap
- BUT poor fit for GPUs: [Singh13, Hechtman14]
 - Directory overhead, transient states, excessive traffic, indirection

Complex coherence, simple consistency

Atomics Background

- Default: Data-race-free-0 (DRF0) [ISCA '90]
 - Identify all races as synchronization accesses (C++: atomics)

// each thread for i = 0:n

ADD R4, A[i], R1synch (atomic)ADD R5, B[i], R1synch (atomic)

- All atomics örder data accesses
- Atomics order other atomics

...

 \Rightarrow Ensures SC semantics if no data races

Atomics Background (Cont.)

- Default: Data-race-free-0 (DRF0) [ISCA '90]
 - All atomics order data accesses
 - Atomics order other atomics
 ⇒Ensures SC semantics if no data races
- Data-race-free-1 (DRF1): unpaired atomics [TPDS '93]
 - + Unpaired atomics do not order data accesses
 - Atomics order other atomics
 - \Rightarrow Ensures SC semantics if no data races
- Relaxed atomics [PLDI '08]
 - + Do not order data or other atomics

 \Rightarrow But can violate SC and no formal specification

Traditional GPU Coherence



Each thread accesses independent data (no races) No data reuse or data sharing Coarse-grained synchronization Optimized for streaming, data parallel applications

GPU Memory Consistency Model

- Active area of research
- Tightly tied in with coherence protocol
- Provides very weak guarantees
 - Respect program order within a single thread
 - Easy to design hardware
 - Programmers add *fences* to provide extra guarantees
 - Fence guarantee all previous accesses are visible before proceeding
 - ... usually
- Most GPUs use a *scoped* memory consistency model
 - ____threadfence__block local synchronization (usually at L1)
 - _____threadfence GPU global synchronization (usually at L2)
 - _____threadfence___system CPU-GPU global synchronization (flush GPU)

GPU Coherence with DRF



- With data-race-free (DRF) memory model
 - No data races; synchs must be explicitly distinguished
 - Synchronization accesses (atomics) go to last level cache (LLC)
 - Synchronization points are expensive, preclude reuse

Simple but inefficient coherence, simple consistency

GPU Coherence with HRF



- New memory model: Heterogeneous-race-free (HRF) [ASPLOS '14]
 - Adds scoped synchronization

GPU Coherence with HRF



- New memory model: Heterogeneous-race-free (HRF) [ASPLOS '14]
 - Adds scoped synchronization
 - No overhead for locally scoped synchronizations
- But higher programming complexity

More efficient coherence, complex consistency

GPU Coherence with HRF

- heterogeneous HRF
 With data-race-free (DRF) memory model [ASPLOS '14]
 - No data races; synchs must be explicitly distinguished heterogeneous and their scopes
 - At all synch points

global

- Flush all dirty data: Unnecessary writethroughs
- Invalidate all data: Can't reuse data across synch points Global
 - Synchronization accesses must go to last level cache (LLC)
 - No overhead for locally scoped synchs
- But higher programming complexity

DeNovo Coherence with DRF



- Reuse dirty data across synch points more data reuse
- Synchronization accesses can be performed at L1 synch reuse

3% area overhead vs. GPU Coherence + HRF

Efficient coherence, simple consistency (except relaxed atomics)

Additional Topics

- Extending Coherence across accelerators
 - CCIX, ACE, Spandex, CXL,
 - Challenge: Different accelerators have different coherence requirements
 - E.g., different data widths
 - Especially important as number of accelerators increases
- Multi-GPU/Multi-Chiplet Coherence
 - Challenge: NUMA accesses (ala SMPs), page migration, partitioning
 - Potentially can extend existing mechanisms
- DRFrlx [Sinclair ISCA '17]
 - Extend MCM to provide sane semantics for relaxed atomics
- Timestamps
 - Avoid traditional coherence overheads, but some bookkeeping/delay instead

Additional Topics (Cont.)

- Coherence Granularity
 - HSC [Power MICRO '13] large granularity for traditional GPGPU apps
 - hUVM [Koukos TACO '16] page granularity coherence
 - Exploits traditional GPGPU apps nature
 - Some relationship with subsequent work on ACE, CCIX, Spandex, etc.
- Coherent Scratchpads
 - Best of both worlds!
- Remote Scopes [Orr ASPLOS '15]
 - Dynamically vary scope granularity to reduce overhead
- hLRC [Alsop MICRO '17]

Conclusion

- GPU coherence and consistency are hot, recent topics
 - Lots of ongoing research in the area
- Goal: avoid replicating 20-30 years of CPU MCM research
- Idea: evolve from that CPUs already have now
 - Design Goal: keep GPU coherence protocols simple
 - GPU apps don't need more complex coherence protocols
 - This has implications on MCM
- State-of-the-art (products): GPU + HRF
 - CPU-GPU coherence in real products assumes this for the GPU component

Backup

Coherence & Consistency Qualitative Analysis

Coherence + Consistency	Reu	se Data	Do Synchs		
ouncrence - ounsistency	Owned	Valid	at L1		
GPU + DRF(GD)	X	X	X		
GPU + HRF(GH)	local	local	local		
DeNovo + DRF (DD)	~	X	\checkmark		
DeNovo-RO + DRF (DD+RO)	\checkmark	read-only	\checkmark		
DeNovo + HRF (DH)	\checkmark	local	\checkmark		

g Other	Use Cases Into DRFrlx									
Work Queas	4. CKS									
Flags	Ref Counters									
How to incorporate SC violations in approximable applications? How to incorporate overlapped atomics that do not order? How to incorporate violations of SC that do not affect final result? How to incorporate relaxed atomics that do not order data?										
Category	Semantics									
Unpaired	SC									
Non-Ordering										
Commutative	Final result always SC									
Speculative										
Quantum	SC-centric: non-SC parts isolated									
	g Other Work Queues Flags Flags orate SC violati Corporate relaxe Category Unpaired Non-Ordering Commutative Speculative Quantum									

Review: Cache Coherence Problem



- Processors see different values for u after event 3
- With write back caches, value written back to memory depends on order of which cache writes back value first
- Unacceptable situation for programmers

Coherence Invariants

1. Single-Writer, Multiple-Reader (SWMR) Invariant



2. Data-Value Invariant. The value of the memory location at the start of an epoch is the same as the value of the memory location at the end of its last read-write epoch.

Coherence States

- How to design system satisfying invariants?
- Track "state" of memory block copies and ensure states changes satisfy invariants.
- Typical states: "modified", "shared", "invalid".
- Mechanism for updating block state called a coherence protocol.

Intra-GPU Coherence

•

•



GPU Coherence Challenges

• Challenge 1: Coherence traffic



GPU Coherence Challenges

- Challenge 2: Tracking in-flight requests
 - Significant % of L2



GPU Coherence Challenges

Challenge 3: Complexity

Non-coherent L1



Non-coherent L2

	L1 GETS	WB Data	L2 Atomic	<u>L2</u> Replacement	L2 Replacement clean	Mem Data
<u>NP</u>	q <u>lpB</u> i s a j / <u>ISS</u>	q <u>lpB d i x as</u> j / <u>IM</u>	q <u>lpB</u> <u>d</u> i x as j / <u>IMA</u>			
<u>88</u>	<u>lpR ds set j</u>	<u>f lpW d de mr</u> <u>set j</u>	<u>f lpW d ds a mr</u> <u>set j</u>	<u>f lpE c r /NP</u>	<u>f lpE r /NP</u>	
<u>ISS</u>	<u>lpR s j</u>	<u>Z</u>	<u>Z</u>	<u>Z</u>	Z	<u>m e s o /SS</u>
IM	Z	Z	Z	<u>Z</u>	Z	<u>m mt ee s o</u> / <u>SS</u>
<u>IMA</u>	z	Z	Z	Z	Z	<u>m mt e a s</u> o / <u>SS</u>

MESI L2 States

		L1 GET INSTR	L1 GETS	L1 GETX	L1 UPGRADE	L1 PUTX	L1 PUTX old	<u>L2</u> Replacement	<u>L2</u> Replacement <u>clean</u>	<u>Mem</u> Data	<u>WB</u> Data	<u>WB</u> Data clean	<u>Ack</u>	<u>Ack</u> <u>all</u>	<u>Unblock</u>	<u>Unblock</u> <u>Cancel</u>	Exclusive Unblock
N	<u>NP</u>	qln is auj/ <u>IS</u>	qln isa uj/ <u>ISS</u>	qlixaQu j/⊡		t j	t j										
	<u>SS</u>	<u>ds n u</u> set j	<u>ds n u set</u> j	d <u>fwm</u> u set j / <u>SS</u> MB	<u>fwm ts u set</u> j / <u>SS MB</u>	ţj	ţj	if r/<u>SI</u>	if r/∐								
N	M	d <u>n u set</u> j / <u>SS</u>	<u>dd u set</u> j / <u>MT MB</u>	<u>d u set j</u> / <u>MT MB</u>		t j	t j	i <u>c r s /NP</u>	i <u>r s /NP</u>								
1	MT	<u>b u set j</u> / <u>MT IIB</u>	<u>b u set</u> j / <u>MT IIB</u>	<u>b u set j</u> / <u>MT MB</u>		lmrtj / <u>M</u>	<u>t j</u>	ifr/ <u>MTI</u>	i <u>f r /MCT I</u>								
5	MT I	<u>zz</u>	<u>zz</u>	<u>ZZ</u>	<u>ZZ</u>	<u>77</u>	<u>77</u>				q q ct s o / <u>NP</u>	<u>s o</u> / <u>NP</u>		<u>ct s</u> 0 / <u>NP</u>			
E	MCT I	<u>ZZ</u>	<u>zz</u>	<u>zz</u>	<u>ZZ</u>	ZZ	<u>ZZ</u>				q q <u>ct s</u> <u>o</u> / <u>NP</u>	<u>s o</u> / <u>NP</u>		<u>s o</u> / <u>NP</u>			
N	Ш	<u>ZZ</u>	zz	<u>ZZ</u>	<u>ZZ</u>	t j	t j						<u>q o</u>	<u>s o</u> / <u>NP</u>			
Ľ	<u>s i</u>	ZZ	ZZ	zz	ZZ	t j	ţį						d o	<u>ct s</u> o /NP			
I	<u>188</u>	<u>n s u j</u> / <u>IS</u>	nsuj / <u>IS</u>	<u>zz</u>		t j	t j	<u>zz</u>	<u>zz</u>	<u>md ex s</u> od / <u>MT</u> MB							
<u>S</u>	<u>15</u>	<u>n s u j</u>	nsuj	ZZ		t j	t j	ZZ	ZZ	<u>md e s od</u> / <u>SS</u>							
IS	IM	<u>zz</u>	zz	ZZ		t j	ţj	ZZ	<u>ZZ</u>	<u>md ee s</u> od / <u>MT</u> <u>MB</u>							
M	SS MB	<u>zz</u>	<u>zz</u>	<u>zz</u>	<u>zz</u>	zz	<u>zz</u>	<u>ZZ</u>	<u>77</u>							<u>k</u> / <u>SS</u>	<u>mu k</u> / <u>MT</u>
<u>E</u>	<u>MT</u> <u>MB</u>	<u>77</u>	zz	<u>zz</u>	<u>ZZ</u>	<u>zz</u>	<u>ZZ</u>	<u>ZZ</u>	<u>ZZ</u>							<u>k</u> / <u>MT</u>	<u>mu k</u> / <u>MT</u>
	<u>М</u> <u>МВ</u>	<u>77</u>	zz	<u>zz</u>	<u>ZZ</u>	t j	t j	<u>ZZ</u>	<u>ZZ</u>								<u>mu k</u> / <u>MT</u>
	MT IIB	<u>zz</u>	zz	<u>ZZ</u>	<u>ZZ</u>	<u>ZZ</u>	<u>zz</u>	<u>ZZ</u>	<u>ZZ</u>		<u>m o</u> / <u>MT</u> <u>SB</u>	<u>m o</u> / <u>MT</u> <u>SB</u>			<u>nu k</u> / <u>MT IB</u>		
	<u>MT</u> <u>IB</u>	ZZ	ZZ	ZZ	ZZ	ţj	tj	<u>ZZ</u>	<u>77</u>		<u>m o</u> / <u>SS</u>	<u>m o /SS</u>				<u>k</u> / <u>MT</u>	
	<u>MT</u> <u>SB</u>	<u>ZZ</u>	ZZ	<u>zz</u>	<u>ZZ</u>	t j	t j	<u>ZZ</u>	<u>ZZ</u>						<u>nu k</u> / <u>SS</u>		

Coherence Challenges

- Challenges of introducing coherence messages on a GPU
 - 1. Traffic: transferring messages
 - 2. Storage: tracking message
 - 3. Complexity: managing races between messages
- GPU cache coherence without coherence messages?
 - YES using global time

Temporal Coherence

Related: Library Cache Coherence



Temporal Coherence Example





Lifetime Predictor

- One prediction value per L2 bank
- Events local to L2 bank update prediction value







- Reduces traffic by 53% over MESI and 23% over GPU-VI for intra-workgroup applications
- Lower traffic than 16x-sized 32-way directory

Performance





 TC-Weak with simple predictor performs 85% better than disabling L1 caches

CPU-GPU Coherence?

- Many vendors have introduced chips with both CPU and GPU (e.g., AMD Fusion, Intel Core i7, NVIDIA Tegra, etc...)
- What are the challenges with maintaining coherence across CPU and GPU?
- One important one: GPU has higher cache miss rate than CPU. Can place pressure on directory impacting performance.
- Power et al., Heterogeneous System Coherence for Integrated CPU-GPU Systems, MICRO 2013: Use "region coherence" to reduce number of GPU requests that need to access directory.

Synchronization

- Locks are not encouraged in current GPGPU programming manuals.
- Interaction with SIMT stack can easily cause deadlocks:

```
while( atomicCAS(&lock[a[tid]],0,1) != 0 )
; // deadLock here if a[i] = a[j] for any i,j = tid in
warp
```

// critical section goes here

```
atomicExch (&lock[a[tid]], 0);
```

Correct way to write critical section for GPGPU:

```
done = false;
while( !done ) {
    if( atomicCAS (&lock[a[tid]], 0 , 1 )==0 ) {
        // critical section goes here
        atomicExch(&lock[a[tid]], 0 ) ;
    }
}
```

Most current GPGPU programs use barriers within thread blocks and/or lock-free data structures.

This leads to the following picture...

• Lifetime of GPU Application Development

