

Architecting and Exploiting Asymmetry in Multi-Core Architectures

Onur Mutlu

onur@cmu.edu

July 23, 2013

BSC/UPC

SAFARI Carnegie Mellon

Overview of My Group's Research

- Heterogeneous systems, accelerating bottlenecks
- Memory (and storage) systems
 - Scalability, energy, latency, parallelism, performance
 - Compute in/near memory
- Predictable performance, QoS
- Efficient interconnects
- Bioinformatics algorithms and architectures
- Acceleration of important applications, software/hardware co-design

Three Key Problems in Future Systems

■ Memory system

- ❑ Many important existing and future applications are increasingly data intensive → require bandwidth and capacity
- ❑ Data storage and movement limits performance & efficiency

■ Efficiency (performance and energy) → scalability

- ❑ Enables scalable systems → new applications
- ❑ Enables better user experience → new usage models

■ Predictability and robustness

- ❑ Resource sharing and unreliable hardware causes QoS issues
- ❑ Predictable performance and QoS are first class constraints

Readings and Videos

Mini Course: Multi-Core Architectures

- Lecture 1.1: Multi-Core System Design
 - <http://users.ece.cmu.edu/~omutlu/pub/onur-Bogazici-June-6-2013-lecture1-1-multicore-and-asymmetry-afterlecture.pptx>
- Lecture 1.2: Cache Design and Management
 - <http://users.ece.cmu.edu/~omutlu/pub/onur-Bogazici-June-7-2013-lecture1-2-cache-management-afterlecture.pptx>
- Lecture 1.3: Interconnect Design and Management
 - <http://users.ece.cmu.edu/~omutlu/pub/onur-Bogazici-June-10-2013-lecture1-3-interconnects-afterlecture.pptx>

Mini Course: Memory Systems

- Lecture 2.1: DRAM Basics and DRAM Scaling
 - <http://users.ece.cmu.edu/~omutlu/pub/onur-Bogazici-June-13-2013-lecture2-1-dram-basics-and-scaling-afterlecture.pptx>
- Lecture 2.2: Emerging Technologies and Hybrid Memories
 - <http://users.ece.cmu.edu/~omutlu/pub/onur-Bogazici-June-14-2013-lecture2-2-emerging-memory-afterlecture.pptx>
- Lecture 2.3: Memory QoS and Predictable Performance
 - <http://users.ece.cmu.edu/~omutlu/pub/onur-Bogazici-June-17-2013-lecture2-3-memory-qos-afterlecture.pptx>

Readings for Today

- Required – Symmetric and Asymmetric Multi-Core Systems
 - Suleman et al., “[Accelerating Critical Section Execution with Asymmetric Multi-Core Architectures](#),” ASPLOS 2009, IEEE Micro 2010.
 - Suleman et al., “[Data Marshaling for Multi-Core Architectures](#),” ISCA 2010, IEEE Micro 2011.
 - Joao et al., “[Bottleneck Identification and Scheduling for Multithreaded Applications](#),” ASPLOS 2012.
 - Joao et al., “[Utility-Based Acceleration of Multithreaded Applications on Asymmetric CMPs](#),” ISCA 2013.
- Recommended
 - Amdahl, “[Validity of the single processor approach to achieving large scale computing capabilities](#),” AFIPS 1967.
 - Olukotun et al., “[The Case for a Single-Chip Multiprocessor](#),” ASPLOS 1996.
 - Mutlu et al., “[Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-order Processors](#),” HPCA 2003, IEEE Micro 2003.
 - Mutlu et al., “[Techniques for Efficient Processing in Runahead Execution Engines](#),” ISCA 2005, IEEE Micro 2006.

Videos for Today

■ Multiprocessors

□ Basics:

http://www.youtube.com/watch?v=7ozCK_Mgxfk&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=31

□ Correctness and Coherence:

<http://www.youtube.com/watch?v=U-VZKMgItDM&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=32>

□ Heterogeneous Multi-Core:

<http://www.youtube.com/watch?v=r6r2NJxj3kI&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=34>

■ Runahead Execution

□ [http://www.youtube.com/watch?](http://www.youtube.com/watch?v=z8YpjqXQJIA&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=28)

[v=z8YpjqXQJIA&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=28](http://www.youtube.com/watch?v=z8YpjqXQJIA&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=28)

Online Lectures and More Information

■ Online Computer Architecture Lectures

- <http://www.youtube.com/playlist?list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ>

■ Online Computer Architecture Courses

- Intro: <http://www.ece.cmu.edu/~ece447/s13/doku.php>
- Advanced: <http://www.ece.cmu.edu/~ece740/f11/doku.php>
- Advanced: <http://www.ece.cmu.edu/~ece742/doku.php>

■ Recent Research Papers

- <http://users.ece.cmu.edu/~omutlu/projects.htm>
- <http://scholar.google.com/citations?user=7XyGUGkAAAAJ&hl=en>

Architecting and Exploiting Asymmetry in Multi-Core Architectures

Warning

- This is an asymmetric talk
- But, we do not need to cover all of it...
- Component 1: A case for asymmetry *everywhere*
- Component 2: A deep dive into mechanisms to exploit asymmetry in processing cores
- Component 3: Asymmetry in memory controllers
- Asymmetry = heterogeneity
 - A way to enable specialization/customization

The Setting

- Hardware resources are shared among many threads/apps in a many-core system
 - Cores, caches, interconnects, memory, disks, power, lifetime, ...
- Management of these resources is a very difficult task
 - When optimizing parallel/multiprogrammed workloads
 - Threads interact unpredictably/unfairly in shared resources
- Power/energy consumption is arguably the most valuable shared resource
 - Main limiter to efficiency and performance

Shield the Programmer from Shared Resources

- Writing even sequential software is hard enough
 - Optimizing code for a complex shared-resource parallel system will be a nightmare for most programmers
- Programmer should not worry about (hardware) resource management
 - What should be executed where with what resources
- Future computer architectures should be designed to
 - Minimize programmer effort to optimize (parallel) programs
 - Maximize runtime system's effectiveness in automatic shared resource management

Shared Resource Management: Goals

- Future many-core systems should manage power and performance automatically across threads/applications
- Minimize energy/power consumption
- While satisfying performance/SLA requirements
 - Provide predictability and Quality of Service
- Minimize programmer effort
 - In creating optimized parallel programs
- Asymmetry and configurability in system resources essential to achieve these goals

Asymmetry Enables Customization

c	c	c	c
c	c	c	c
c	c	c	c
c	c	c	c

Symmetric

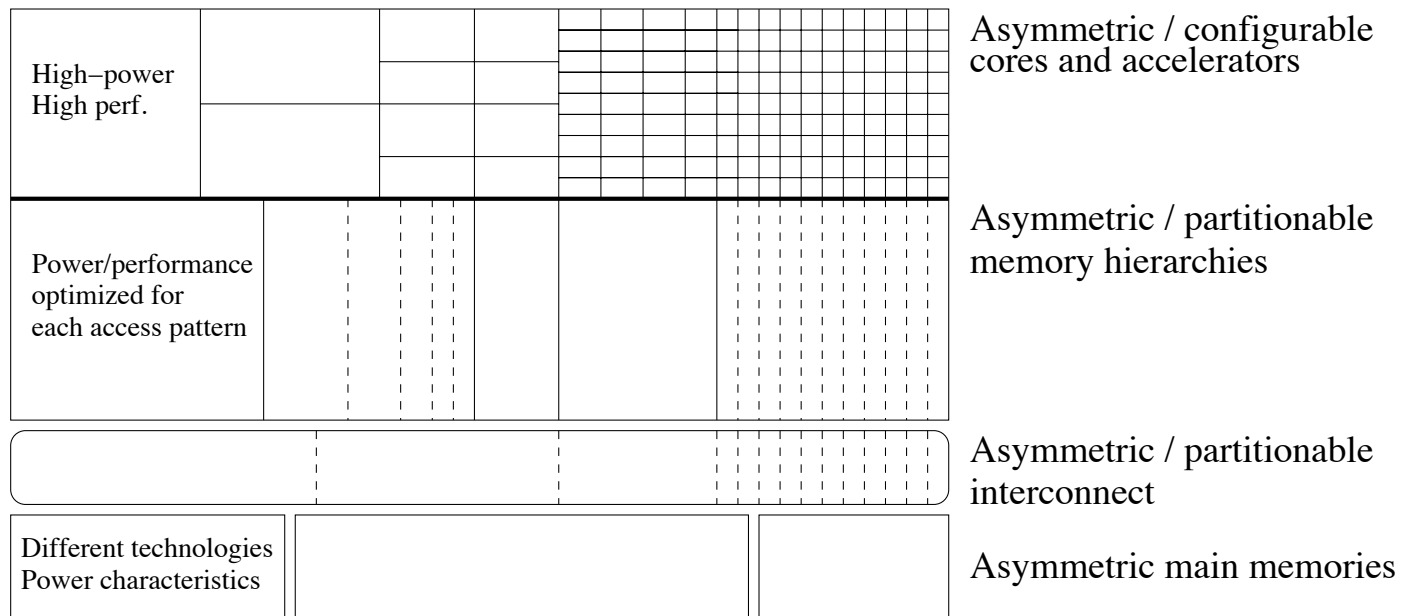
C1		C2	
		C3	
C4	C4	C4	C4
C5	C5	C5	C5

Asymmetric

- Symmetric: One size fits all
 - Energy and performance suboptimal for different phase behaviors
- Asymmetric: Enables tradeoffs and customization
 - Processing requirements vary across applications and phases
 - Execute code on best-fit resources (minimal energy, adequate perf.)

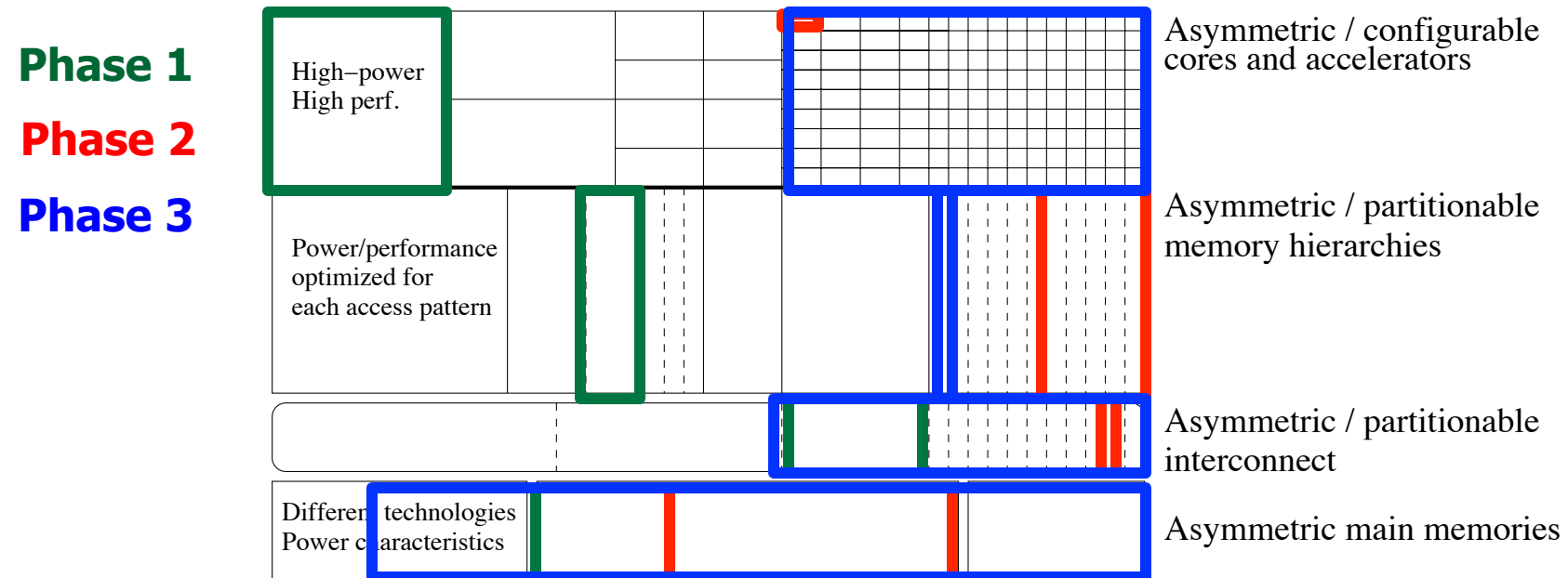
Thought Experiment: Asymmetry Everywhere

- Design each hardware resource with **asymmetric, (re-)configurable, partitionable components**
 - ❑ Different power/performance/reliability characteristics
 - ❑ To fit different computation/access/communication patterns



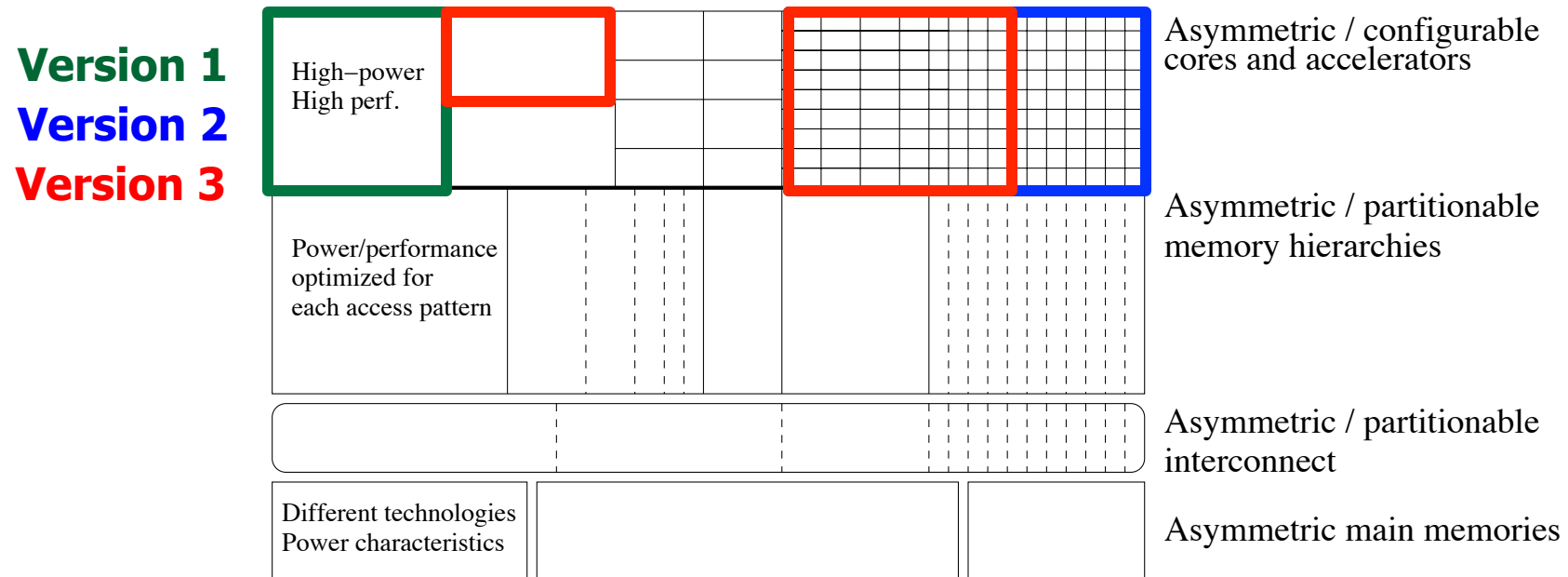
Thought Experiment: Asymmetry Everywhere

- Design the **runtime system (HW & SW)** to **automatically choose** the best-fit components for each phase
 - Satisfy performance/SLA with minimal energy
 - Dynamically stitch together the “best-fit” chip for each phase



Thought Experiment: Asymmetry Everywhere

- **Morph software components** to match asymmetric HW components
 - Multiple versions for different resource characteristics



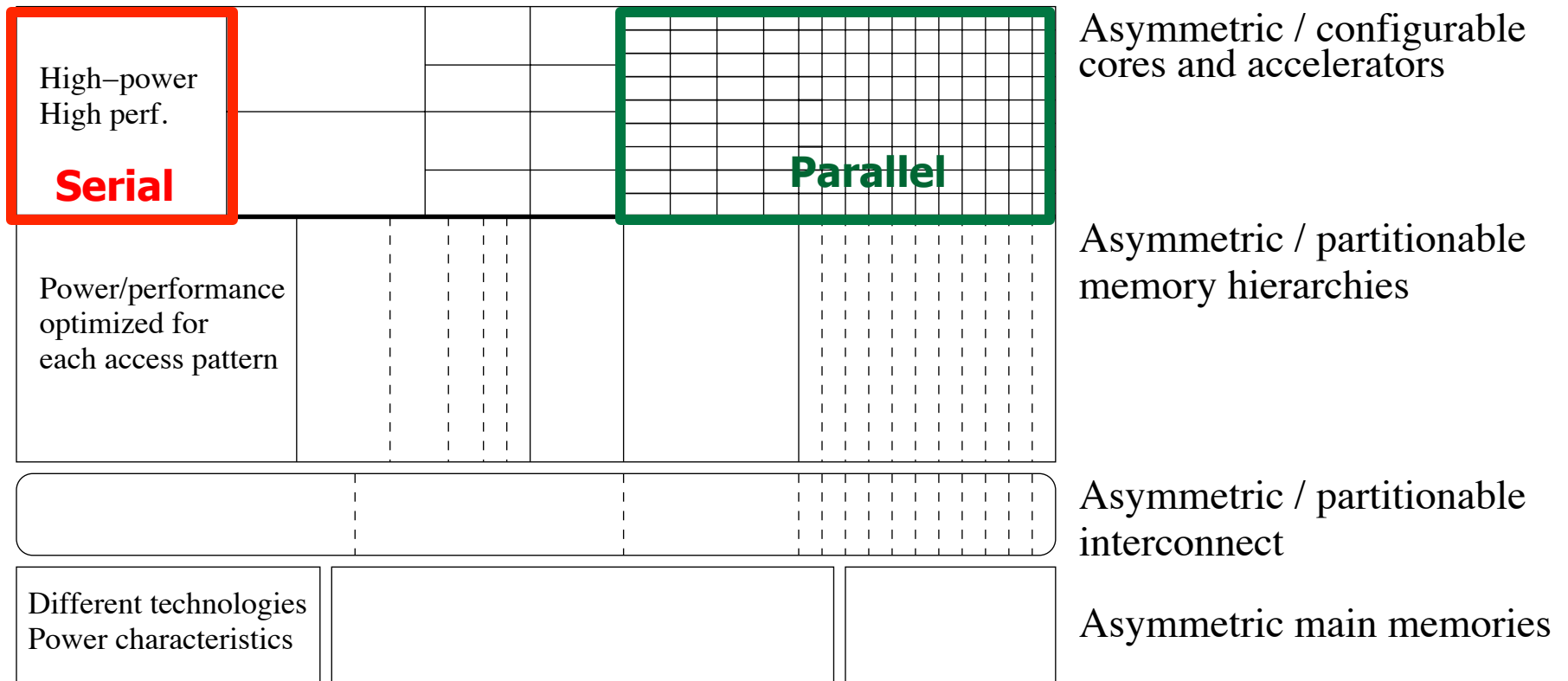
Many Research and Design Questions

- How to design asymmetric components?
 - Fixed, partitionable, reconfigurable components?
 - What types of asymmetry? Access patterns, technologies?
- What monitoring to perform cooperatively in HW/SW?
 - Automatically discover phase/task requirements
- How to design feedback/control loop between components and runtime system software?
- How to design the runtime to automatically manage resources?
 - Track task behavior, pick “best-fit” components for the entire workload

Talk Outline

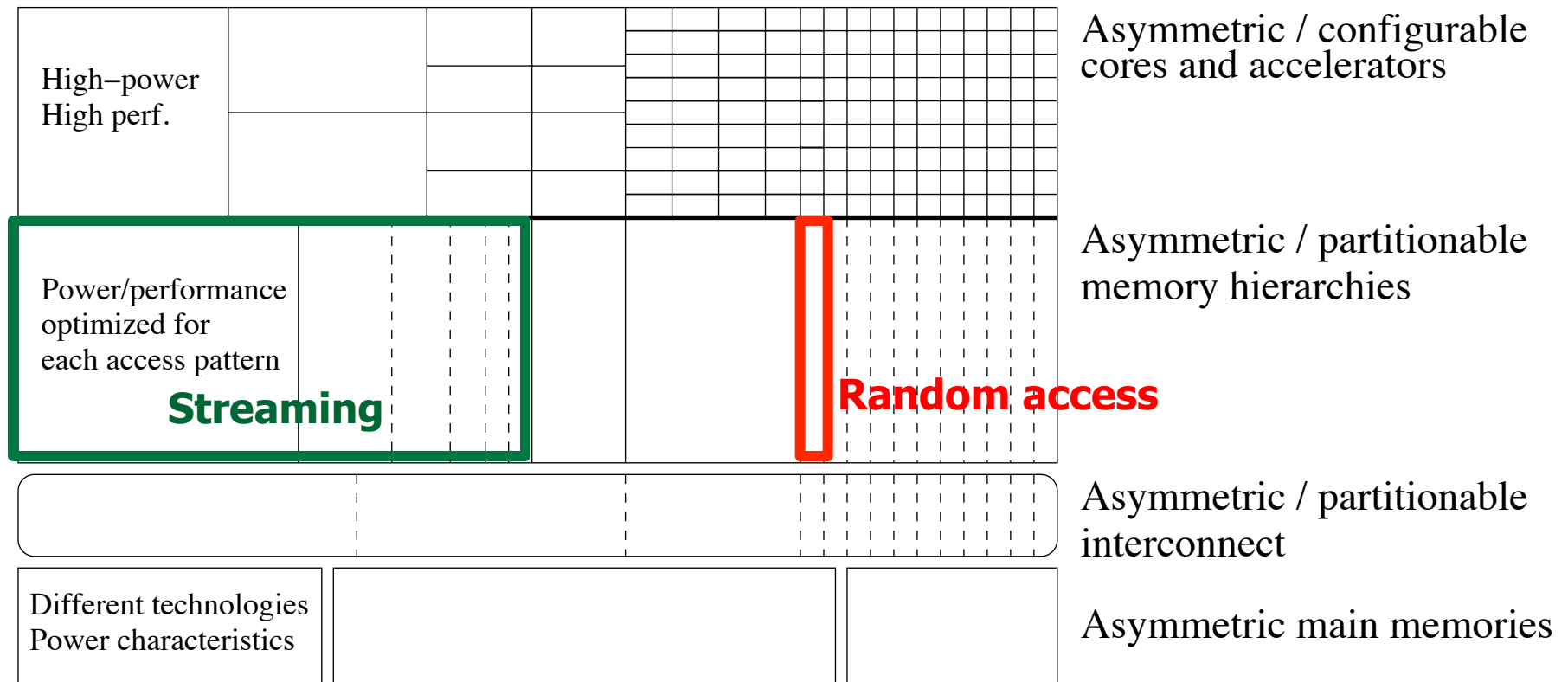
- Problem and Motivation
- How Do We Get There: Examples
- Accelerated Critical Sections (ACS)
- Bottleneck Identification and Scheduling (BIS)
- Staged Execution and Data Marshaling
- Thread Cluster Memory Scheduling (if time permits)
- Ongoing/Future Work
- Conclusions

Exploiting Asymmetry: Simple Examples



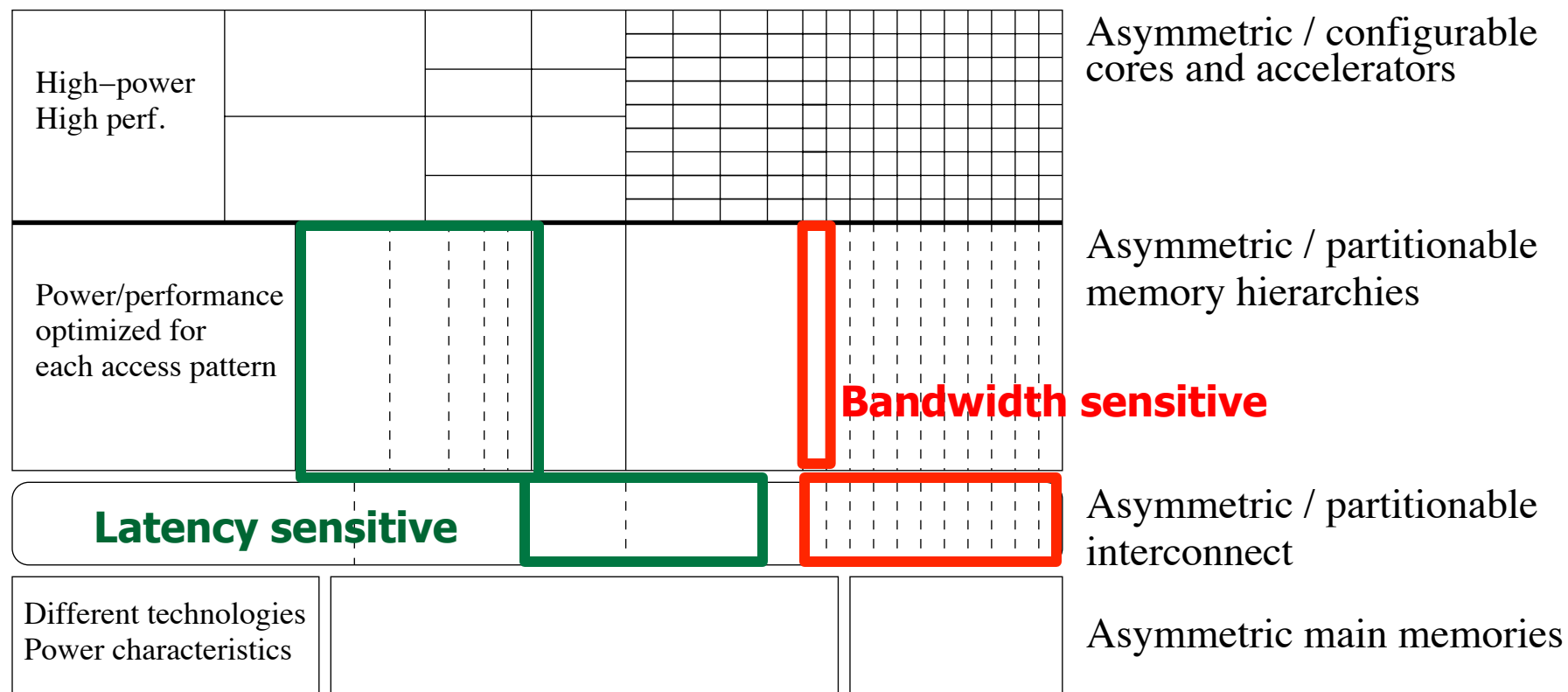
- Execute critical/serial sections on high-power, high-performance cores/resources [Suleman+ ASPLOS'09, ISCA'10, Top Picks'10'11, Joao+ ASPLOS'12]
 - Programmer can write less optimized, but more likely correct programs

Exploiting Asymmetry: Simple Examples



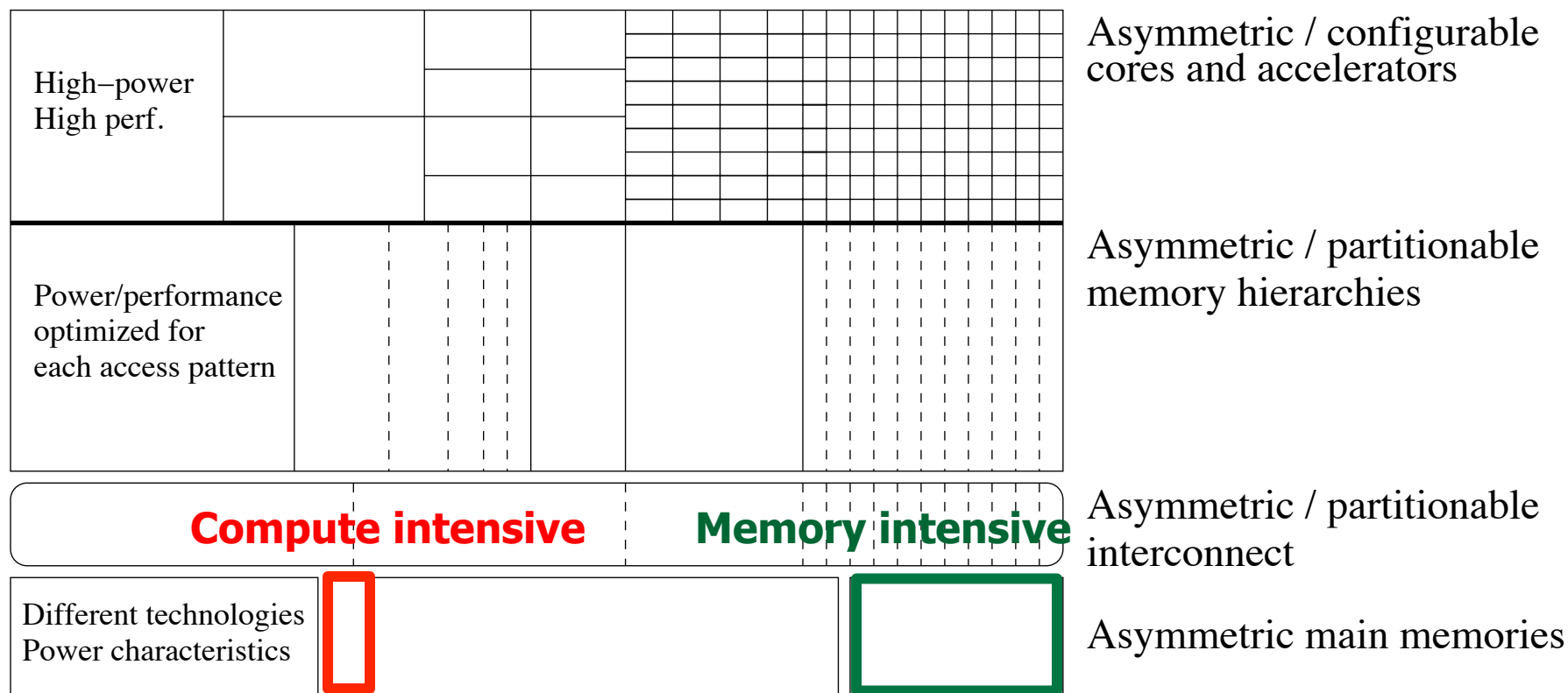
- Execute streaming “memory phases” on streaming-optimized cores and memory hierarchies
 - More efficient and higher performance than general purpose hierarchy

Exploiting Asymmetry: Simple Examples



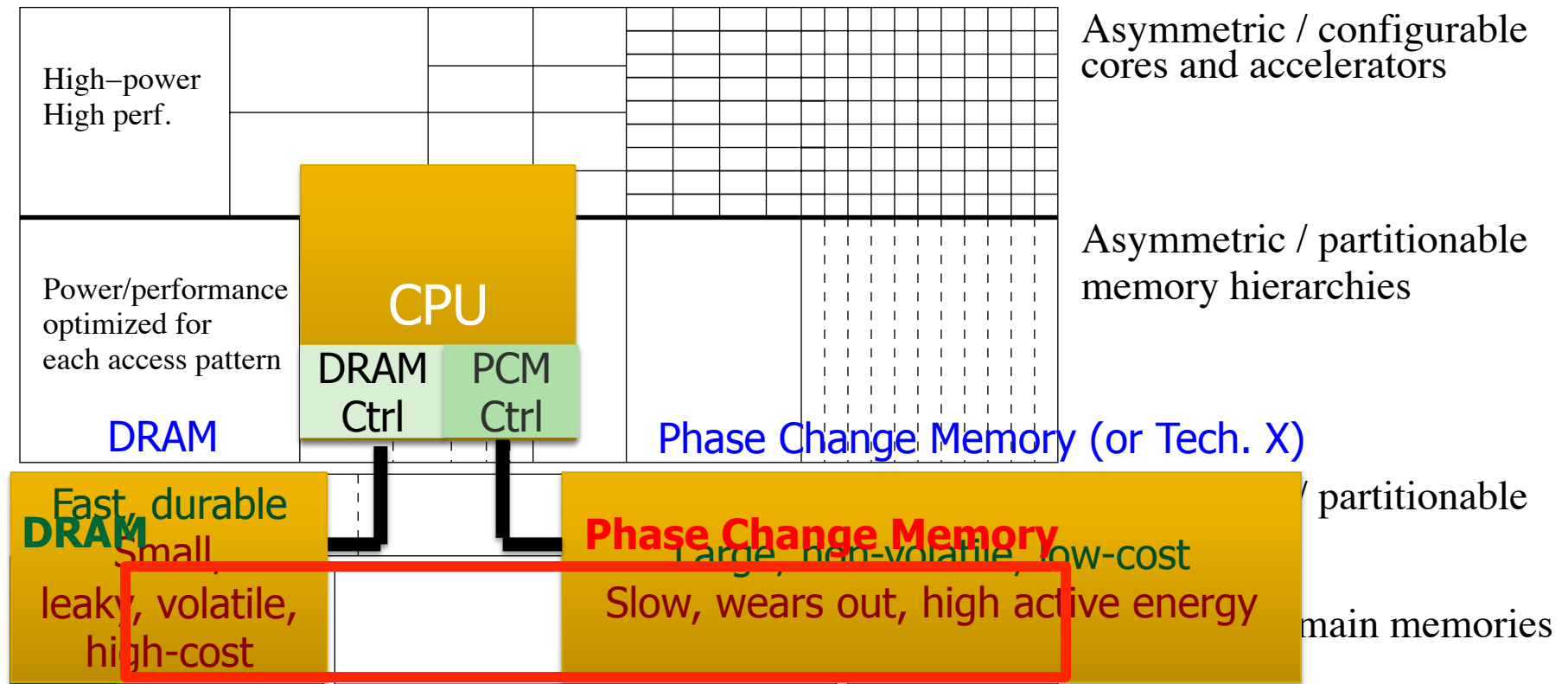
- Partition memory controller and on-chip network bandwidth asymmetrically among threads [Kim+ HPCA 2010, MICRO 2010, Top Picks 2011] [Nychis+ HotNets 2010] [Das+ MICRO 2009, ISCA 2010, Top Picks 2011]
 - Higher performance and energy-efficiency than symmetric/free-for-all

Exploiting Asymmetry: Simple Examples



- Have multiple different memory scheduling policies apply them to different sets of threads based on thread behavior [Kim+ MICRO 2010, Top Picks 2011] [Ausavarungnirun, ISCA 2012]
 - Higher performance and fairness than a homogeneous policy

Exploiting Asymmetry: Simple Examples



- Build main memory with different technologies with different characteristics (energy, latency, wear, bandwidth) [Meza+ IEEE CAL'12]
 - Map pages/applications to the best-fit memory resource
 - Higher performance and energy-efficiency than single-level memory

Talk Outline

- Problem and Motivation
- How Do We Get There: Examples
- Accelerated Critical Sections (ACS)
- Bottleneck Identification and Scheduling (BIS)
- Staged Execution and Data Marshaling
- Thread Cluster Memory Scheduling (if time permits)
- Ongoing/Future Work
- Conclusions

Serialized Code Sections in Parallel Applications

- Multithreaded applications:
 - Programs split into threads
- Threads execute concurrently on multiple cores
- Many parallel programs cannot be parallelized completely
- Serialized code sections:
 - Reduce performance
 - Limit scalability
 - Waste energy

Causes of Serialized Code Sections

- Sequential portions (Amdahl's "serial part")
- Critical sections
- Barriers
- Limiter stages in pipelined programs

Bottlenecks in Multithreaded Applications

Definition: any code segment for which threads contend (i.e. wait)

Examples:

- Amdahl's serial portions

- Only one thread exists → on the critical path

- Critical sections

- Ensure mutual exclusion → likely to be on the critical path if contended

- Barriers

- Ensure all threads reach a point before continuing → the latest thread arriving is on the critical path

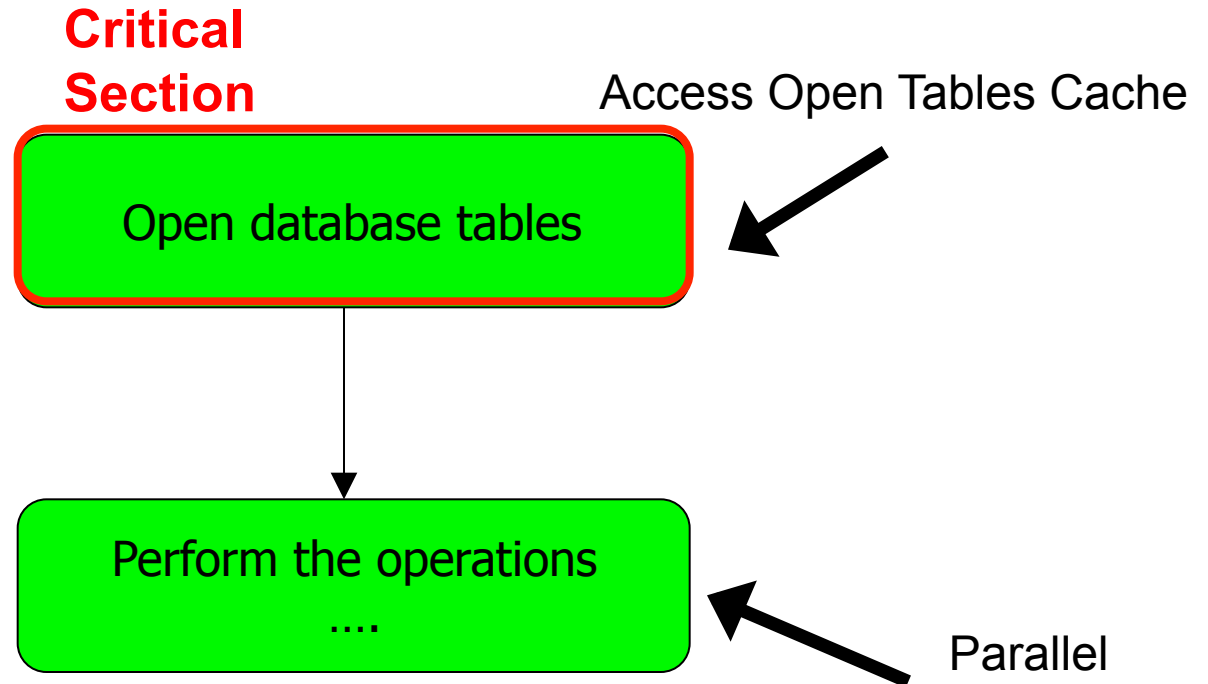
- Pipeline stages

- Different stages of a loop iteration may execute on different threads, slowest stage makes other stages wait → on the critical path

Critical Sections

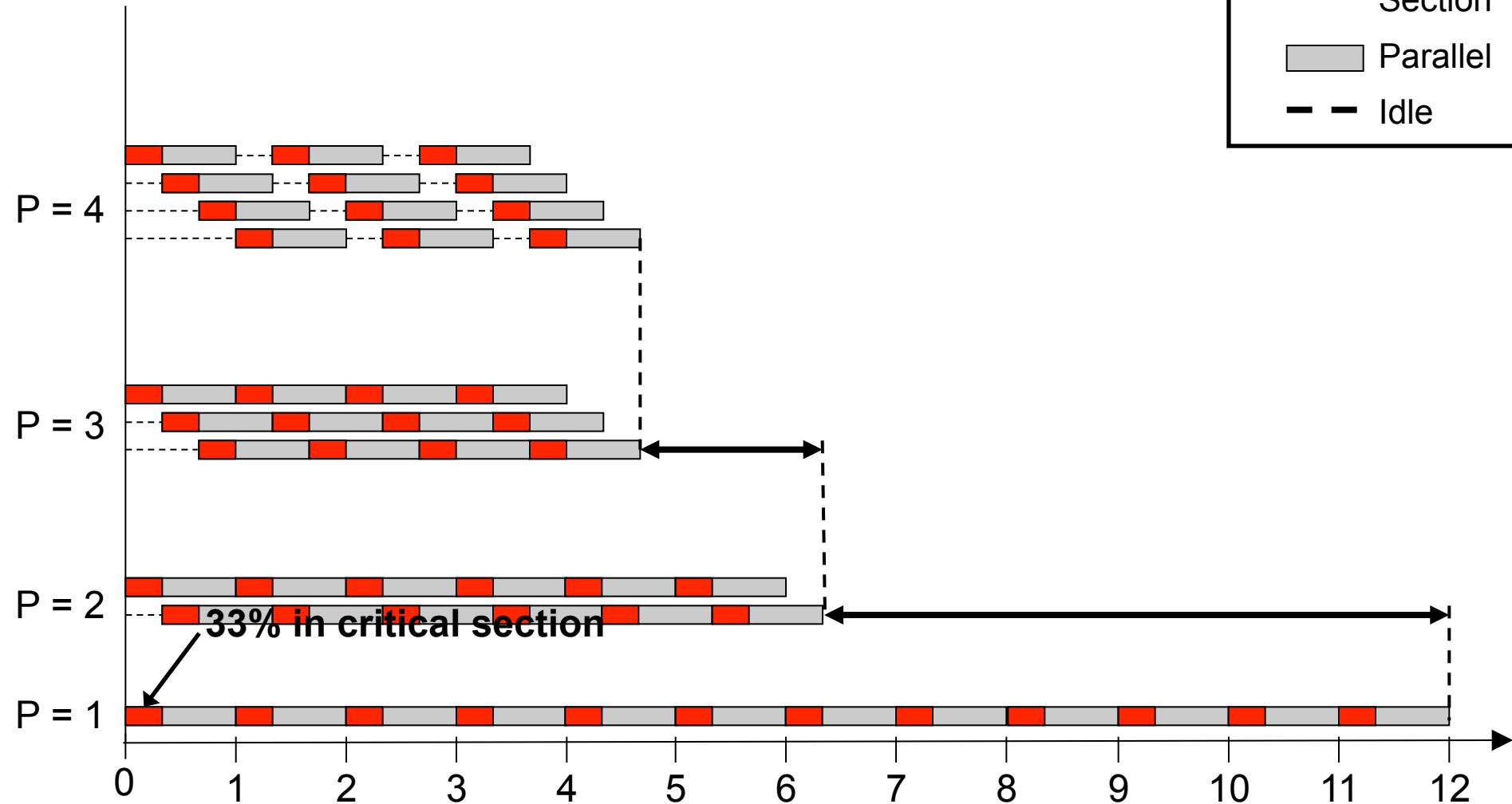
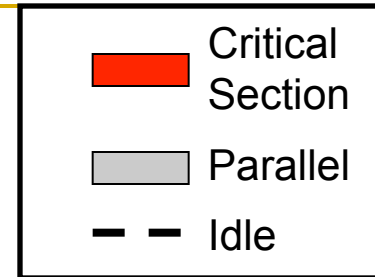
- Threads are not allowed to update shared data concurrently
 - For correctness (mutual exclusion principle)
- Accesses to shared data are encapsulated inside ***critical sections***
- Only one thread can execute a critical section at a given time

Example from MySQL



Contention for Critical Sections

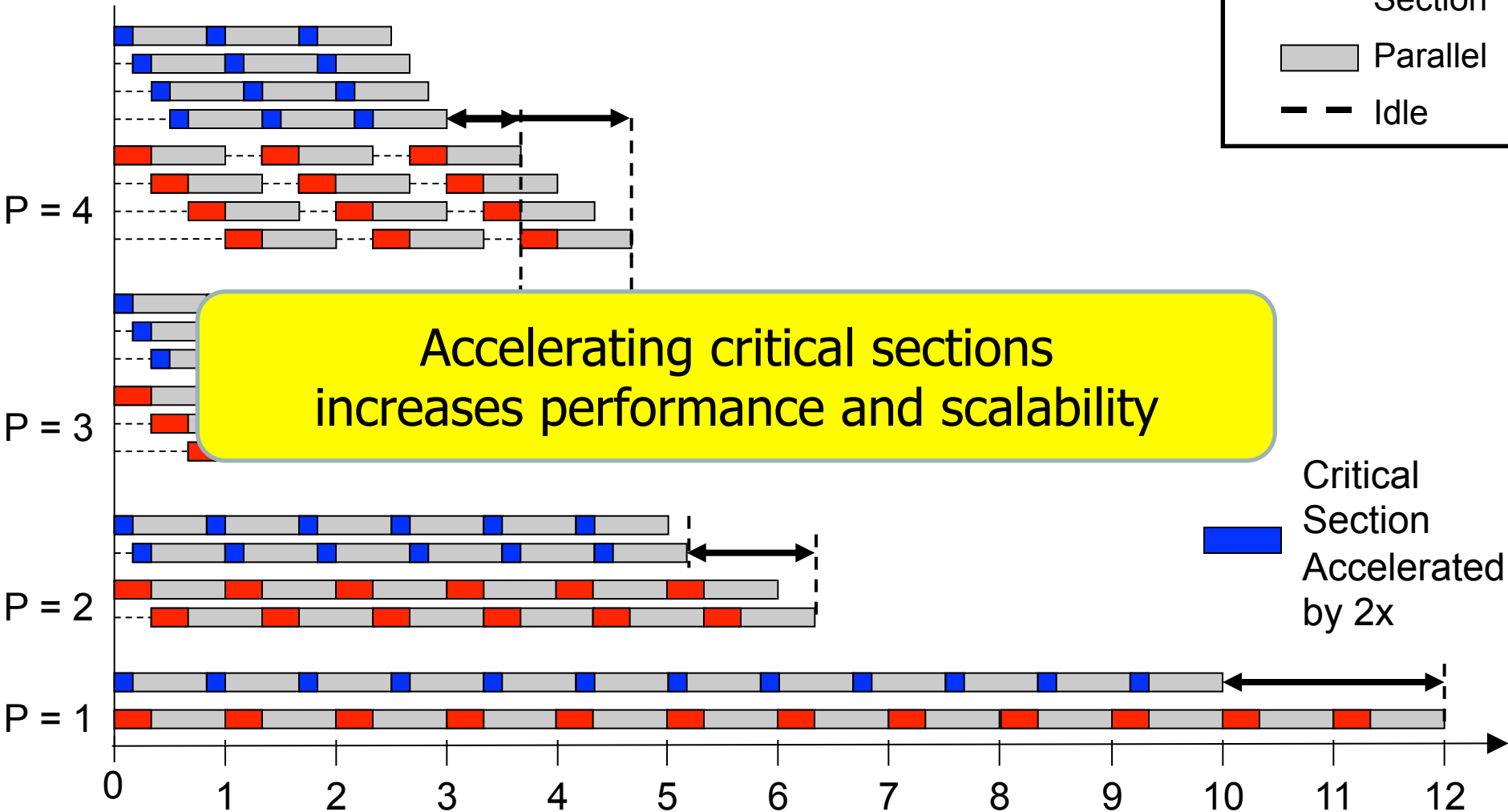
12 iterations, 33% instructions inside the critical section



Contention for Critical Sections

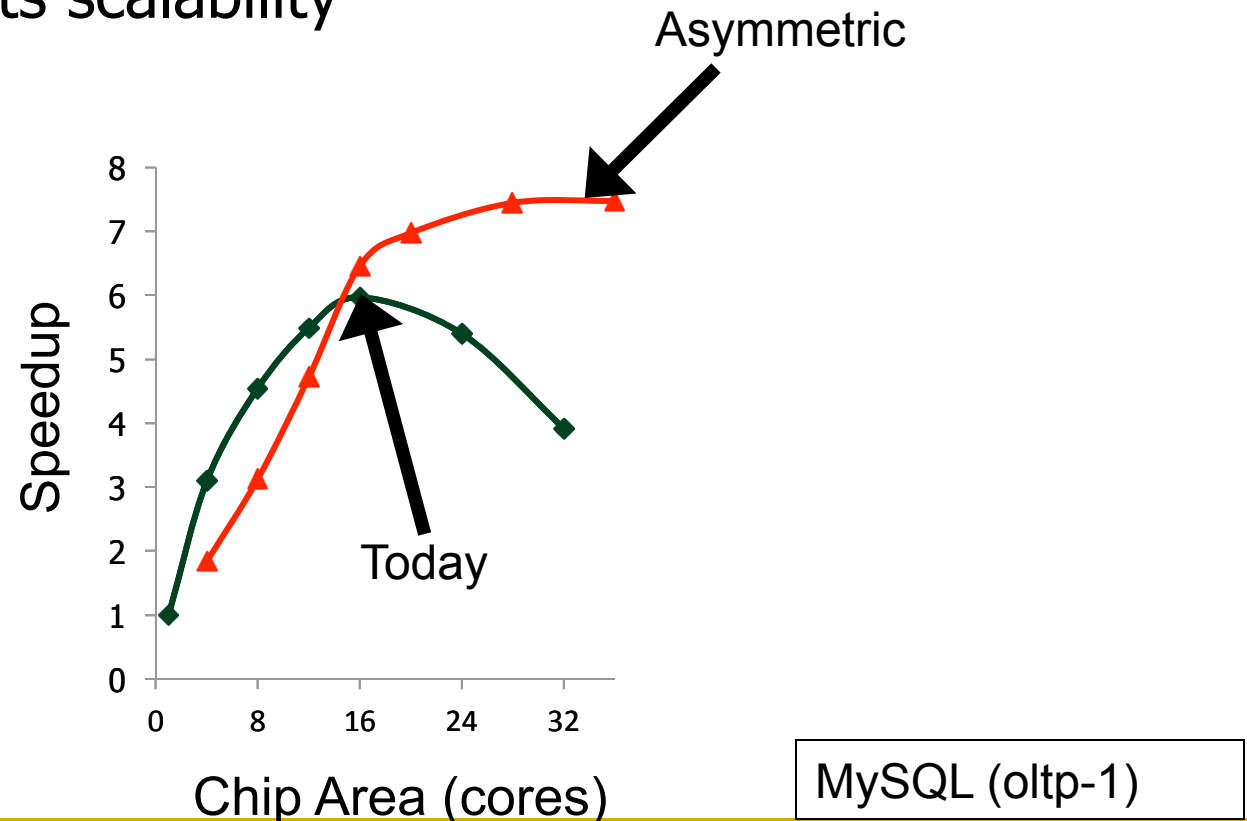
12 iterations, 33% instructions inside the critical section

Critical Section
Parallel
Idle



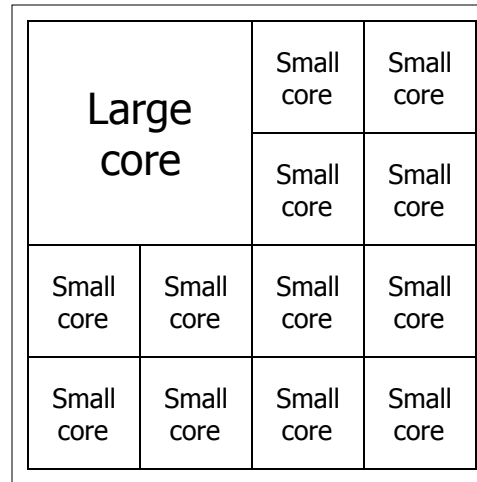
Impact of Critical Sections on Scalability

- Contention for critical sections leads to serial execution (serialization) of threads in the parallel program portion
- Contention for critical sections increases with the number of threads and limits scalability



A Case for Asymmetry

- Execution time of sequential kernels, critical sections, and limiter stages must be short
- It is difficult for the programmer to shorten these serialized sections
 - Insufficient domain-specific knowledge
 - Variation in hardware platforms
 - Limited resources
- Goal: A mechanism to shorten serial bottlenecks without requiring programmer effort
- Idea: Accelerate serialized code sections by shipping them to powerful cores in an asymmetric multi-core (ACMP)



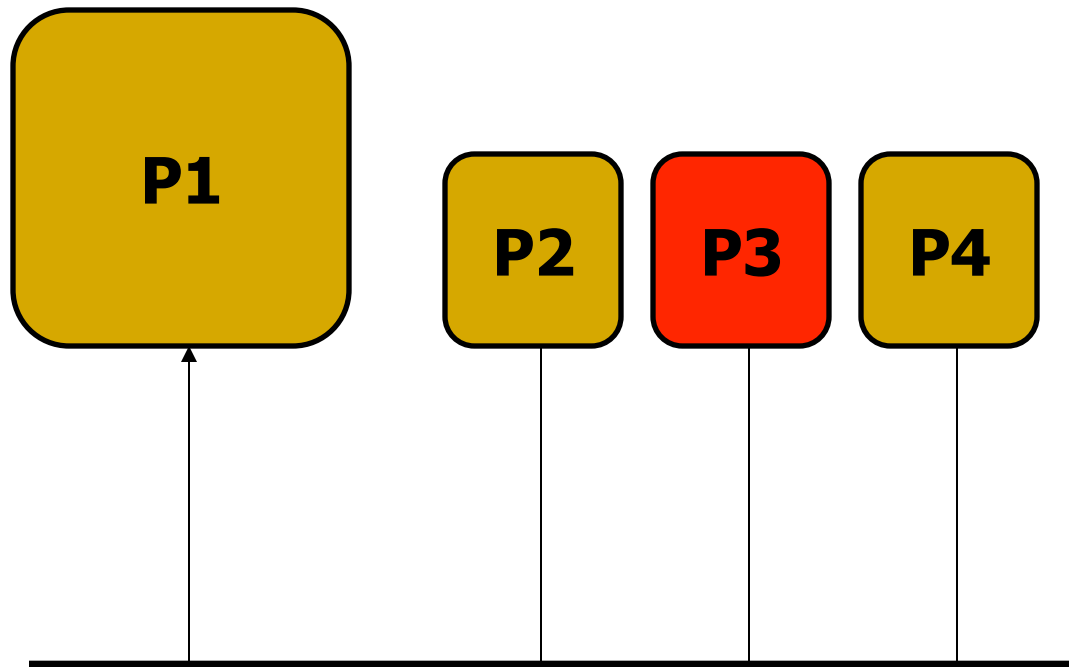
ACMP

- Provide one large core and many small cores
- Execute parallel part on small cores for high throughput
- Accelerate serialized sections using the large core
 - Baseline: Amdahl's serial part accelerated [Morad+ CAL 2006, Suleman+, UT-TR 2007]

Conventional ACMP

```
EnterCS()  
    PriorityQ.insert(...)  
LeaveCS()
```

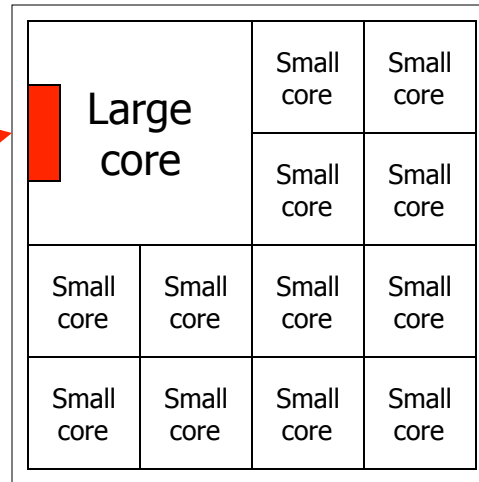
1. P2 encounters a Critical Section
2. Sends a request for the lock
3. Acquires the lock
4. Executes Critical Section
5. Releases the lock



On-chip
Interconnect

Accelerated Critical Sections (ACS)

**Critical Section
Request Buffer
(CSRB)**



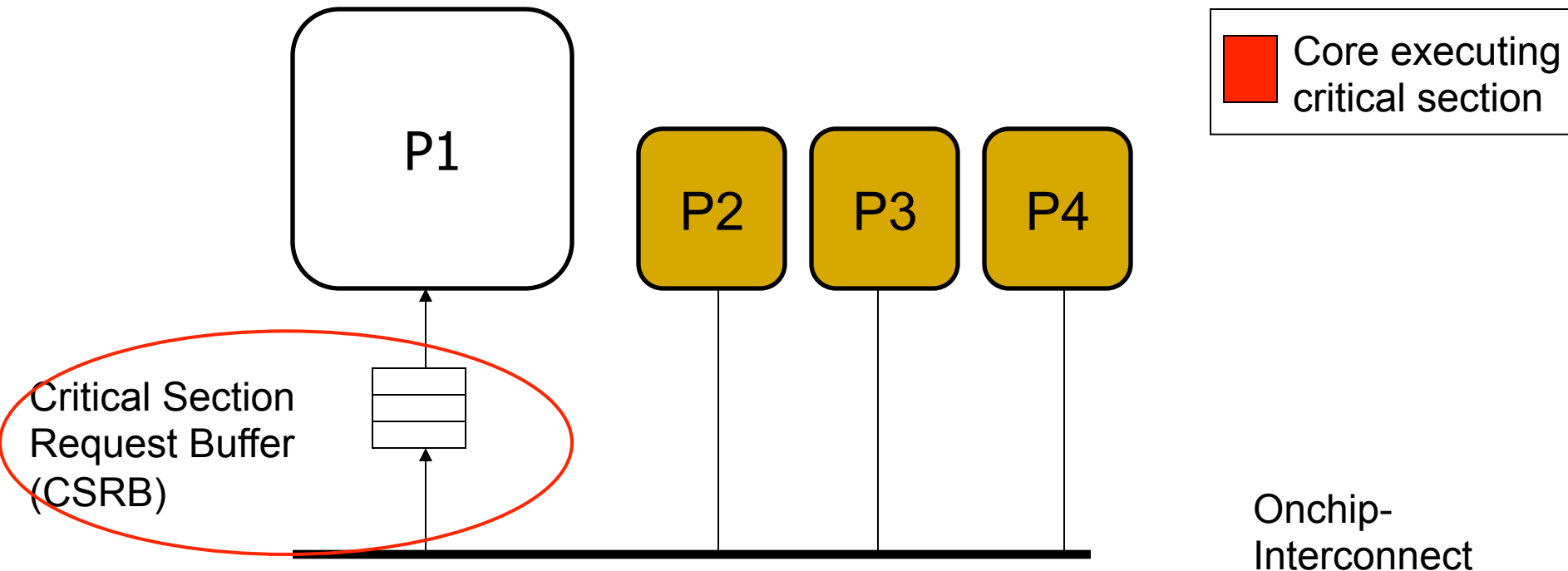
ACMP

- Accelerate Amdahl's serial part **and critical sections** using the large core
 - Suleman et al., "[Accelerating Critical Section Execution with Asymmetric Multi-Core Architectures](#)," ASPLOS 2009, IEEE Micro Top Picks 2010.

Accelerated Critical Sections (ACS)

```
EnterCS()  
    PriorityQ.insert(...)  
LeaveCS()
```

1. P2 encounters a critical section (CSCALL)
2. P2 sends CSCALL Request to CSRB
3. P1 executes Critical Section
4. P1 sends CSDONE signal



ACS Architecture Overview

- ISA extensions
 - `CSCALL LOCK_ADDR, TARGET_PC`
 - `CSRET LOCK_ADDR`
- Compiler/Library inserts `CSCALL/CSRET`
- On a `CSCALL`, the small core:
 - Sends a `CSCALL` request to the large core
 - Arguments: Lock address, Target PC, Stack Pointer, Core ID
 - Stalls and waits for `CSDONE`
- Large Core
 - Critical Section Request Buffer (CSRB)
 - Executes the critical section and sends `CSDONE` to the requesting core

Accelerated Critical Sections (ACS)

Small Core

A = compute()

LOCK X

result = CS(A)

UNLOCK X

print result

Small Core

A = compute()

PUSH A

CSCALL X, Target PC

...

...

...

...

...

...

POP result

print result

Large Core

...

...

...

TPC: Acquire X

POP A

result = CS(A)

PUSH result

Release X

CSRET X

Waiting in
Critical Section
Request Buffer
(CSRB)

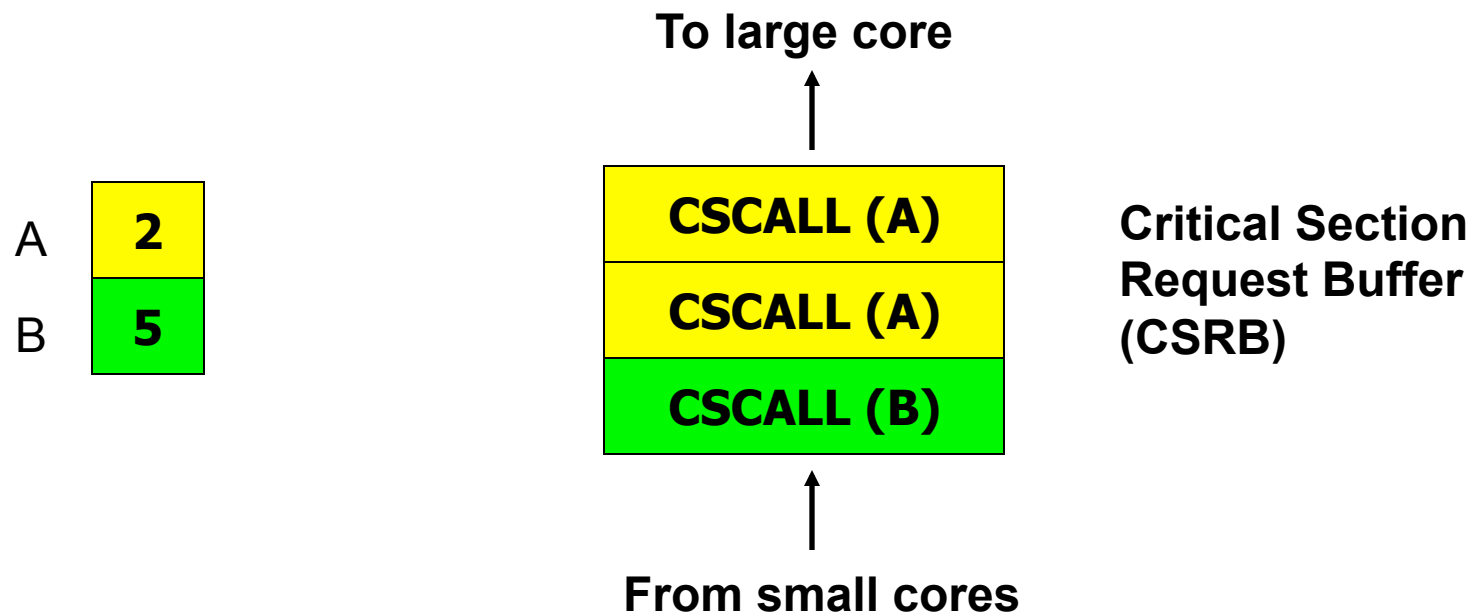
CSCALL Request

Send X, TPC,
STACK_PTR, CORE_ID

CSDONE Response

False Serialization

- ACS can serialize independent critical sections
- Selective Acceleration of Critical Sections (SEL)
 - Saturating counters to track false serialization



ACS Performance Tradeoffs

■ Pluses

- + Faster critical section execution
- + Shared locks stay in one place: better lock locality
- + Shared data stays in large core's (large) caches: better shared data locality, less ping-ponging

■ Minuses

- Large core dedicated for critical sections: reduced parallel throughput
- CSCALL and CSDONE control transfer overhead
- Thread-private data needs to be transferred to large core: worse private data locality

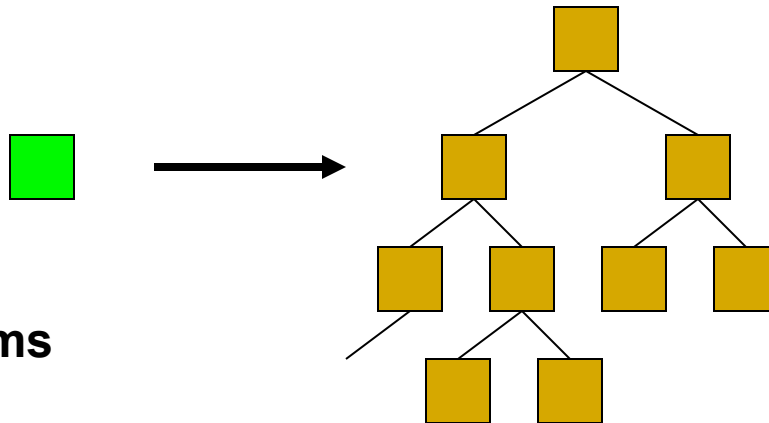
ACS Performance Tradeoffs

- ***Fewer parallel threads vs. accelerated critical sections***
 - Accelerating critical sections offsets loss in throughput
 - As the number of cores (threads) on chip increase:
 - Fractional loss in parallel performance decreases
 - Increased contention for critical sections makes acceleration more beneficial
- ***Overhead of CSCALL/CSDONE vs. better lock locality***
 - ACS avoids “ping-ponging” of locks among caches by keeping them at the large core
- ***More cache misses for private data vs. fewer misses for shared data***

Cache Misses for Private Data

PriorityHeap.insert(NewSubProblems)

Private Data:
NewSubProblems



Shared Data:
The priority heap

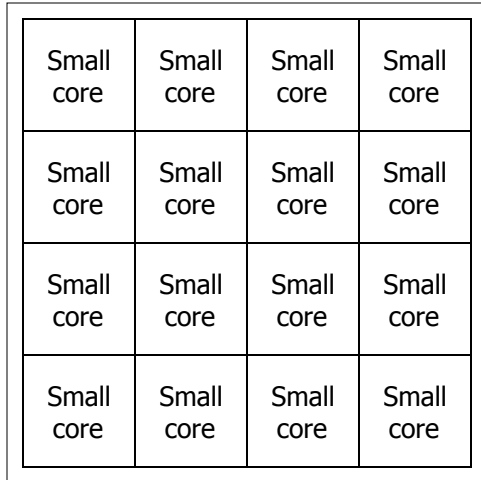
Puzzle Benchmark

ACS Performance Tradeoffs

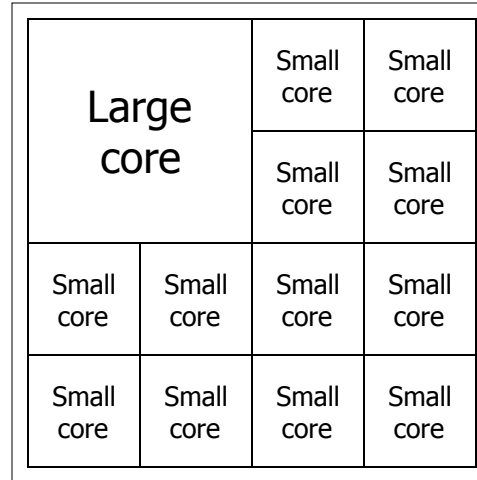
- ***Fewer parallel threads vs. accelerated critical sections***
 - Accelerating critical sections offsets loss in throughput
 - As the number of cores (threads) on chip increase:
 - Fractional loss in parallel performance decreases
 - Increased contention for critical sections makes acceleration more beneficial
- ***Overhead of CSCALL/CSDONE vs. better lock locality***
 - ACS avoids “ping-ponging” of locks among caches by keeping them at the large core
- ***More cache misses for private data vs. fewer misses for shared data***
 - Cache misses reduce if shared data > private data

We will get back to this

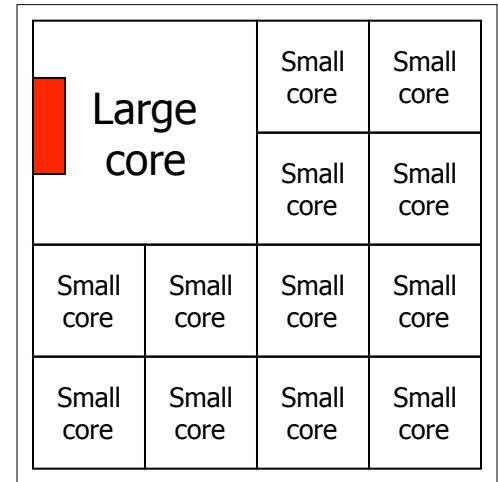
ACS Comparison Points



SCMP



ACMP



ACS

- Conventional locking
 - Large core executes Amdahl's serial part
- Conventional locking
 - Large core executes Amdahl's serial part
- Large core executes Amdahl's serial part and critical sections

Accelerated Critical Sections: Methodology

- Workloads: 12 critical section intensive applications
 - Data mining kernels, sorting, database, web, networking
- Multi-core x86 simulator
 - 1 large and 28 small cores
 - Aggressive stream prefetcher employed at each core
- Details:
 - Large core: 2GHz, out-of-order, 128-entry ROB, 4-wide, 12-stage
 - Small core: 2GHz, in-order, 2-wide, 5-stage
 - Private 32 KB L1, private 256KB L2, 8MB shared L3
 - On-chip interconnect: Bi-directional ring, 5-cycle hop latency

ACS Performance

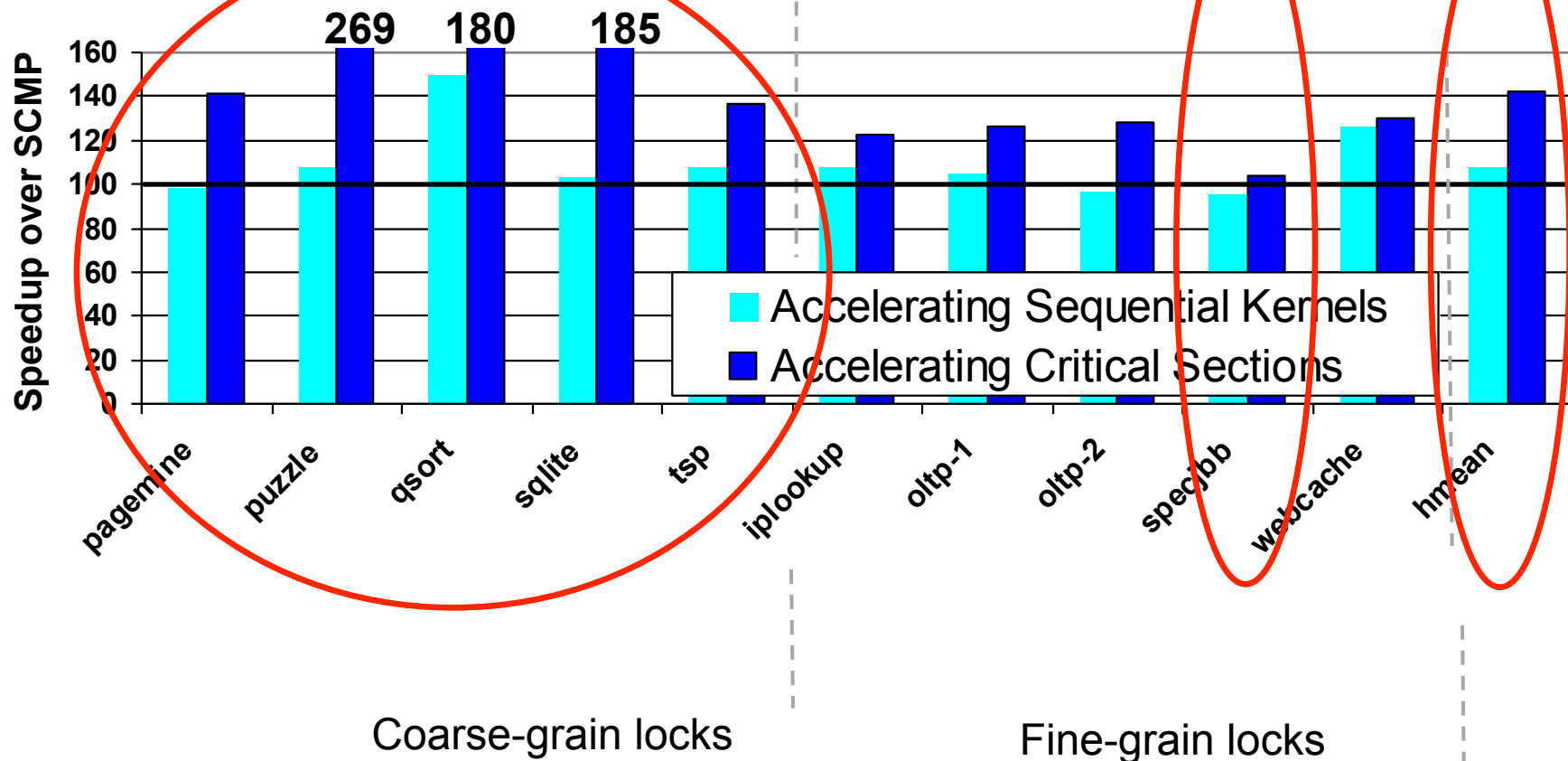
Chip Area = 32 small cores

SCMP = 32 small cores

ACMP = 1 large and 28 small cores

Equal-area comparison

Number of threads = *Best threads*

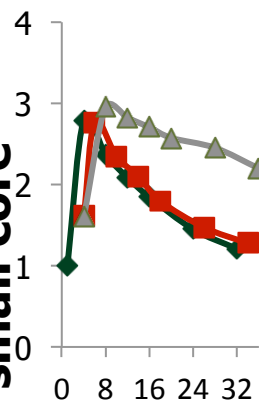


Equal-Area Comparisons

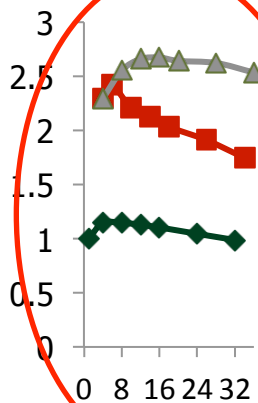
----- **SCMP**
----- **ACMP**
----- **ACS**

Number of threads = No. of cores

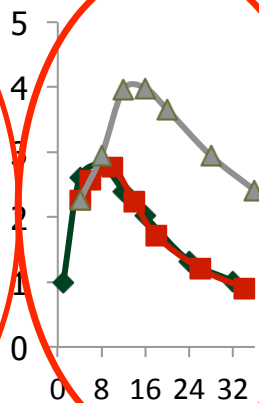
Speedup over a small core



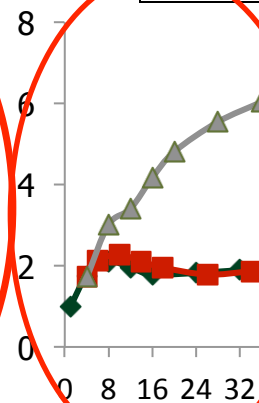
(a) ep



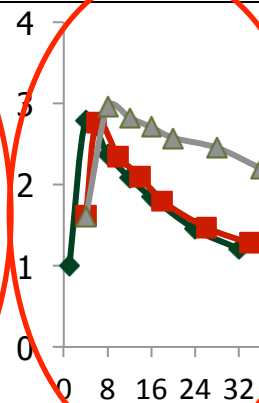
(b) is



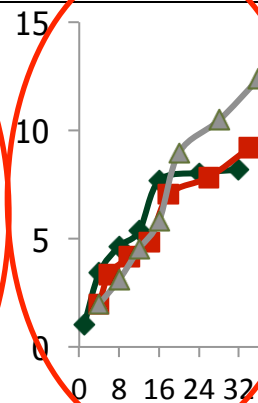
(c) pagemine



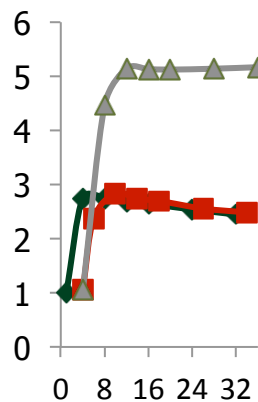
(d) puzzle



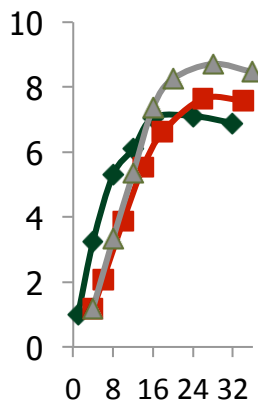
(e) qsort



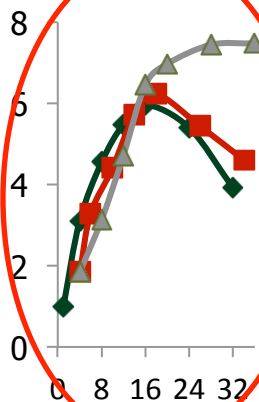
(f) tsp



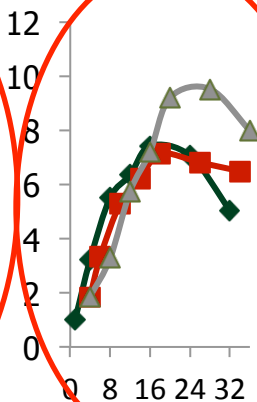
(g) sqlite



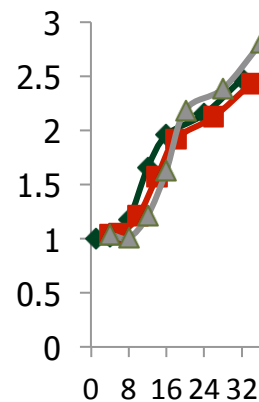
(h) iplookup



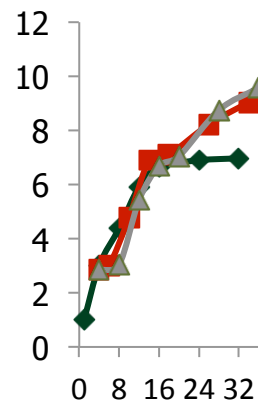
(i) oltp-1



(j) oltp-2



(k) specjbb



(l) webcache

Chip Area (small cores)

ACS Summary

- Critical sections reduce performance and limit scalability
- Accelerate critical sections by executing them on a powerful core
- ACS reduces average execution time by:
 - 34% compared to an equal-area SCMP
 - 23% compared to an equal-area ACMP
- ACS improves scalability of 7 of the 12 workloads
- Generalizing the idea: Accelerate all bottlenecks (“critical paths”) by executing them on a powerful core

Talk Outline

- Problem and Motivation
- How Do We Get There: Examples
- Accelerated Critical Sections (ACS)
- Bottleneck Identification and Scheduling (BIS)
- Staged Execution and Data Marshaling
- Thread Cluster Memory Scheduling (if time permits)
- Ongoing/Future Work
- Conclusions

BIS Summary

- **Problem:** Performance and scalability of multithreaded applications are limited by serializing bottlenecks
 - ❑ different types: critical sections, barriers, slow pipeline stages
 - ❑ importance (criticality) of a bottleneck can change over time
- **Our Goal:** Dynamically identify the most important bottlenecks and accelerate them
 - ❑ How to identify the most critical bottlenecks
 - ❑ How to efficiently accelerate them
- **Solution:** Bottleneck Identification and Scheduling (BIS)
 - ❑ Software: annotate bottlenecks (BottleneckCall, BottleneckReturn) and implement waiting for bottlenecks with a special instruction (BottleneckWait)
 - ❑ Hardware: identify bottlenecks that cause the most thread waiting and accelerate those bottlenecks on large cores of an asymmetric multi-core system
- Improves multithreaded application performance and scalability, outperforms previous work, and performance improves with more cores

Bottlenecks in Multithreaded Applications

Definition: any code segment for which threads contend (i.e. wait)

Examples:

- **Amdahl's serial portions**
 - Only one thread exists → on the critical path
- **Critical sections**
 - Ensure mutual exclusion → likely to be on the critical path if contended
- **Barriers**
 - Ensure all threads reach a point before continuing → the latest thread arriving is on the critical path
- **Pipeline stages**
 - Different stages of a loop iteration may execute on different threads, slowest stage makes other stages wait → on the critical path

Observation: Limiting Bottlenecks Change Over Time

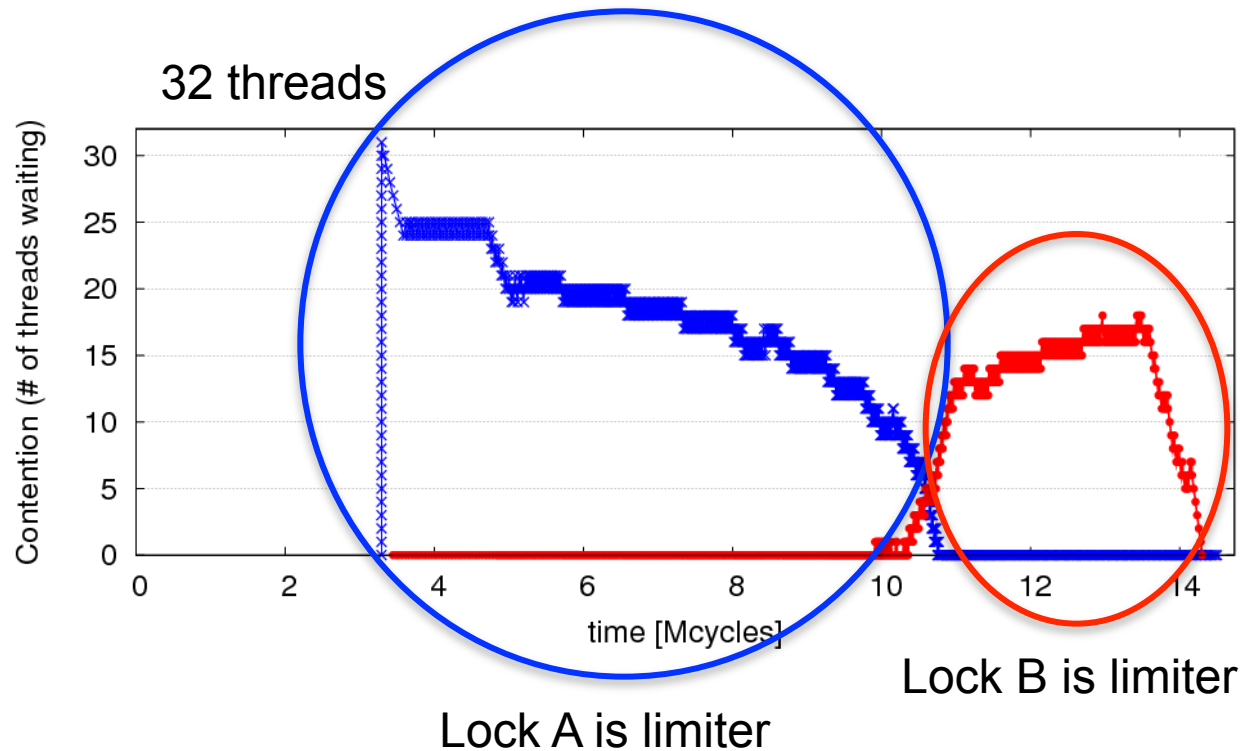
A=full linked list; B=empty linked list
repeat

Lock A
 Traverse list A
 Remove X from A
Unlock A

Compute on X

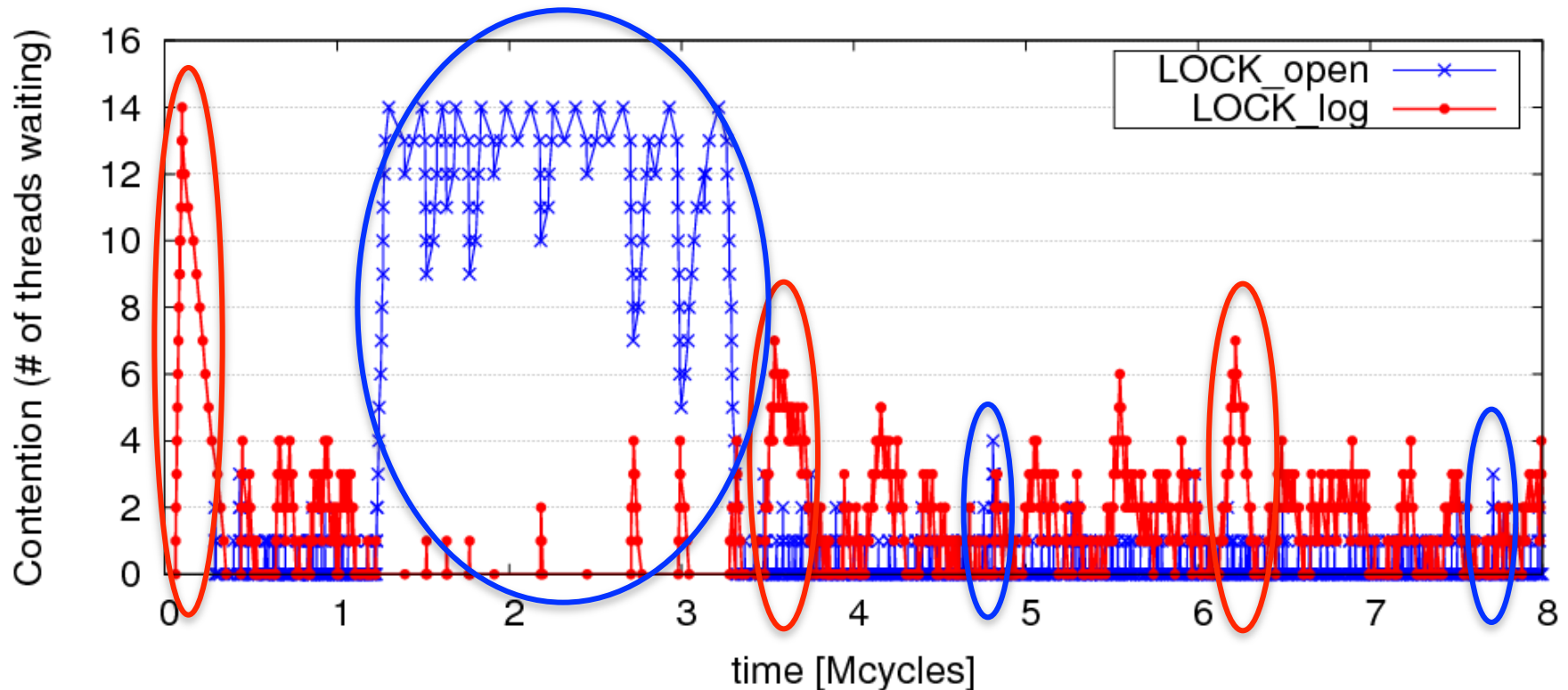
Lock B
 Traverse list B
 Insert X into B
Unlock B

until A is empty



Limiting Bottlenecks Do Change on Real Applications

MySQL running Sysbench queries, 16 threads



Previous Work on Bottleneck Acceleration

- Asymmetric CMP (ACMP) proposals [Annavaram+, ISCA'05]
[Morad+, Comp. Arch. Letters'06] [Suleman+, Tech. Report'07]
 - Accelerate **only the Amdahl's bottleneck**
- Accelerated Critical Sections (ACS) [Suleman+, ASPLOS'09]
 - Accelerate **only critical sections**
 - **Does not take into account importance** of critical sections
- Feedback-Directed Pipelining (FDP) [Suleman+, PACT'10 and PhD thesis'11]
 - Accelerate **only stages with lowest throughput**
 - **Slow to adapt** to phase changes (software based library)

No previous work can accelerate all three types of bottlenecks or quickly adapts to fine-grain changes in the *importance* of bottlenecks

Our goal: general mechanism to identify performance-limiting bottlenecks of any type and accelerate them on an ACMP

Bottleneck Identification and Scheduling (BIS)

- Key insight:
 - Thread waiting reduces parallelism and is likely to reduce performance
 - Code causing the most thread waiting
→ likely critical path
- Key idea:
 - Dynamically identify bottlenecks that cause the most thread waiting
 - Accelerate them (using powerful cores in an ACMP)

Bottleneck Identification and Scheduling (BIS)

Compiler/Library/Programmer

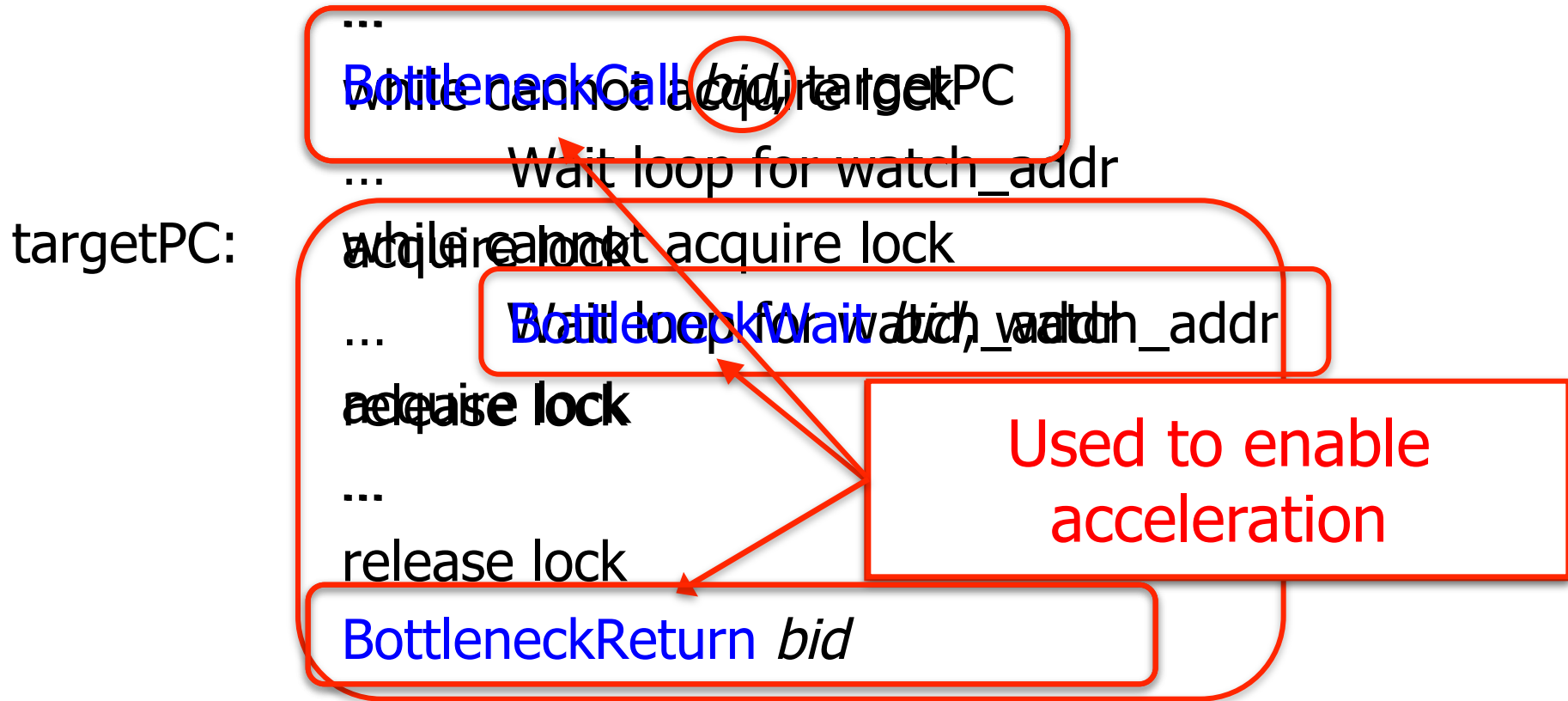
1. Annotate *bottleneck* code
2. Implement *waiting* for bottlenecks

Binary containing
BIS instructions

Hardware

1. Measure *thread waiting cycles (TWC)* for each bottleneck
2. Accelerate bottleneck(s) with the highest TWC

Critical Sections: Code Modifications



Barriers: Code Modifications

...

BottleneckCall *bid*, targetPC

enter barrier

while not all threads in barrier

BottleneckWait *bid*, watch_addr

exit barrier

...

targetPC: code running for the barrier

...

BottleneckReturn *bid*

Pipeline Stages: Code Modifications

BottleneckCall *bid*, targetPC

...

targetPC:

while not done

while empty queue

BottleneckWait prev_bid

dequeue work

do the work ...

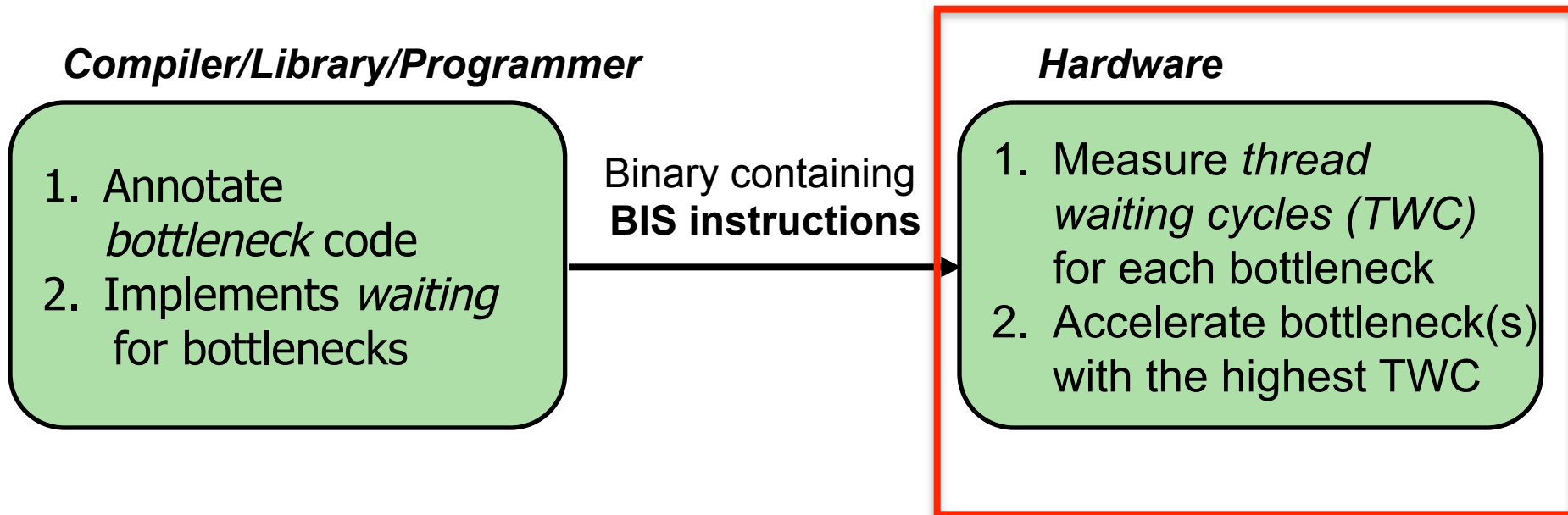
while full queue

BottleneckWait next_bid

enqueue next work

BottleneckReturn *bid*

Bottleneck Identification and Scheduling (BIS)

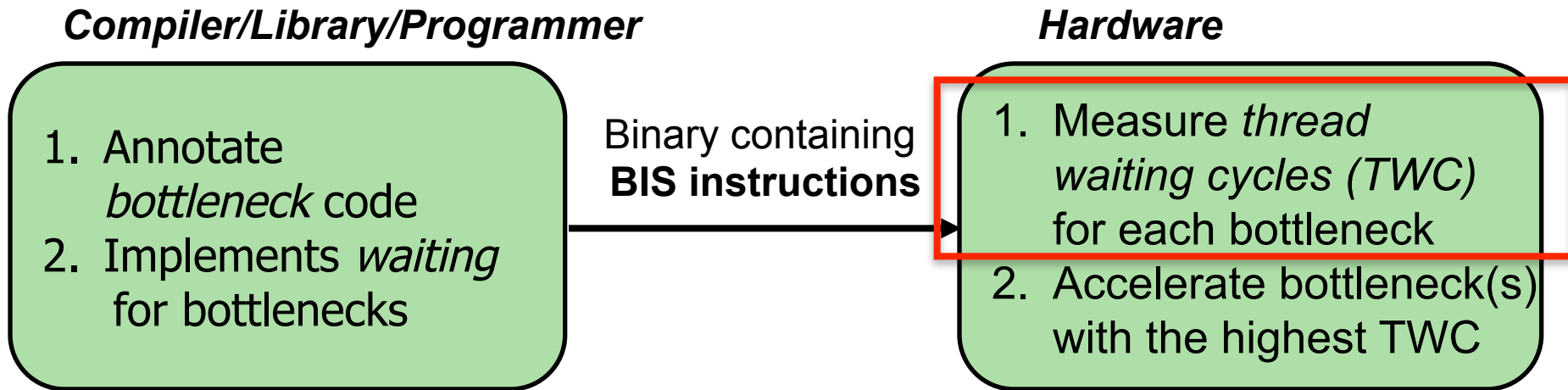


BIS: Hardware Overview

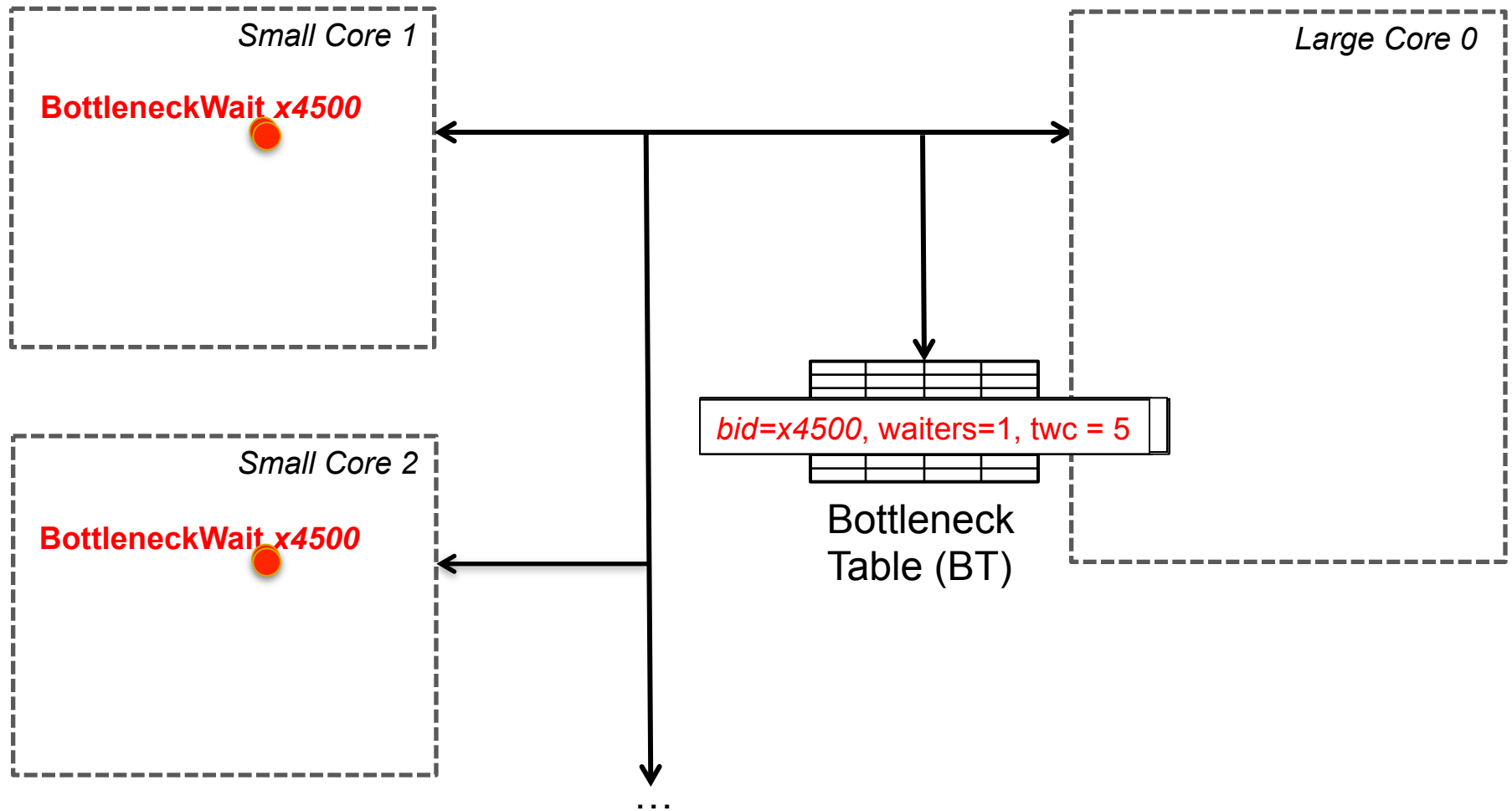
- Performance-limiting bottleneck **identification and acceleration are independent tasks**
- Acceleration can be accomplished in multiple ways
 - ❑ Increasing core frequency/voltage
 - ❑ Prioritization in shared resources [Ebrahimi+, MICRO'11]
 - ❑ **Migration to faster cores in an Asymmetric CMP**

Small core	Small core	Large core	
Small core	Small core		
Small core	Small core	Small core	Small core
Small core	Small core	Small core	Small core

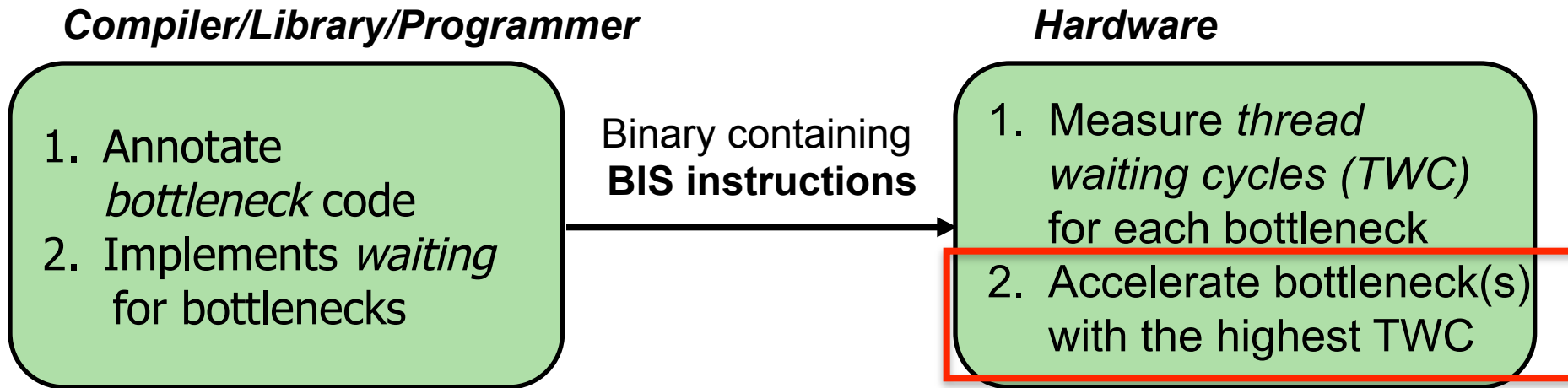
Bottleneck Identification and Scheduling (BIS)



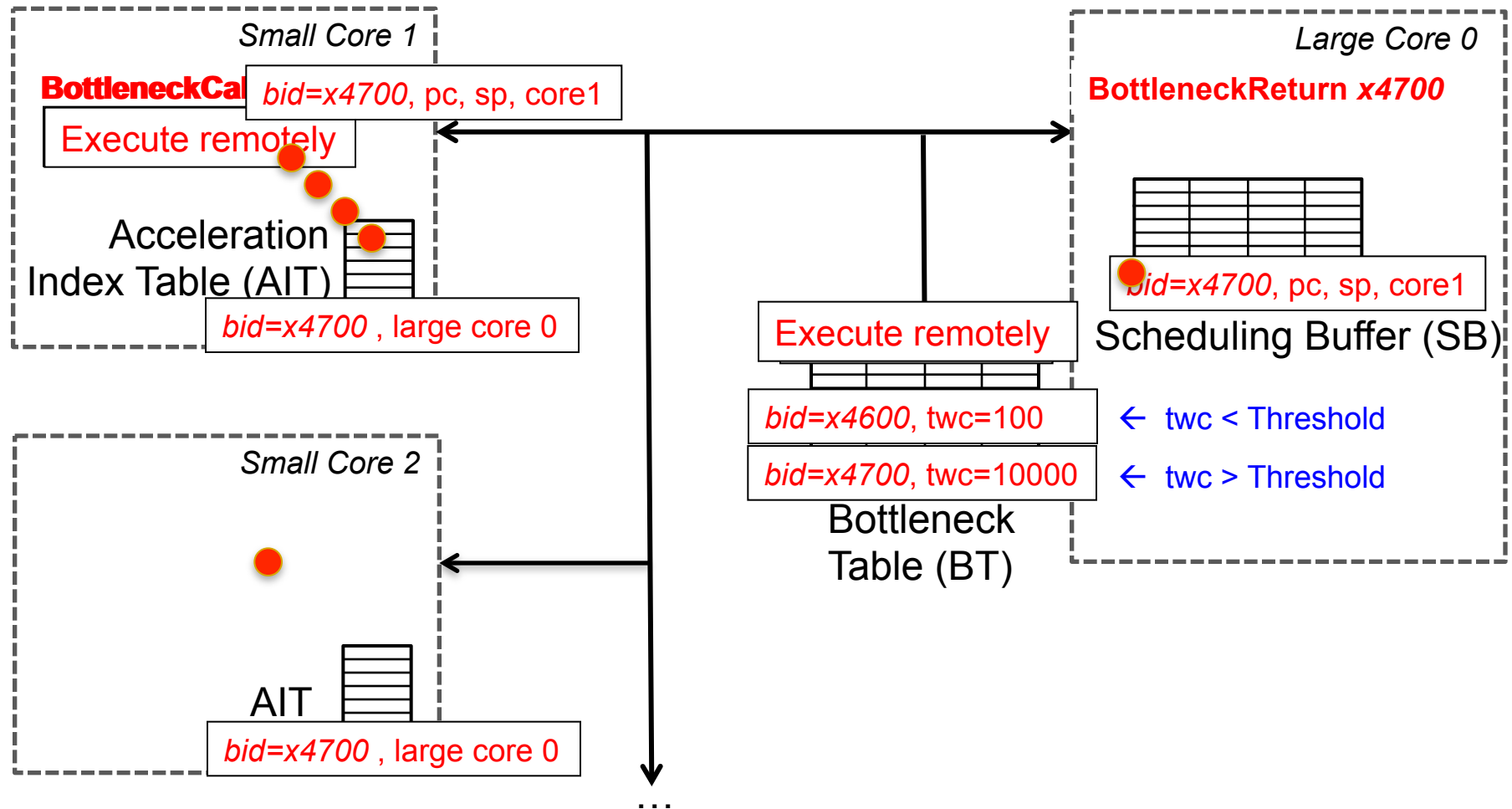
Determining Thread Waiting Cycles for Each Bottleneck



Bottleneck Identification and Scheduling (BIS)



Bottleneck Acceleration



BIS Mechanisms

- Basic mechanisms for BIS:
 - Determining Thread Waiting Cycles ✓
 - Accelerating Bottlenecks ✓
- Mechanisms to improve performance and generality of BIS:
 - Dealing with false serialization
 - Preemptive acceleration
 - Support for multiple large cores

False Serialization and Starvation

- **Observation:** Bottlenecks are picked from Scheduling Buffer in Thread Waiting Cycles order
- **Problem:** An independent bottleneck that is ready to execute has to wait for another bottleneck that has higher thread waiting cycles → **False serialization**
- **Starvation:** Extreme false serialization
- **Solution:** Large core detects when a bottleneck is ready to execute in the Scheduling Buffer but it cannot → sends the bottleneck back to the small core

Preemptive Acceleration

- **Observation:** A bottleneck executing on a small core can become the bottleneck with the highest thread waiting cycles
- **Problem:** This bottleneck should really be accelerated (i.e., executed on the large core)
- **Solution:** The Bottleneck Table detects the situation and sends a preemption signal to the small core. Small core:
 - saves register state on stack, ships the bottleneck to the large core
- Main acceleration mechanism for barriers and pipeline stages

Support for Multiple Large Cores

- **Objective:** to accelerate independent bottlenecks
- Each large core has its own Scheduling Buffer (shared by all of its SMT threads)
- Bottleneck Table assigns each bottleneck to a fixed large core context to
 - preserve cache locality
 - avoid busy waiting
- Preemptive acceleration extended to send multiple instances of a bottleneck to different large core contexts

Hardware Cost

- Main structures:

- Bottleneck Table (BT): global 32-entry associative cache, minimum-Thread-Waiting-Cycle replacement
- Scheduling Buffers (SB): one table per large core, as many entries as small cores
- Acceleration Index Tables (AIT): one 32-entry table per small core

- Off the critical path

- Total storage cost for 56-small-cores, 2-large-cores < 19 KB

BIS Performance Trade-offs

- **Faster bottleneck execution** vs. **fewer parallel threads**
 - ❑ Acceleration offsets loss of parallel throughput with large core counts
- **Better shared data locality** vs. **worse private data locality**
 - ❑ Shared data stays on large core (good)
 - ❑ Private data migrates to large core (bad, but latency hidden with Data Marshaling [Suleman+, ISCA' 10])
- **Benefit of acceleration** vs. **migration latency**
 - ❑ Migration latency usually hidden by waiting (good)
 - ❑ Unless bottleneck not contended (bad, but likely not on critical path)

Methodology

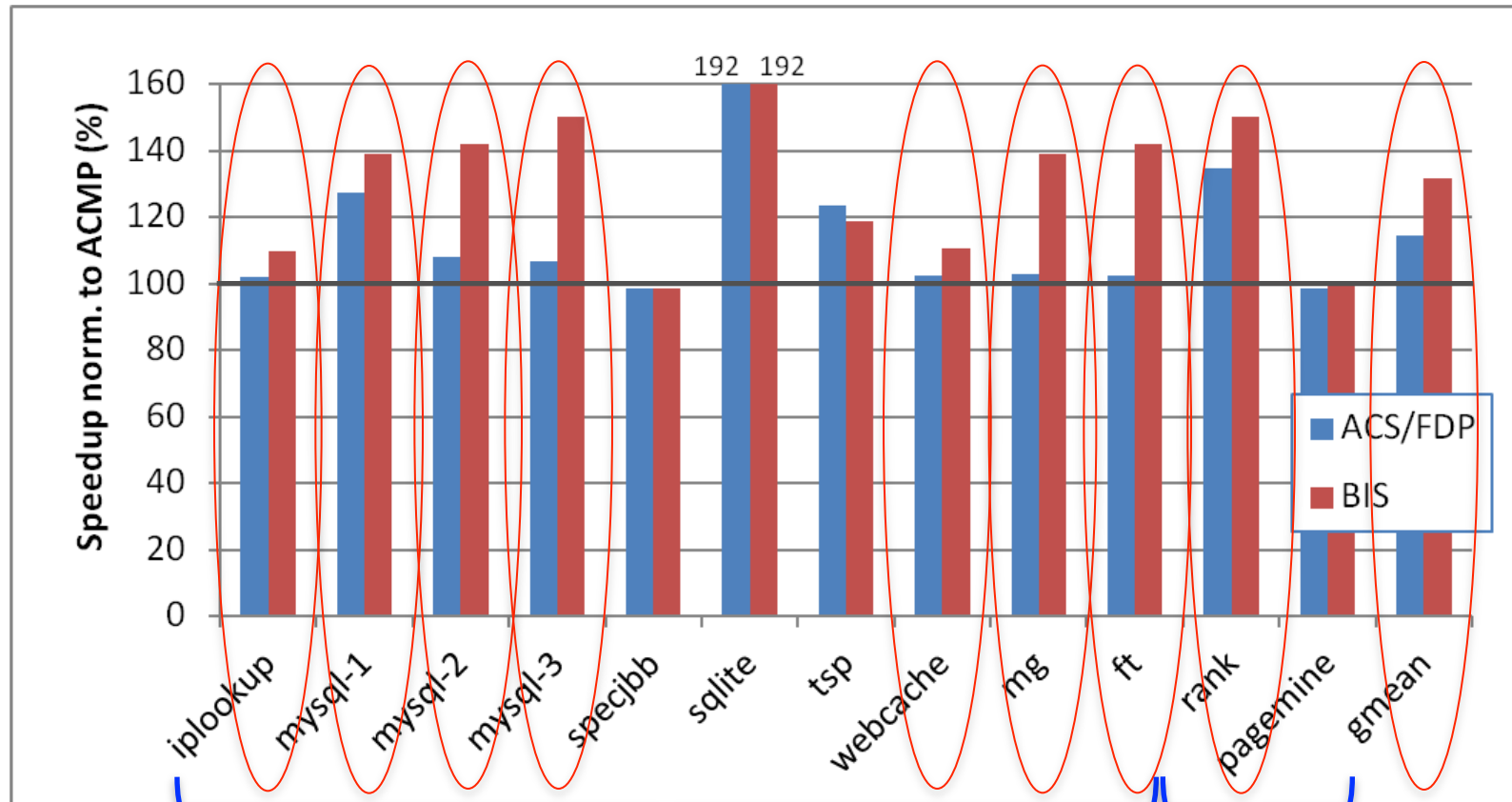
- Workloads: 8 critical section intensive, 2 barrier intensive and 2 pipeline-parallel applications
 - Data mining kernels, scientific, database, web, networking, specjbb
- Cycle-level multi-core x86 simulator
 - 8 to 64 small-core-equivalent area, 0 to 3 large cores, SMT
 - 1 large core is area-equivalent to 4 small cores
- Details:
 - Large core: 4GHz, out-of-order, 128-entry ROB, 4-wide, 12-stage
 - Small core: 4GHz, in-order, 2-wide, 5-stage
 - Private 32KB L1, private 256KB L2, shared 8MB L3
 - On-chip interconnect: Bi-directional ring, 2-cycle hop latency

BIS Comparison Points (Area-Equivalent)

- SCMP (Symmetric CMP)
 - ❑ All small cores
 - ❑ Results in the paper
- **ACMP** (Asymmetric CMP)
 - ❑ Accelerates only Amdahl's serial portions
 - ❑ **Our baseline**
- **ACS** (Accelerated Critical Sections)
 - ❑ Accelerates only critical sections and Amdahl's serial portions
 - ❑ Applicable to multithreaded workloads
(**iplookup, mysql, specjbb, sqlite, tsp, webcache, mg, ft**)
- **FDP** (Feedback-Directed Pipelining)
 - ❑ Accelerates only slowest pipeline stages
 - ❑ Applicable to pipeline-parallel workloads (**rank, pagemine**)

BIS Performance Improvement

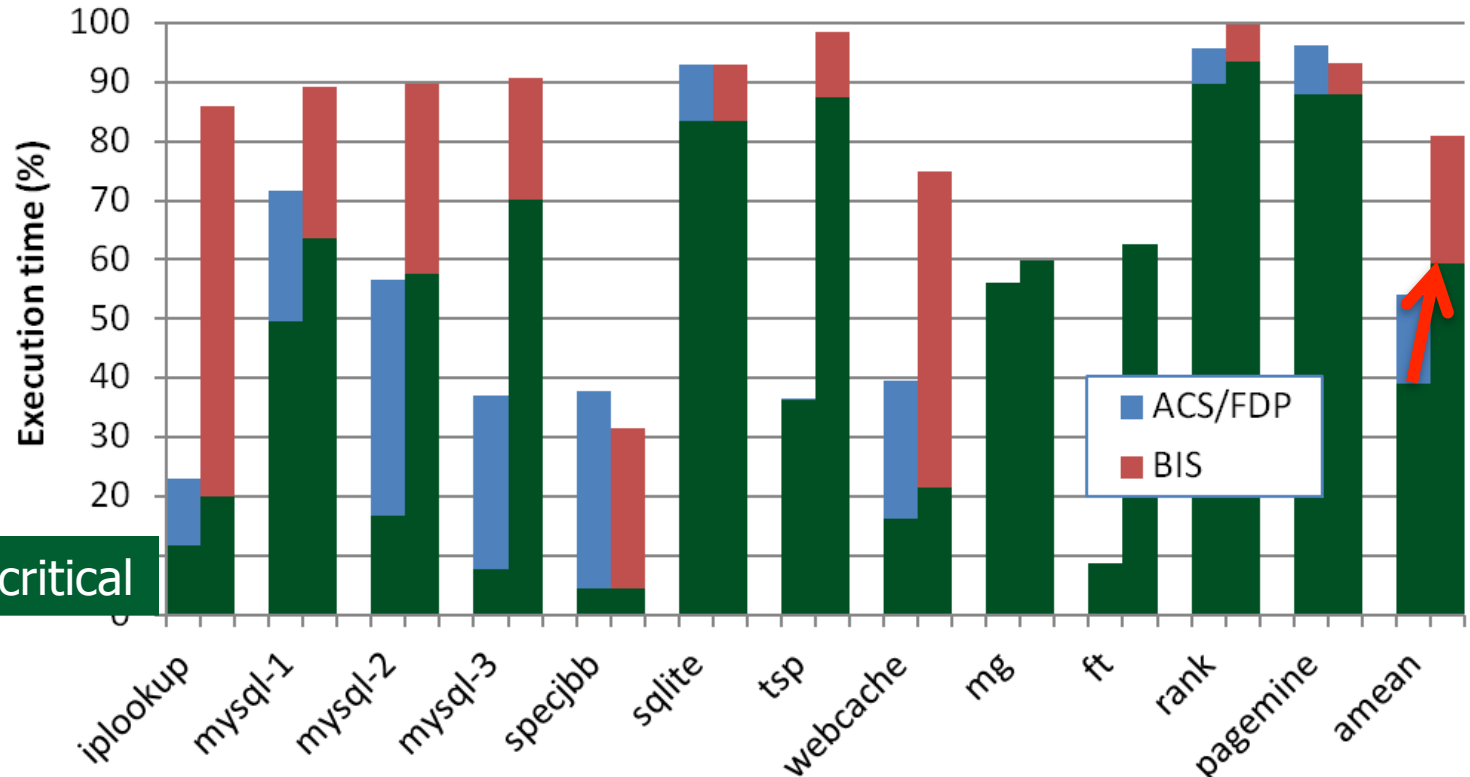
Optimal number of threads, 28 small cores, 1 large core



- BIS outperforms ACS/FDP by 15% and ACMP by 32%
limiting bottlenecks change over time, which ACS cannot accelerate
- BIS improves scalability on 4 of the benchmarks

Why Does BIS Work?

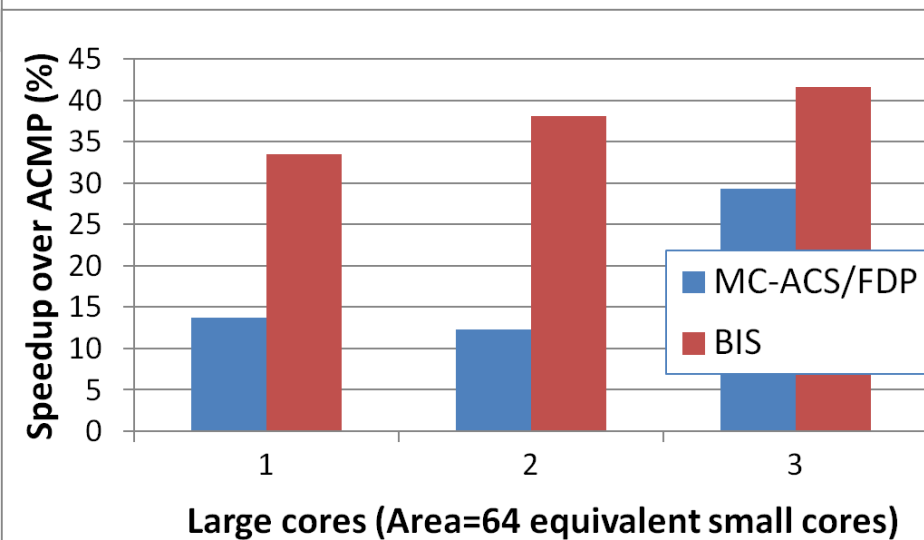
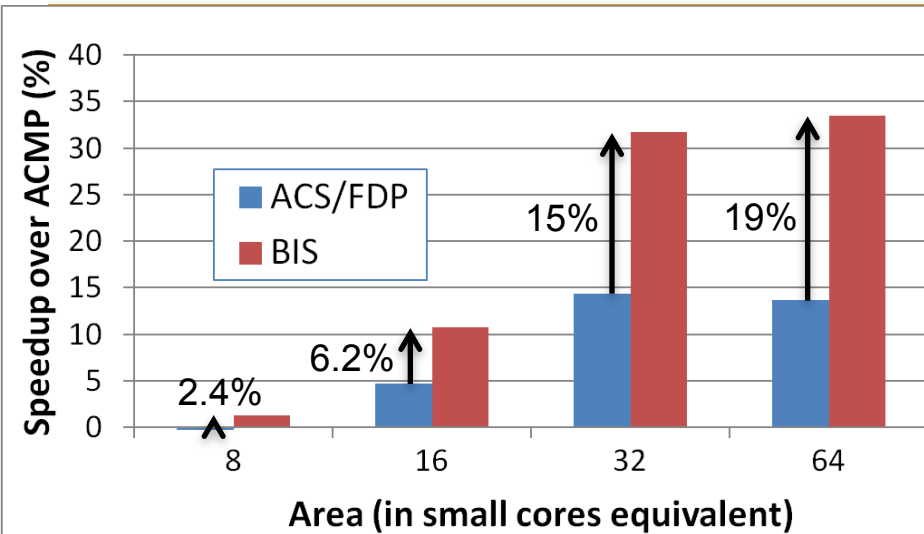
Fraction of execution time spent on predicted-important bottlenecks



Actually critical

- Coverage: fraction of program critical path that is actually identified as bottlenecks
 - 39% (ACS/FDP) to 59% (BIS)
- Accuracy: identified bottlenecks on the critical path over total identified bottlenecks
 - 72% (ACS/FDP) to 73.5% (BIS)

BIS Scaling Results



Performance increases with:

1) More small cores

- Contention due to bottlenecks increases
- Loss of parallel throughput due to large core reduces

2) More large cores

- Can accelerate independent bottlenecks
- *Without reducing parallel throughput (enough cores)*

BIS Summary

- **Serializing bottlenecks of different types** limit performance of multithreaded applications: **Importance changes over time**
- BIS is a hardware/software cooperative solution:
 - ❑ **Dynamically identifies bottlenecks** that cause the **most thread waiting** and **accelerates** them on large cores of an ACMP
 - ❑ Applicable to critical sections, barriers, pipeline stages
- BIS improves application performance and scalability:
 - ❑ 15% speedup over ACS/FDP
 - ❑ Can accelerate multiple independent critical bottlenecks
 - ❑ Performance benefits increase with more cores
- Provides **comprehensive fine-grained bottleneck acceleration for future ACMPs** with little or no programmer effort

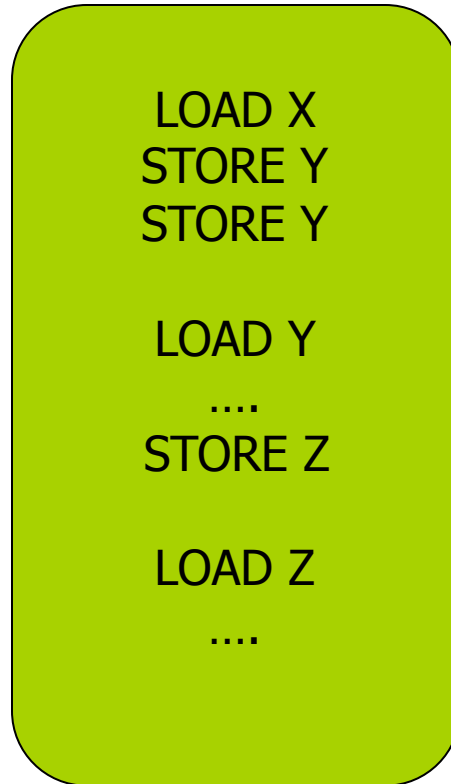
Talk Outline

- Problem and Motivation
- How Do We Get There: Examples
- Accelerated Critical Sections (ACS)
- Bottleneck Identification and Scheduling (BIS)
- **Staged Execution** and Data Marshaling
- Thread Cluster Memory Scheduling (if time permits)
- Ongoing/Future Work
- Conclusions

Staged Execution Model (I)

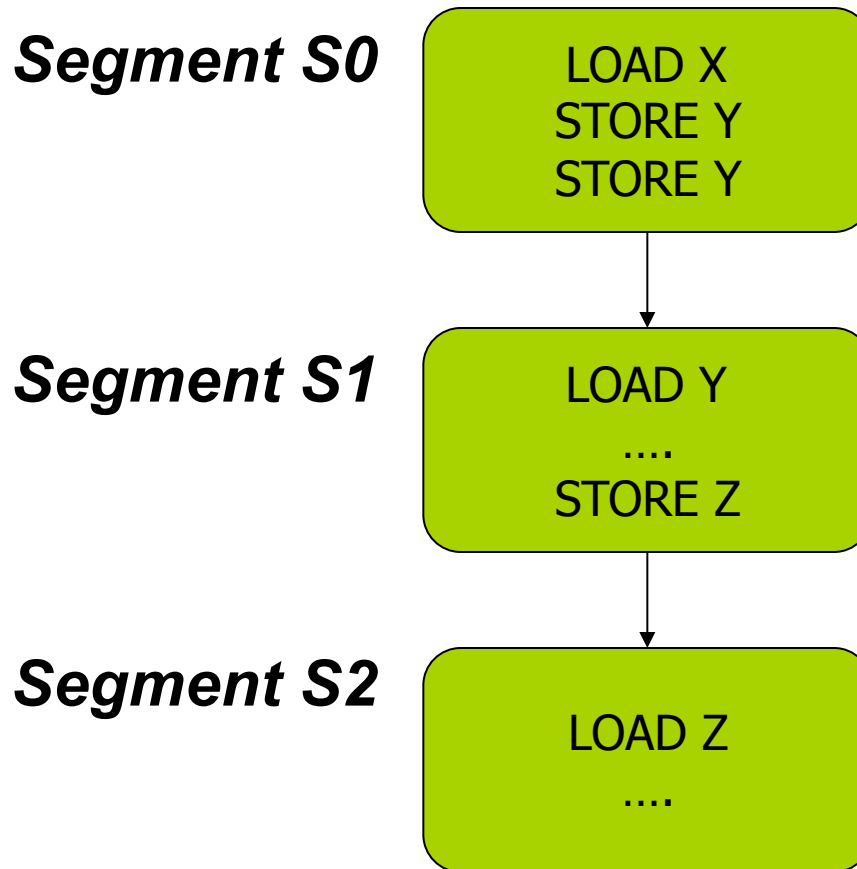
- Goal: speed up a program by dividing it up into pieces
- Idea
 - Split program code into **segments**
 - Run each segment on the core best-suited to run it
 - Each core assigned a work-queue, storing segments to be run
- Benefits
 - Accelerates segments/critical-paths using specialized/heterogeneous cores
 - Exploits inter-segment parallelism
 - Improves locality of within-segment data
- Examples
 - Accelerated critical sections, Bottleneck identification and scheduling
 - Producer-consumer pipeline parallelism
 - Task parallelism (Cilk, Intel TBB, Apple Grand Central Dispatch)
 - Special-purpose cores and functional units

Staged Execution Model (II)

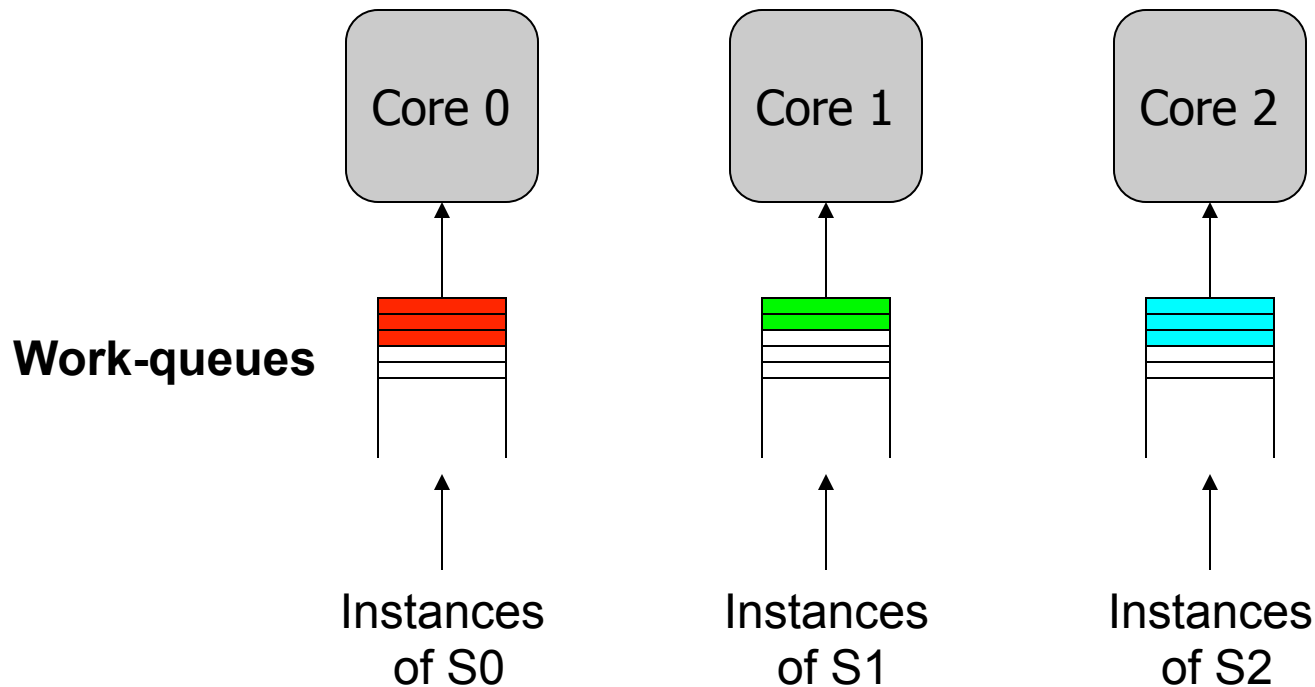


Staged Execution Model (III)

Split code into segments



Staged Execution Model (IV)



Staged Execution Model: Segment Spawning

Core 0

Core 1

Core 2

S0

LOAD X
STORE Y
STORE Y

S1

LOAD Y
....
STORE Z

S2

LOAD Z
....

```
graph LR; S0[S0: LOAD X, STORE Y, STORE Y] --> S1[S1: LOAD Y, ..., STORE Z]; S1 --> S2[S2: LOAD Z, ...];
```

Staged Execution Model: Two Examples

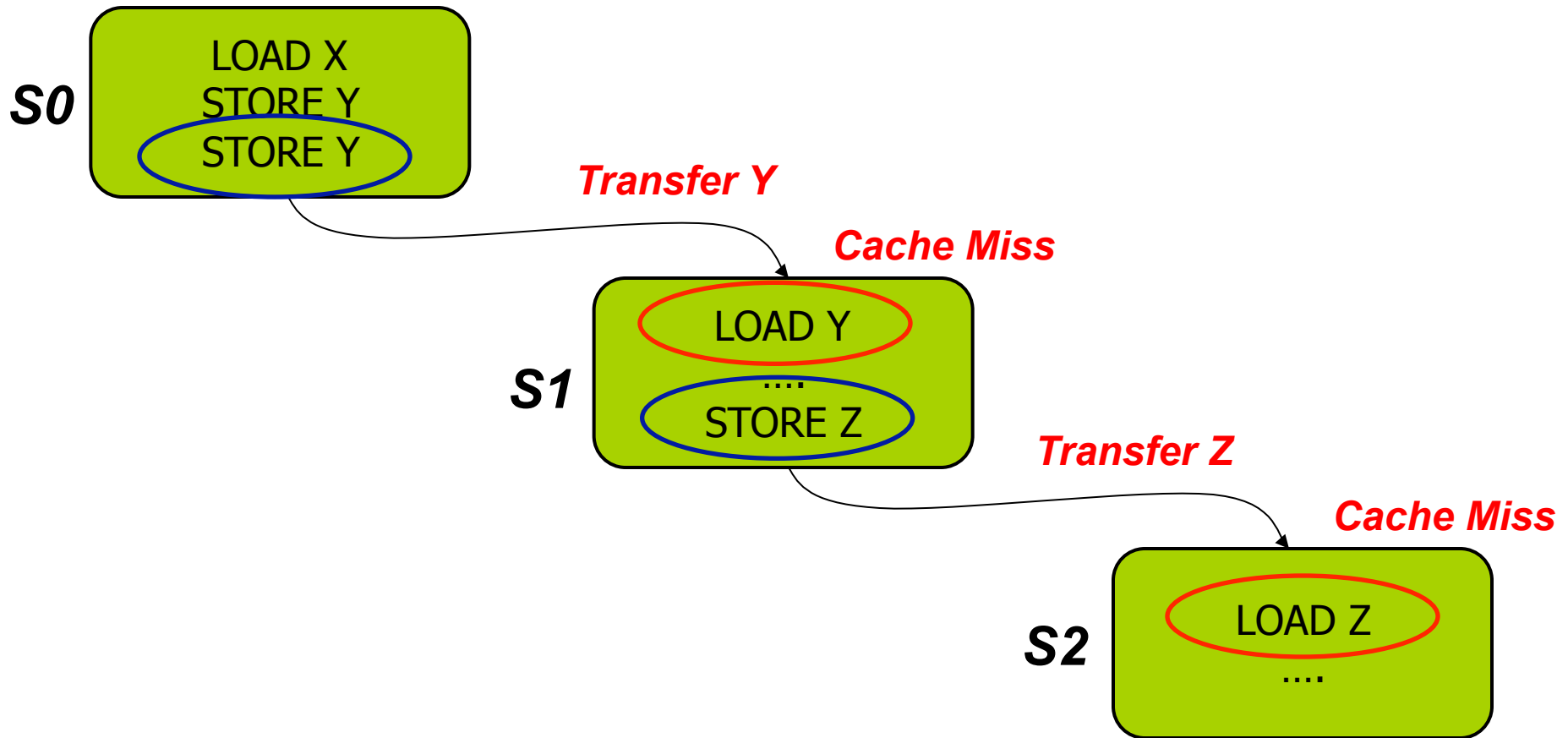
- **Accelerated Critical Sections** [Suleman et al., ASPLOS 2009]
 - Idea: Ship critical sections to a large core in an asymmetric CMP
 - Segment 0: Non-critical section
 - Segment 1: Critical section
 - Benefit: Faster execution of critical section, reduced serialization, improved lock and shared data locality
- **Producer-Consumer Pipeline Parallelism**
 - Idea: Split a loop iteration into multiple “pipeline stages” where one stage consumes data produced by the next stage → each stage runs on a different core
 - Segment N: Stage N
 - Benefit: Stage-level parallelism, better locality → faster execution

Problem: Locality of Inter-segment Data

Core 0

Core 1

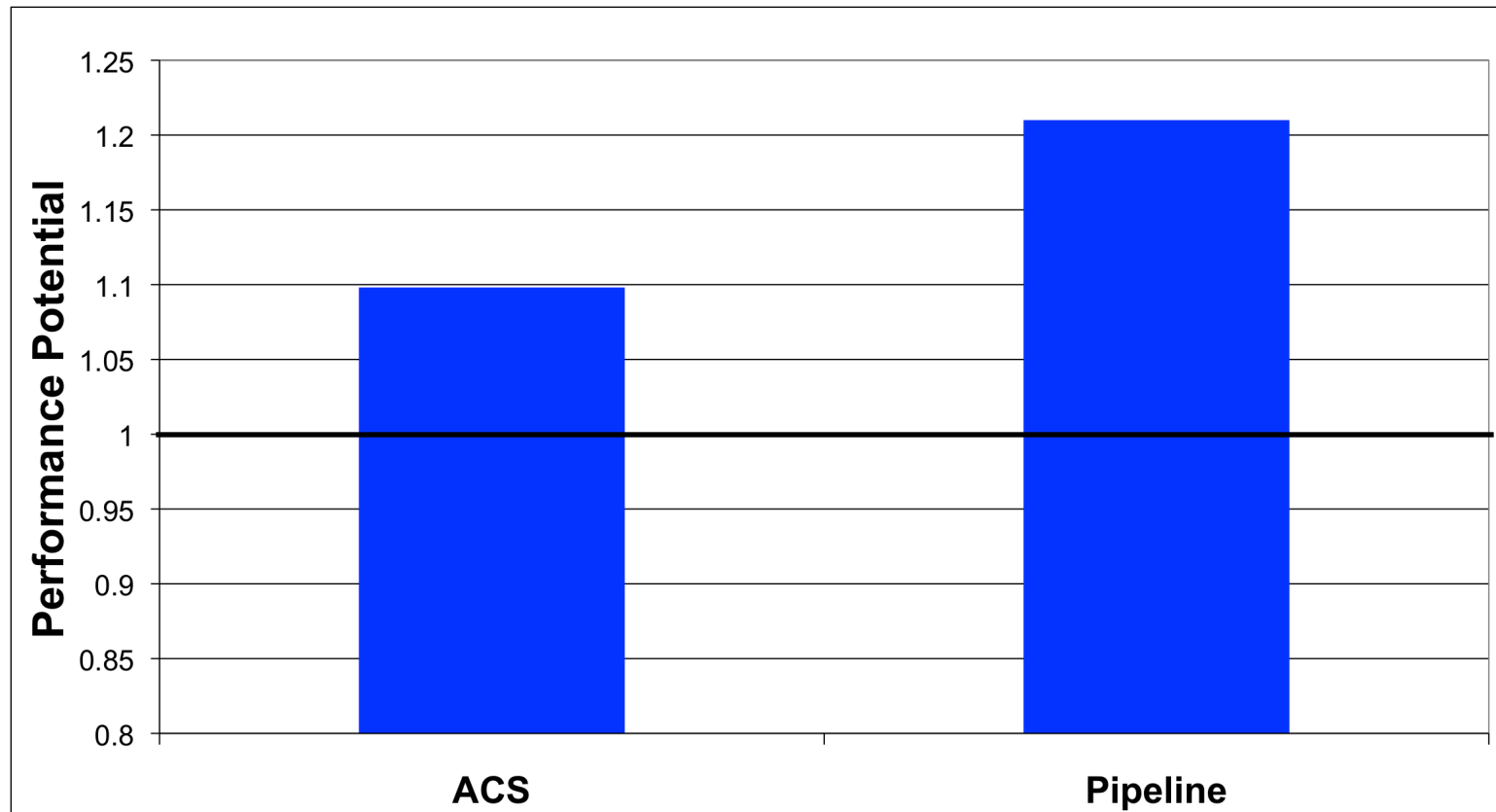
Core 2



Problem: Locality of Inter-segment Data

- Accelerated Critical Sections [Suleman et al., ASPLOS 2010]
 - Idea: Ship critical sections to a large core in an ACMP
 - Problem: Critical section incurs a cache miss when it touches data produced in the non-critical section (i.e., thread private data)
- Producer-Consumer Pipeline Parallelism
 - Idea: Split a loop iteration into multiple “pipeline stages” → each stage runs on a different core
 - Problem: A stage incurs a cache miss when it touches data produced by the previous stage
- Performance of Staged Execution limited by inter-segment cache misses

What if We Eliminated All Inter-segment Misses?



Talk Outline

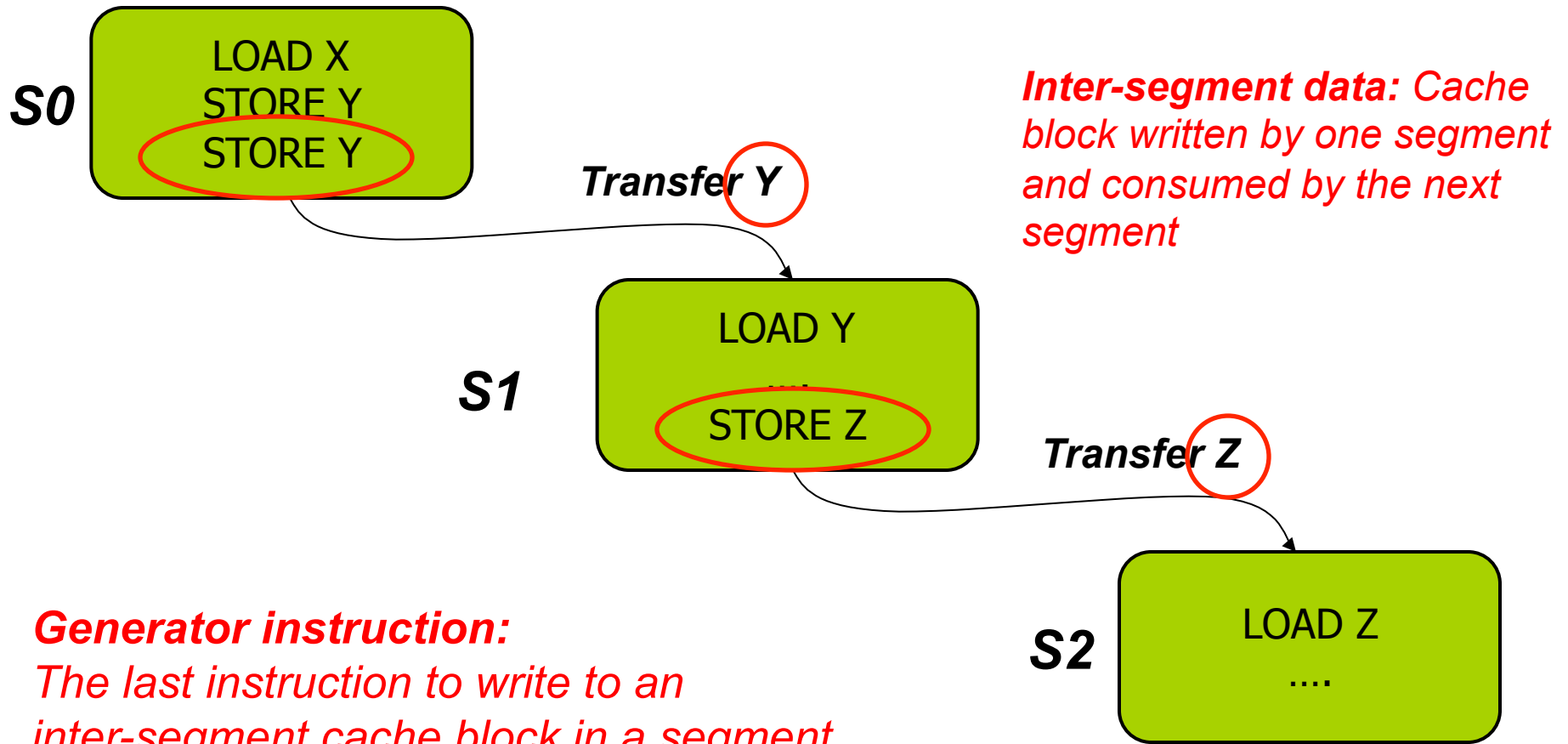
- Problem and Motivation
- How Do We Get There: Examples
- Accelerated Critical Sections (ACS)
- Bottleneck Identification and Scheduling (BIS)
- Staged Execution and **Data Marshaling**
- Thread Cluster Memory Scheduling (if time permits)
- Ongoing/Future Work
- Conclusions

Terminology

Core 0

Core 1

Core 2



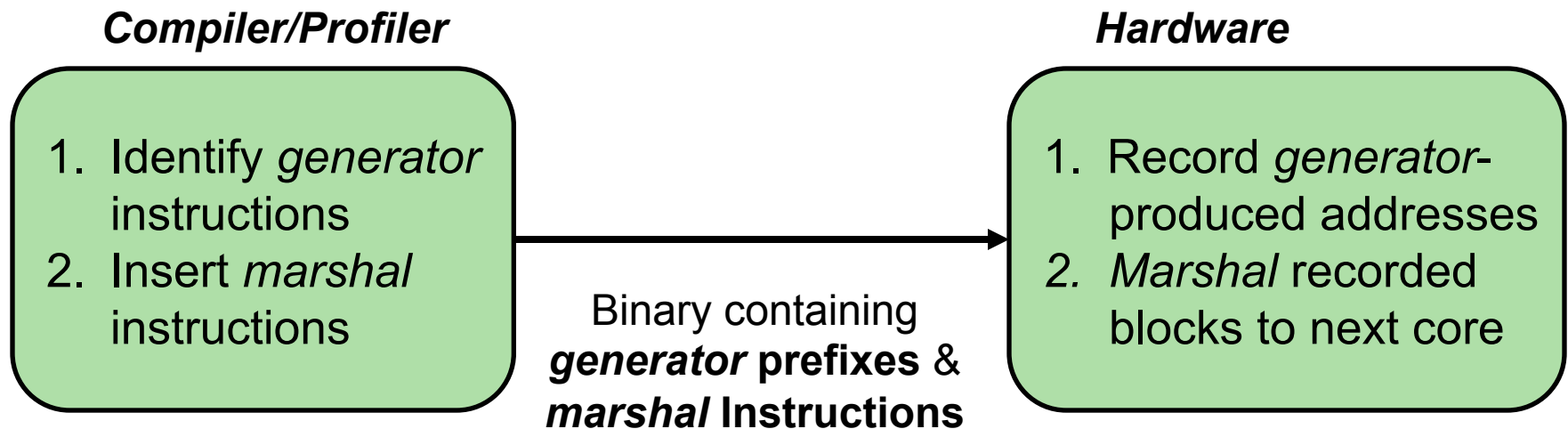
Generator instruction:

The last instruction to write to an inter-segment cache block in a segment

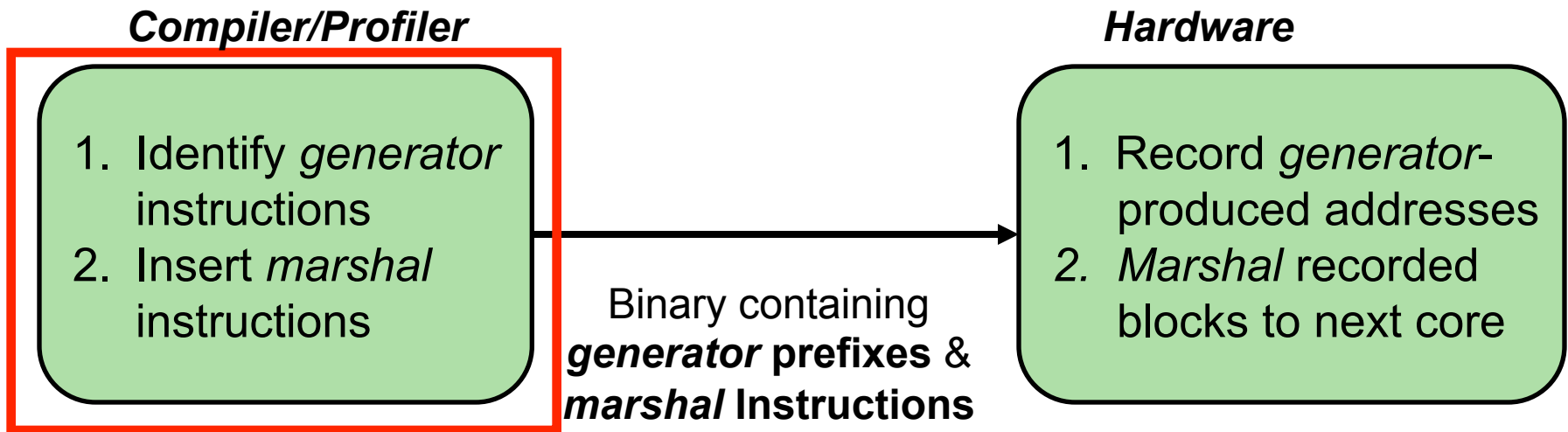
Key Observation and Idea

- Observation: Set of generator instructions is stable over execution time and across input sets
- Idea:
 - Identify the generator instructions
 - Record cache blocks produced by generator instructions
 - Proactively send such cache blocks to the next segment's core before initiating the next segment
- Suleman et al., “Data Marshaling for Multi-Core Architectures,” ISCA 2010, IEEE Micro Top Picks 2011.

Data Marshaling



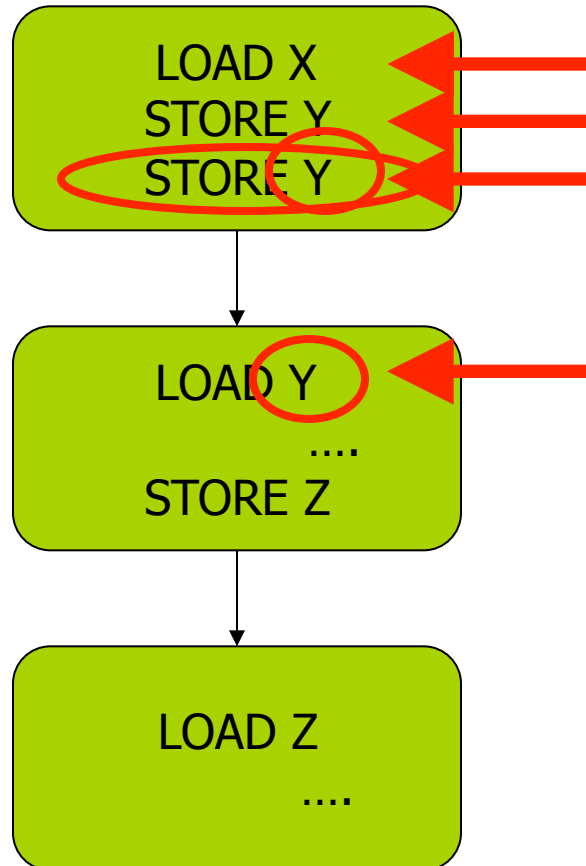
Data Marshaling



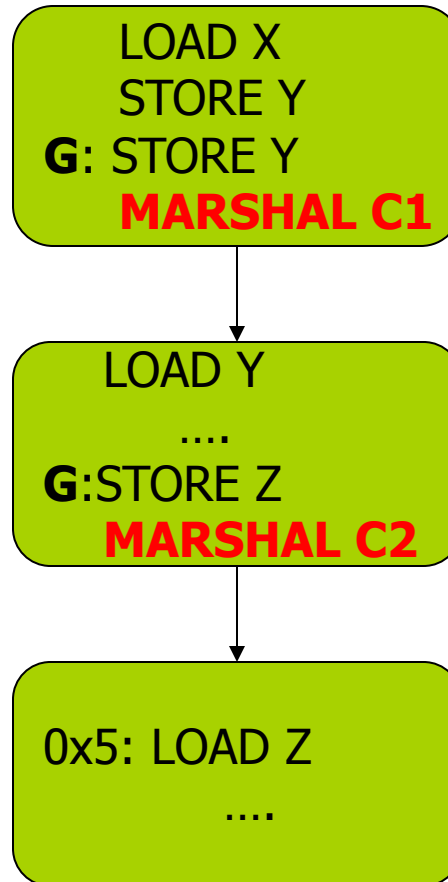
Profiling Algorithm

Inter-segment data

*Mark as Generator
Instruction*



Marshal Instructions



When to send (Marshal)

Where to send (C1)

DM Support/Cost

- Profiler/Compiler: Generators, marshal instructions
- ISA: Generator prefix, marshal instructions
- Library/Hardware: Bind next segment ID to a physical core
- Hardware
 - Marshal Buffer
 - Stores physical addresses of cache blocks to be marshaled
 - 16 entries enough for almost all workloads → 96 bytes per core
 - Ability to execute generator prefixes and marshal instructions
 - Ability to push data to another cache

DM: Advantages, Disadvantages

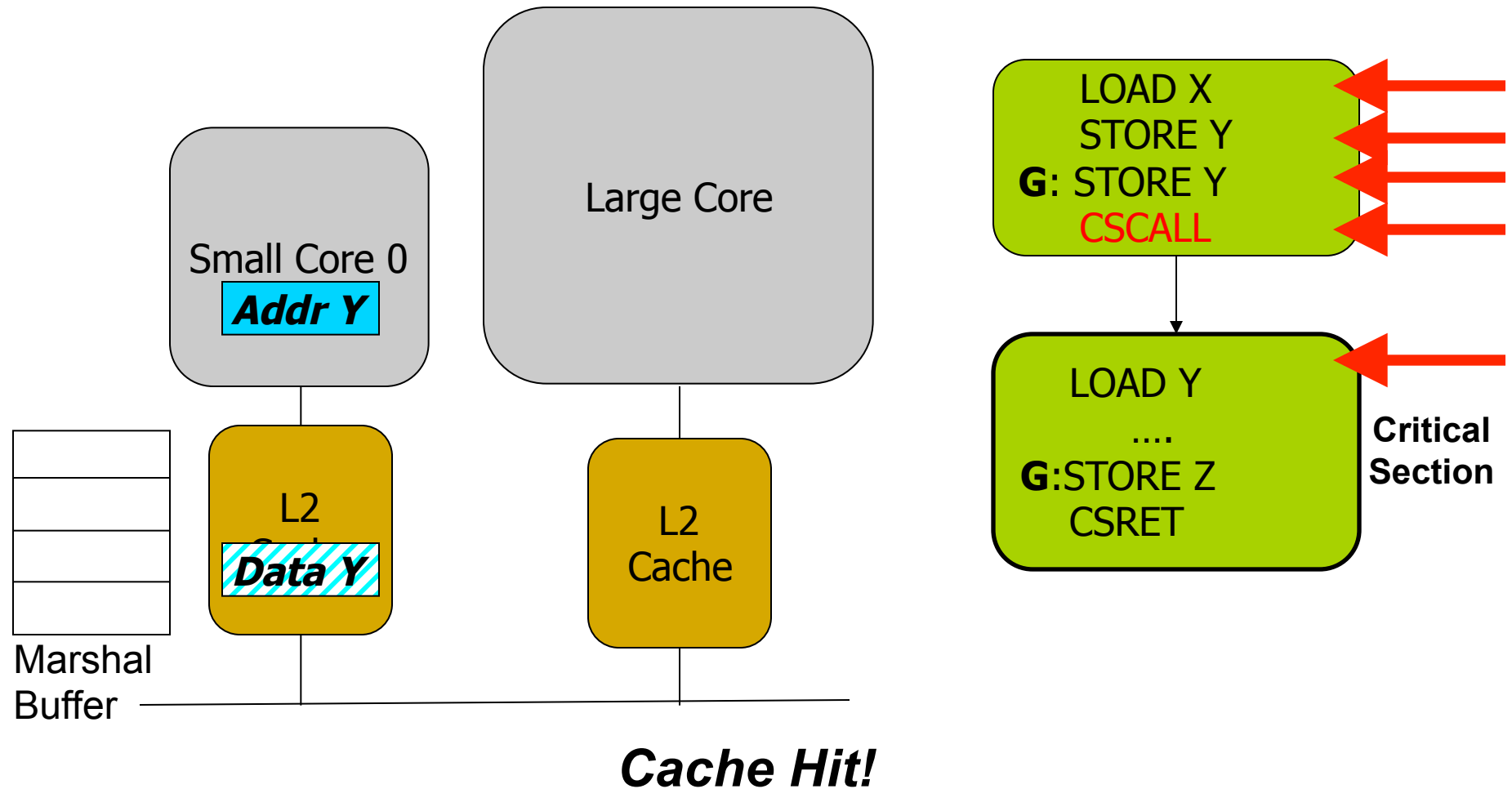
■ Advantages

- ❑ **Timely data transfer**: Push data to core before needed
- ❑ **Can marshal any arbitrary sequence of lines**: Identifies generators, not patterns
- ❑ **Low hardware cost**: Profiler marks generators, no need for hardware to find them

■ Disadvantages

- ❑ **Requires profiler and ISA support**
- ❑ **Not always accurate (generator set is conservative)**: Pollution at remote core, wasted bandwidth on interconnect
 - Not a large problem as number of inter-segment blocks is small

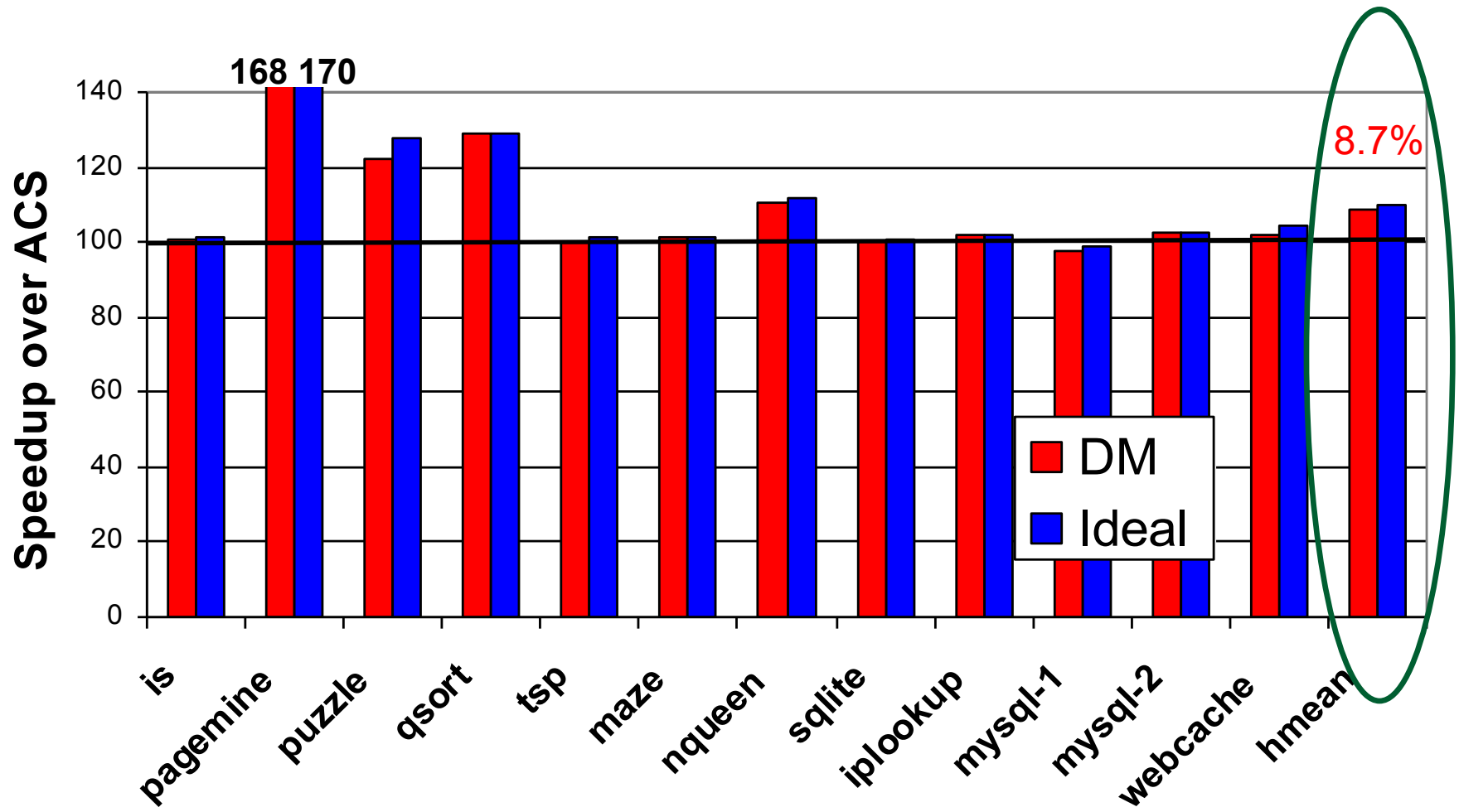
Accelerated Critical Sections with DM



Accelerated Critical Sections: Methodology

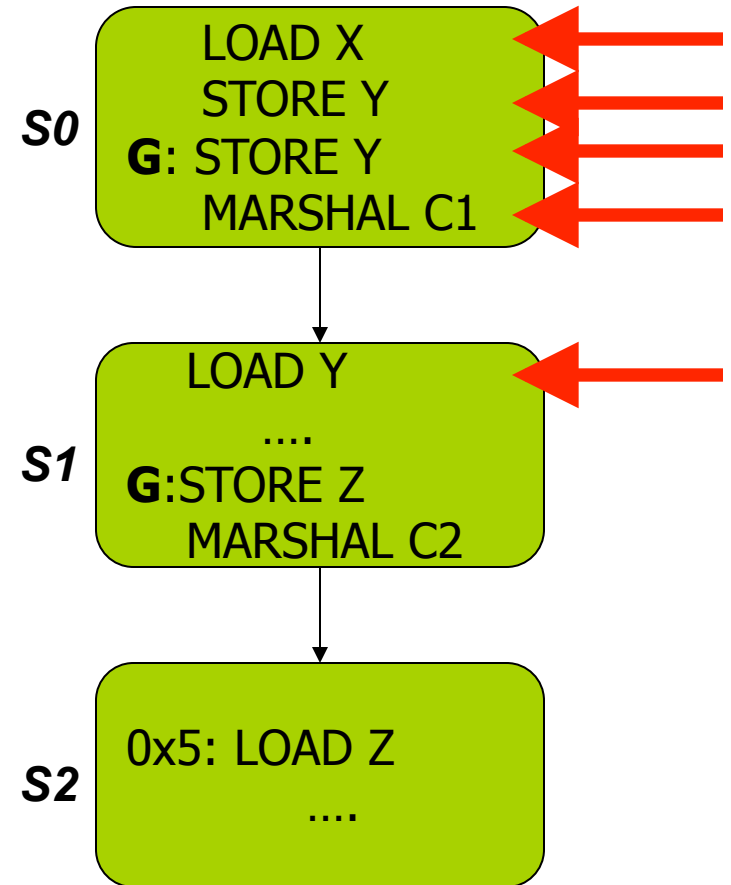
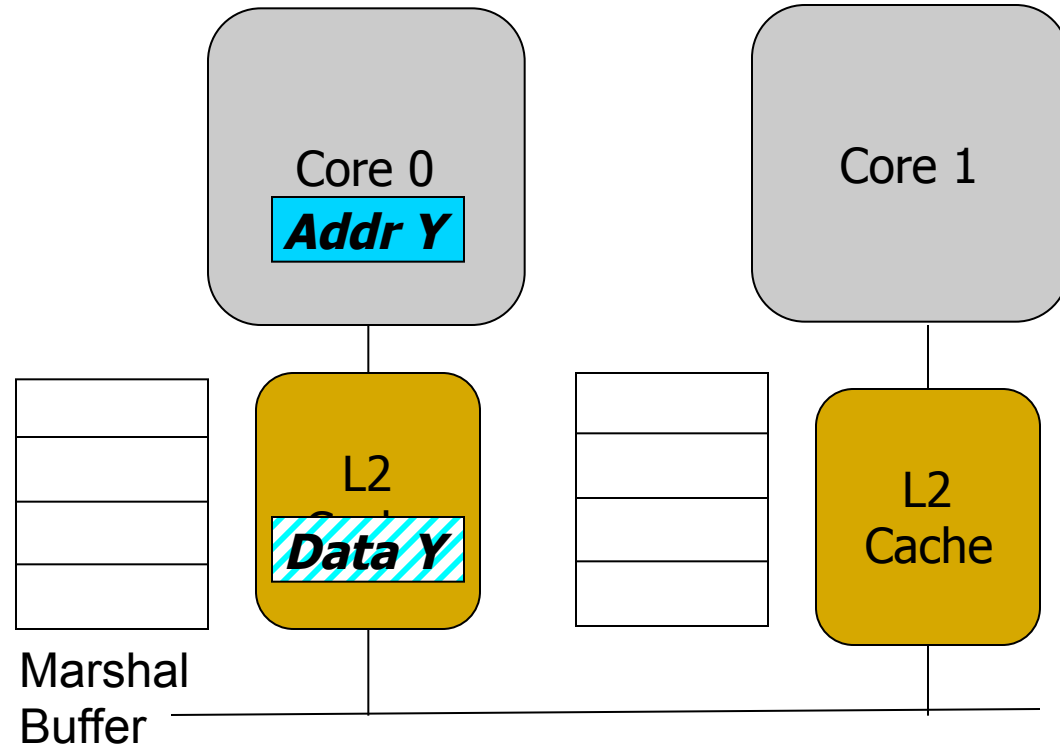
- Workloads: 12 critical section intensive applications
 - Data mining kernels, sorting, database, web, networking
 - Different training and simulation input sets
- Multi-core x86 simulator
 - 1 large and 28 small cores
 - Aggressive stream prefetcher employed at each core
- Details:
 - Large core: 2GHz, out-of-order, 128-entry ROB, 4-wide, 12-stage
 - Small core: 2GHz, in-order, 2-wide, 5-stage
 - Private 32 KB L1, private 256KB L2, 8MB shared L3
 - On-chip interconnect: Bi-directional ring, 5-cycle hop latency

DM on Accelerated Critical Sections: Results



Pipeline Parallelism

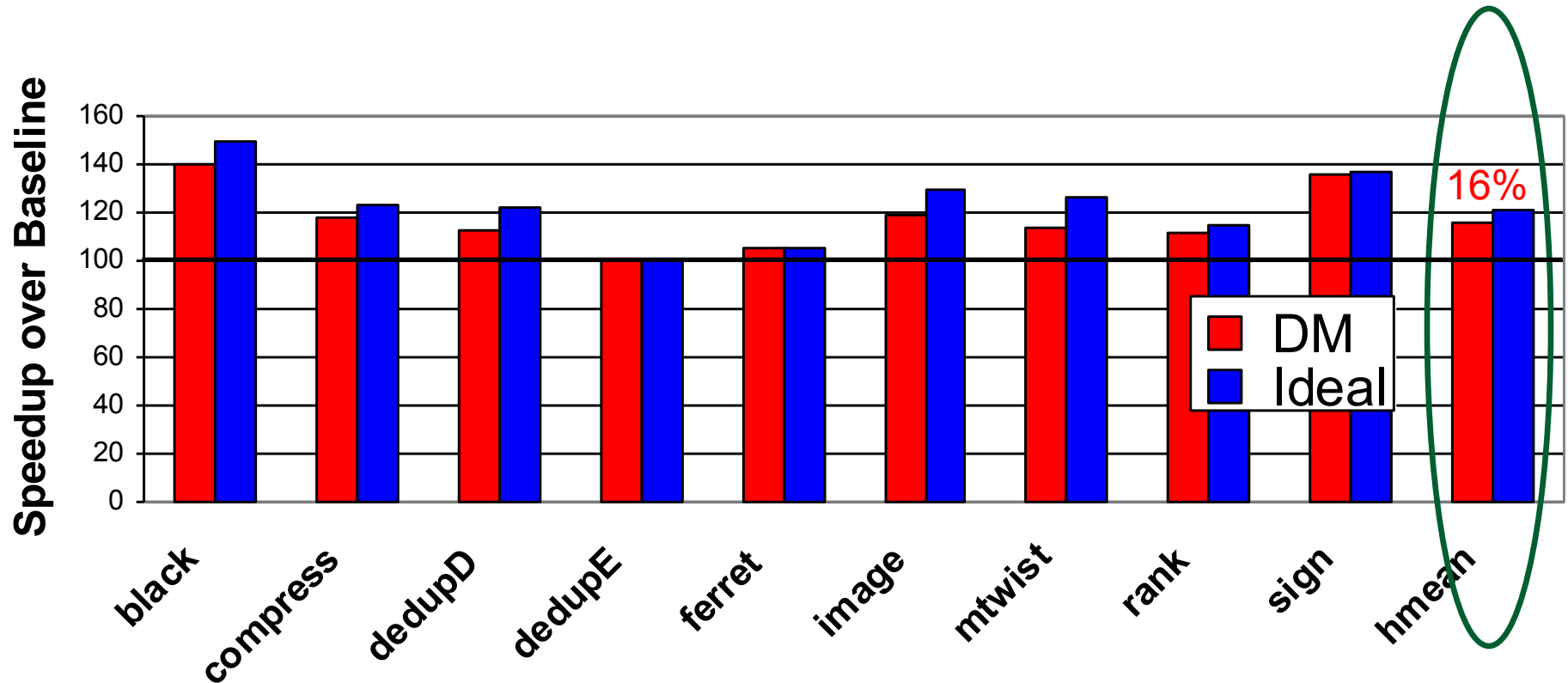
Cache Hit!



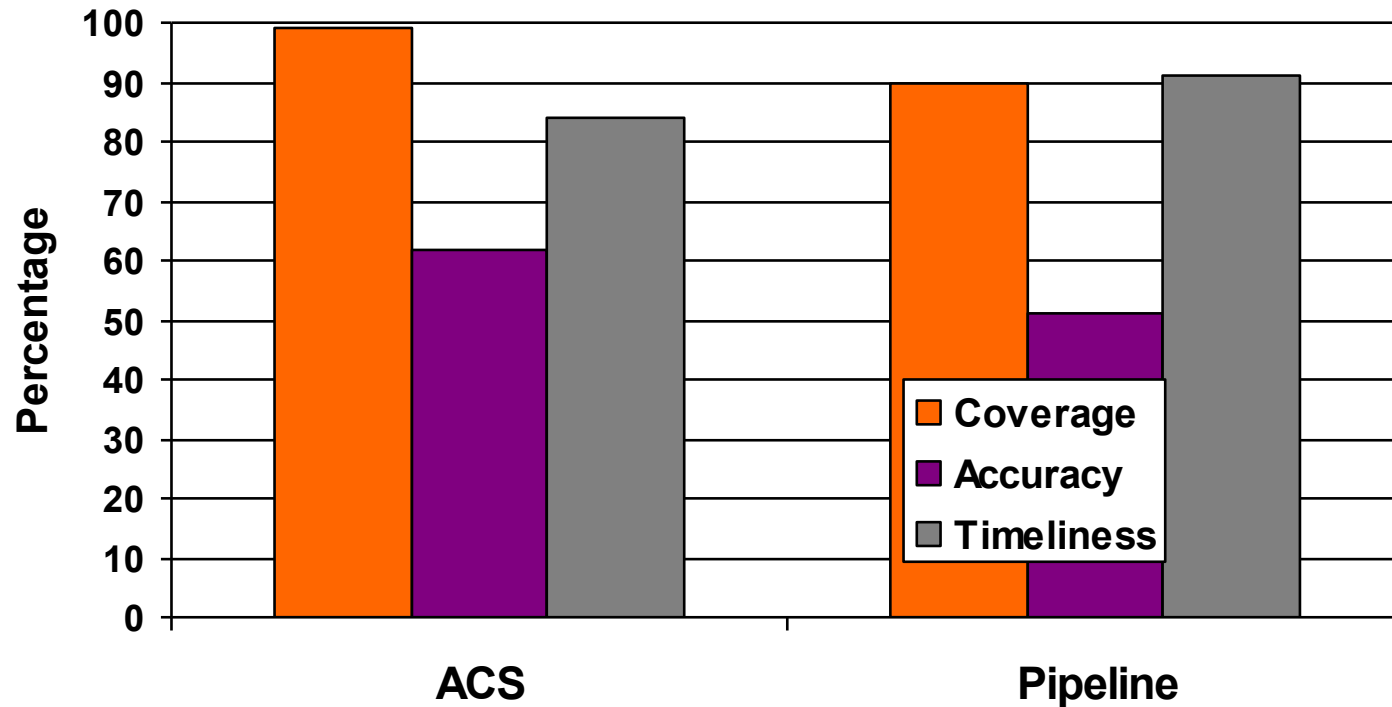
Pipeline Parallelism: Methodology

- Workloads: 9 applications with pipeline parallelism
 - Financial, compression, multimedia, encoding/decoding
 - Different training and simulation input sets
- Multi-core x86 simulator
 - 32-core CMP: 2GHz, in-order, 2-wide, 5-stage
 - Aggressive stream prefetcher employed at each core
 - Private 32 KB L1, private 256KB L2, 8MB shared L3
 - On-chip interconnect: Bi-directional ring, 5-cycle hop latency

DM on Pipeline Parallelism: Results



DM Coverage, Accuracy, Timeliness



- High coverage of inter-segment misses in a timely manner
- Medium accuracy does not impact performance
 - Only 5.0 and 6.8 cache blocks marshaled for average segment

Scaling Results

- DM performance improvement increases with
 - More cores
 - Higher interconnect latency
 - Larger private L2 caches
- Why? Inter-segment data misses become a larger bottleneck
 - More cores → More communication
 - Higher latency → Longer stalls due to communication
 - Larger L2 cache → Communication misses remain

Other Applications of Data Marshaling

- Can be applied to other Staged Execution models
 - Task parallelism models
 - Cilk, Intel TBB, Apple Grand Central Dispatch
 - Special-purpose remote functional units
 - Computation spreading [Chakraborty et al., ASPLOS' 06]
 - Thread motion/migration [e.g., Rangan et al., ISCA' 09]
- Can be an enabler for more aggressive SE models
 - Lowers the cost of data migration
 - an important overhead in remote execution of code segments
 - Remote execution of finer-grained tasks can become more feasible → finer-grained parallelization in multi-cores

Data Marshaling Summary

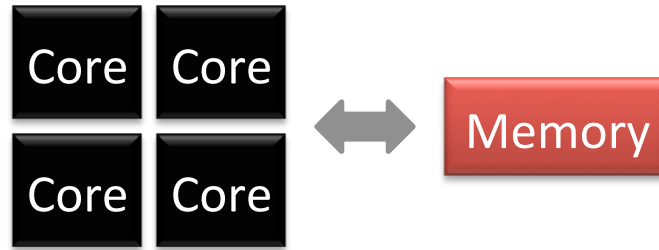
- **Inter-segment data transfers between cores** limit the benefit of promising Staged Execution (SE) models
- Data Marshaling is a hardware/software cooperative solution: **detect inter-segment data generator instructions and push their data to next segment's core**
 - Significantly reduces cache misses for inter-segment data
 - Low cost, high-coverage, timely for arbitrary address sequences
 - Achieves most of the potential of eliminating such misses
- Applicable to several existing Staged Execution models
 - Accelerated Critical Sections: 9% performance benefit
 - Pipeline Parallelism: 16% performance benefit
- Can enable new models → **very fine-grained remote execution**

Talk Outline

- Problem and Motivation
- How Do We Get There: Examples
- Accelerated Critical Sections (ACS)
- Bottleneck Identification and Scheduling (BIS)
- Staged Execution and Data Marshaling
- Thread Cluster Memory Scheduling (if time permits)
- Ongoing/Future Work
- Conclusions

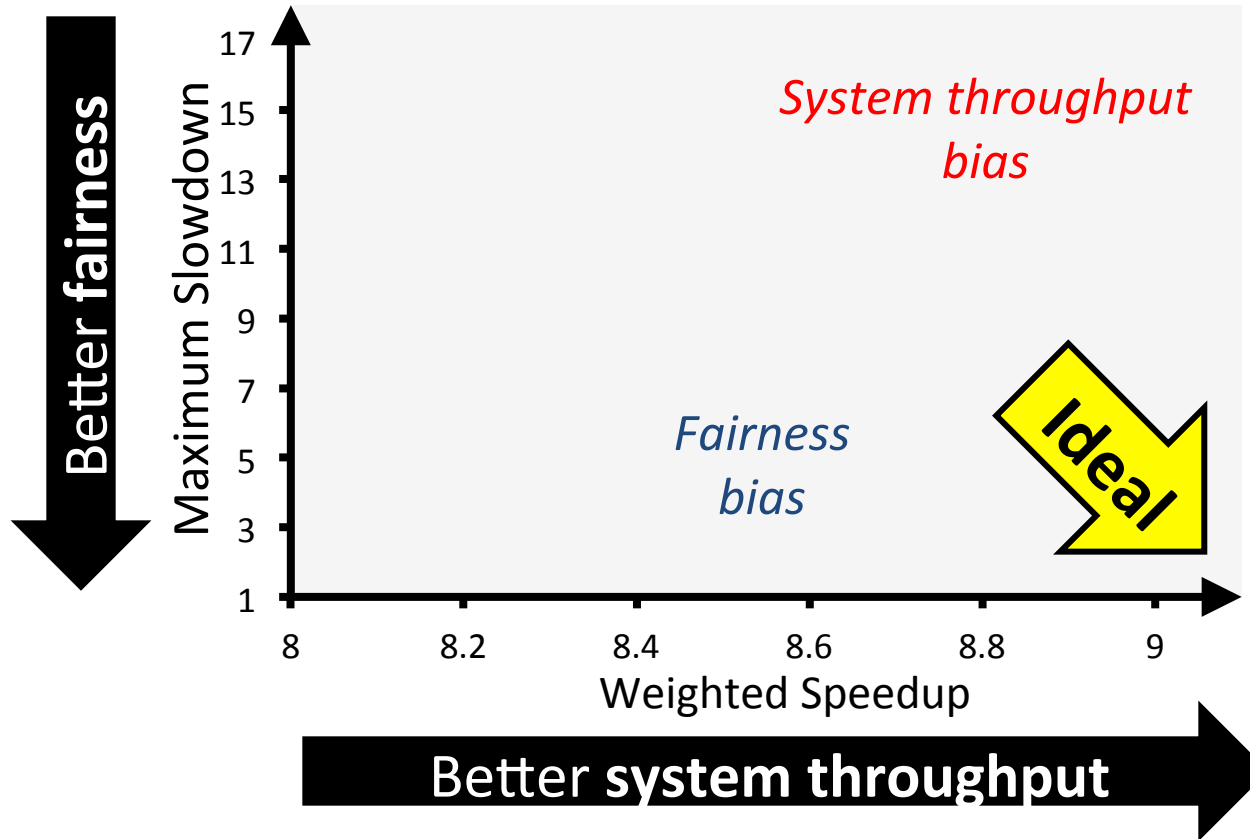
Motivation

- Memory is a shared resource



- Threads' requests contend for memory
 - Degradation in single thread performance
 - Can even lead to starvation
- How to schedule memory requests to increase both system throughput and fairness?

Previous Scheduling Algorithms are Biased



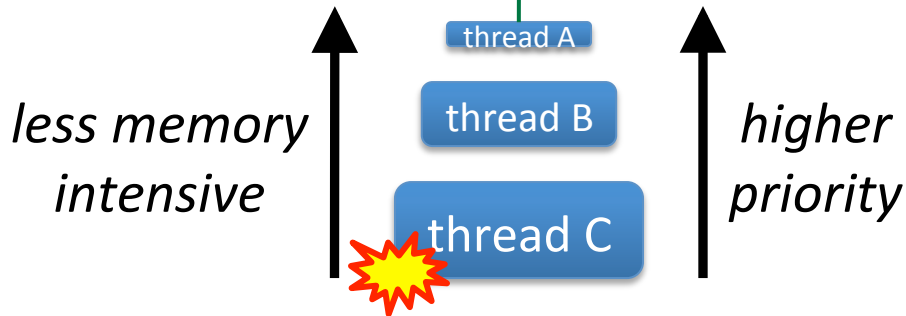
No previous memory scheduling algorithm provides both the best fairness and system throughput

Why do Previous Algorithms Fail?

Throughput biased approach

Prioritize less memory-intensive threads

Good for throughput



starvation → unfairness

Fairness biased approach

Take turns accessing memory

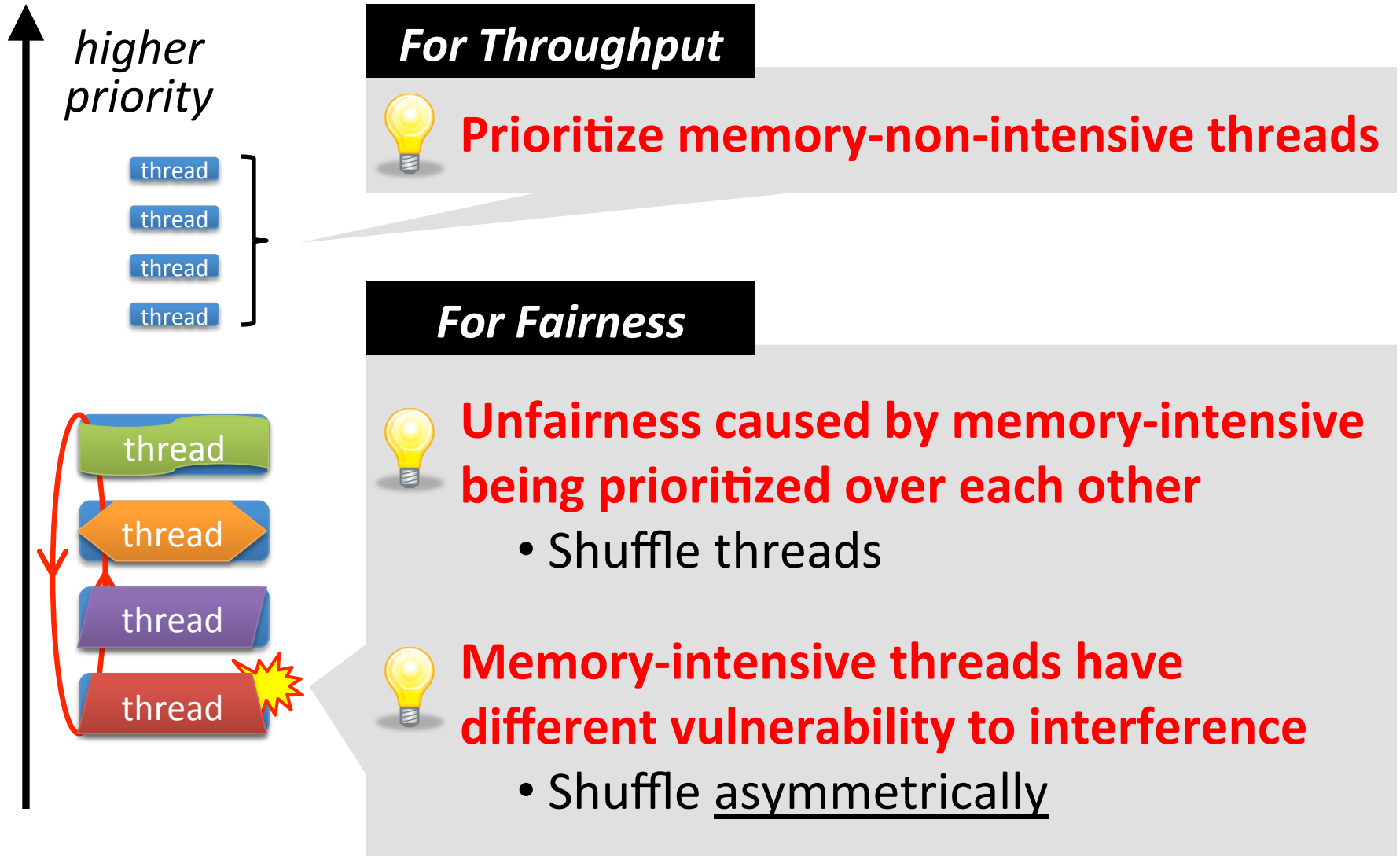
Does not starve



**not prioritized →
reduced throughput**

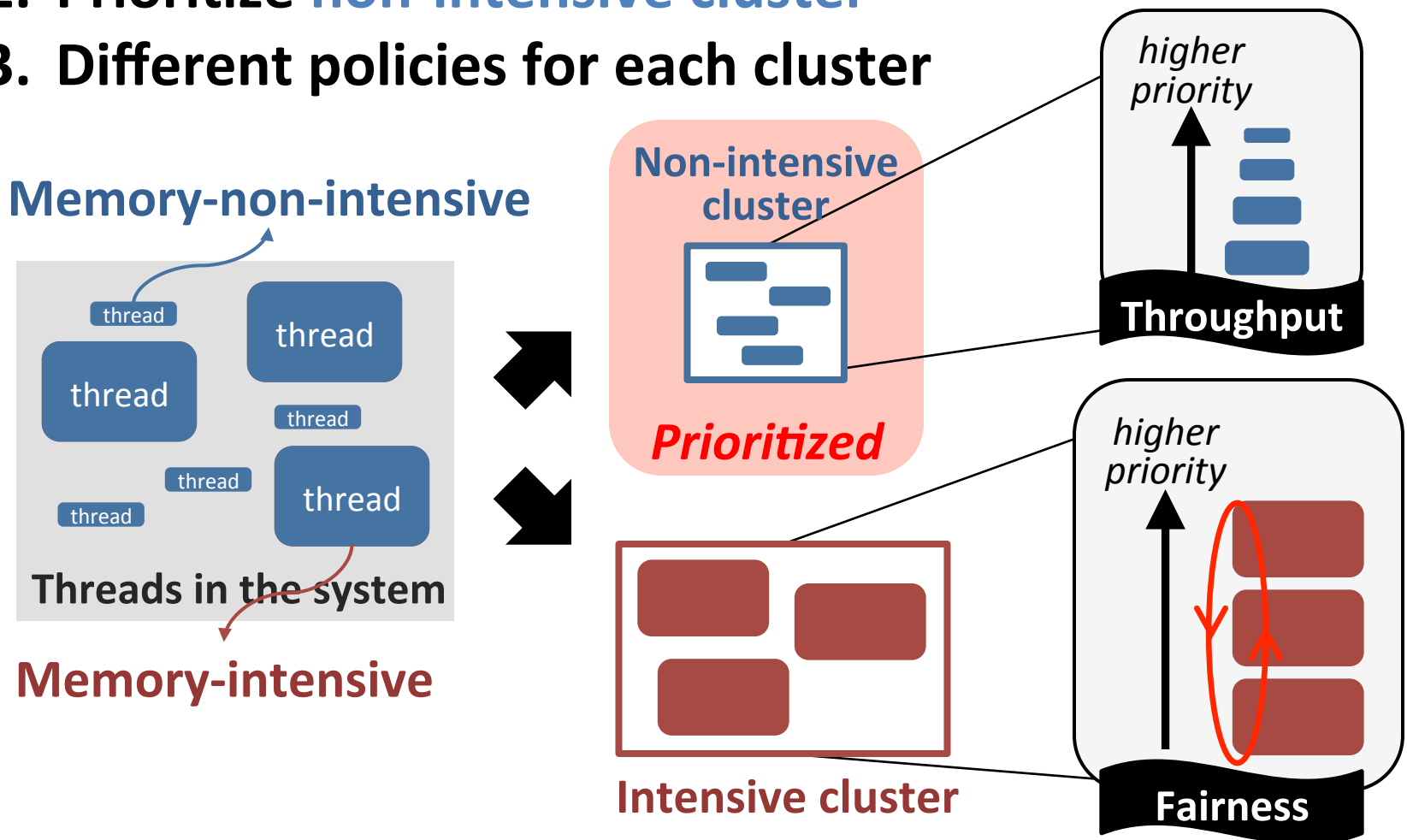
Single policy for all threads is insufficient

Insight: Achieving Best of Both Worlds



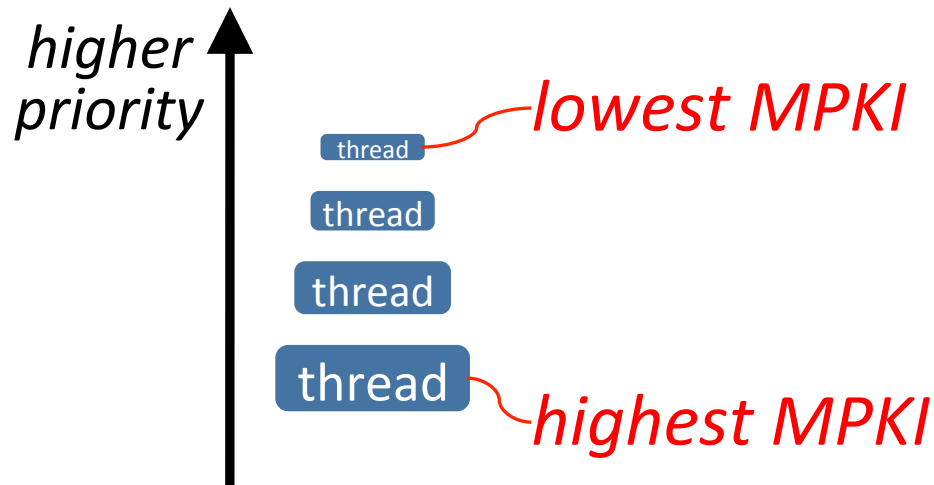
Overview: Thread Cluster Memory Scheduling

1. Group threads into two *clusters*
2. Prioritize *non-intensive cluster*
3. Different policies for each cluster



Non-Intensive Cluster

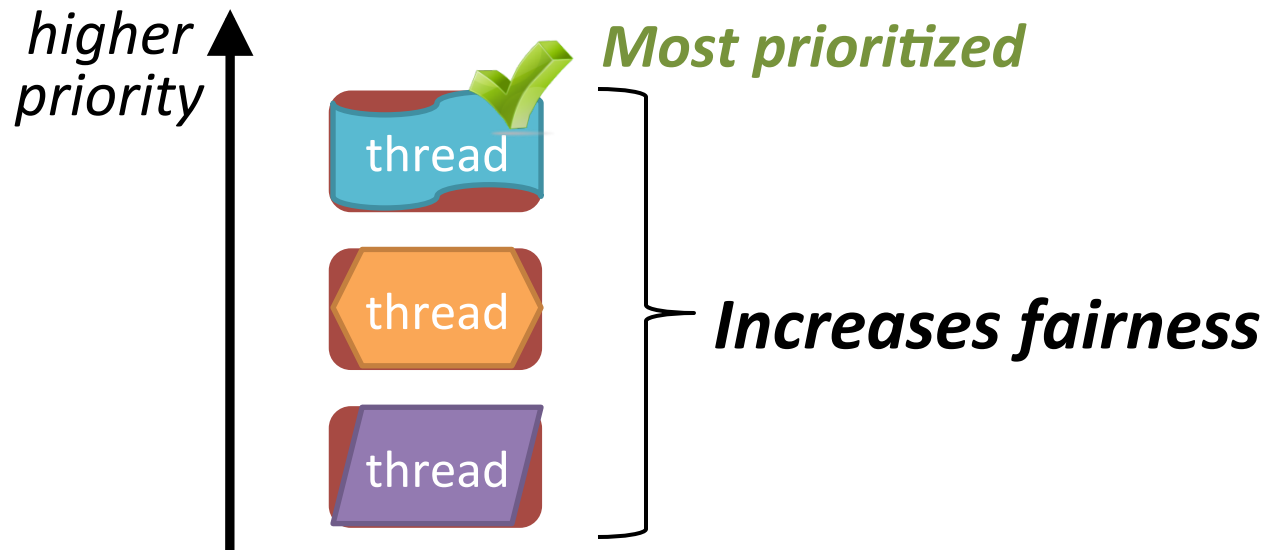
Prioritize threads according to MPKI



- **Increases system throughput**
 - Least intensive thread has the greatest potential for making progress in the processor

Intensive Cluster

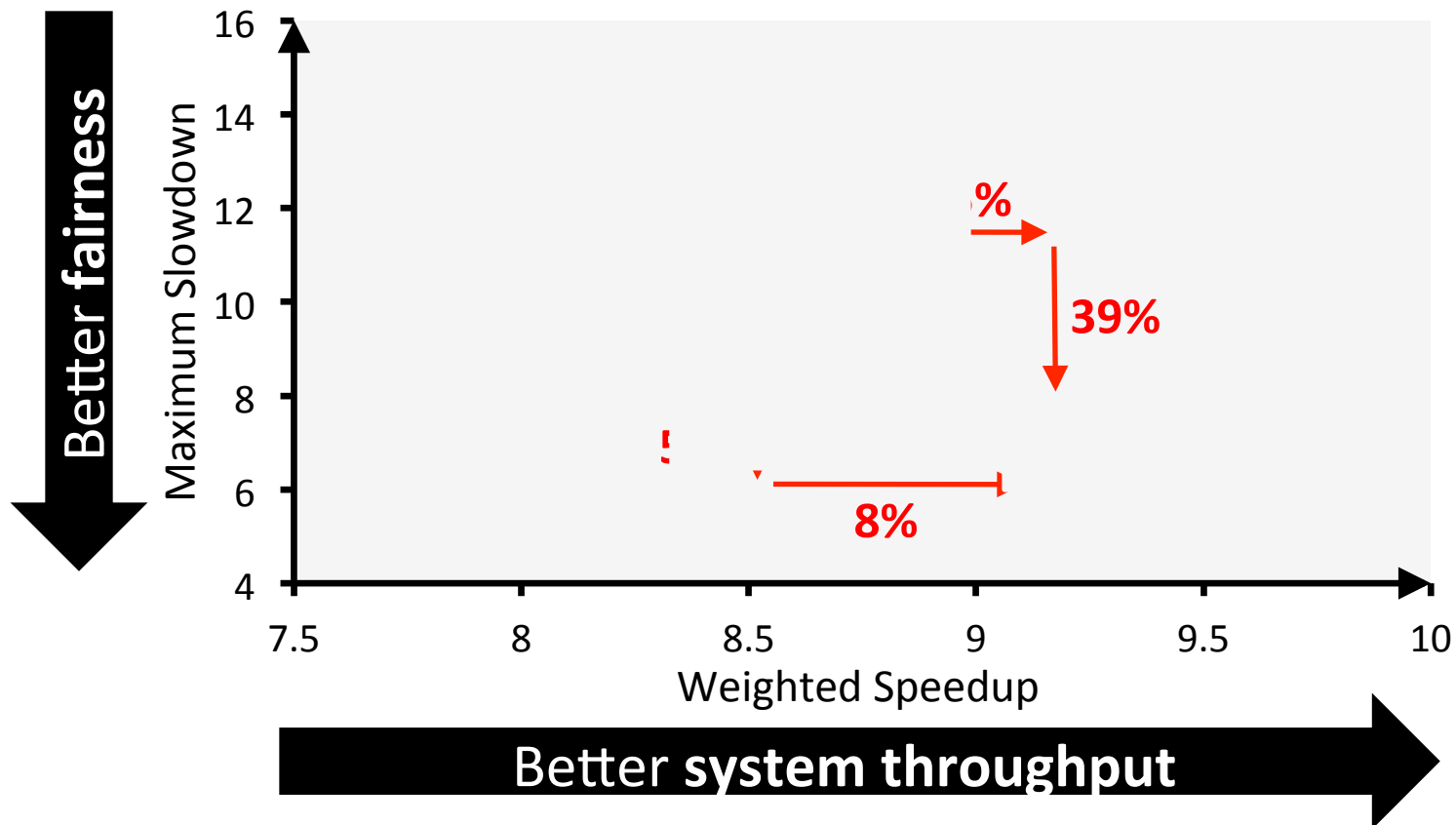
Periodically shuffle the priority of threads



- Is treating all threads equally good enough?
- ***BUT: Equal turns \neq Same slowdown***

Results: Fairness vs. Throughput

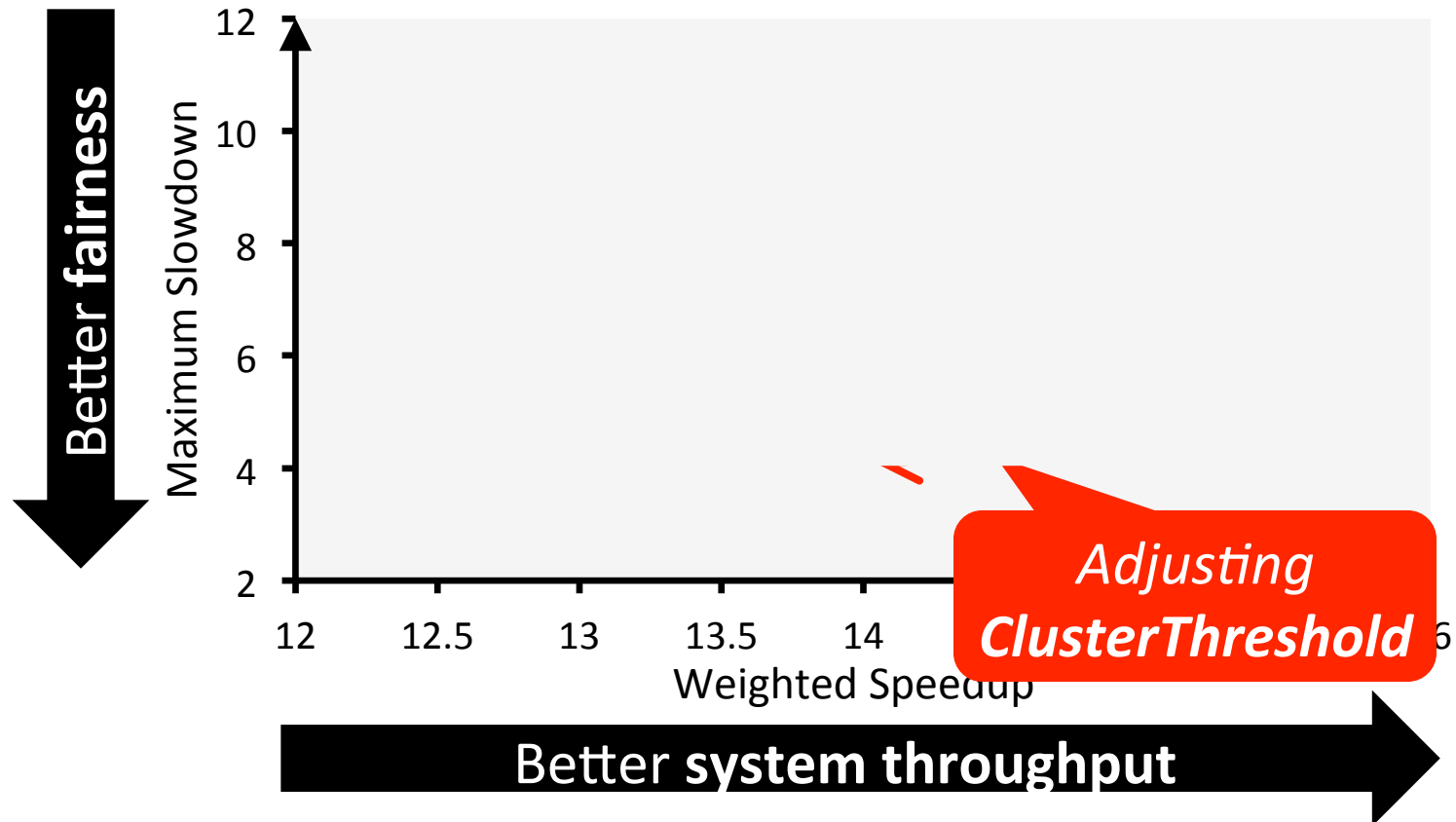
Averaged over 96 workloads



TCM provides best fairness and system throughput

Results: Fairness-Throughput Tradeoff

When configuration parameter is varied...



TCM allows robust fairness-throughput tradeoff

TCM Summary

- No previous memory scheduling algorithm provides both high *system throughput* and *fairness*
 - **Problem:** They use a single policy for all threads
- TCM is a heterogeneous scheduling policy
 1. Prioritize *non-intensive* cluster → throughput
 2. Shuffle priorities in *intensive* cluster → fairness
 3. Shuffling should favor *nice* threads → fairness
- *Heterogeneity in memory scheduling provides the best system throughput and fairness*

More Details on TCM

- Kim et al., “Thread Cluster Memory Scheduling: Exploiting Differences in Memory Access Behavior,” MICRO 2010, Top Picks 2011.

Memory Control in CPU-GPU Systems

- **Observation:** Heterogeneous CPU-GPU systems require memory schedulers with **large request buffers**
- **Problem:** Existing monolithic application-aware memory scheduler designs are **hard to scale** to large request buffer sizes
- **Solution:** Staged Memory Scheduling (SMS)
decomposes the memory controller into three simple stages:
 - 1) Batch formation: maintains row buffer locality
 - 2) Batch scheduler: reduces interference between applications
 - 3) DRAM command scheduler: issues requests to DRAM
- Compared to state-of-the-art memory schedulers:
 - ❑ SMS is significantly simpler and more scalable
 - ❑ SMS provides higher performance and fairness

Asymmetric Memory QoS in a Parallel Application

- Threads in a multithreaded application are inter-dependent
- Some threads can be on the critical path of execution due to synchronization; some threads are not
- How do we schedule requests of inter-dependent threads to maximize multithreaded application performance?
- Idea: **Estimate limiter threads** likely to be on the critical path and prioritize their requests; **shuffle priorities of non-limiter threads** to reduce memory interference among them [Ebrahimi+, MICRO'11]
- Hardware/software cooperative limiter thread estimation:
 - Thread executing the most contended critical section
 - Thread that is falling behind the most in a *parallel for* loop

Talk Outline

- Problem and Motivation
- How Do We Get There: Examples
- Accelerated Critical Sections (ACS)
- Bottleneck Identification and Scheduling (BIS)
- Staged Execution and Data Marshaling
- Thread Cluster Memory Scheduling (if time permits)
- Ongoing/Future Work
- Conclusions

Related Ongoing/Future Work

- Dynamically asymmetric cores
- Memory system design for asymmetric cores
- Asymmetric memory systems
 - Phase Change Memory (or Technology X) + DRAM
 - Hierarchies optimized for different access patterns
- Asymmetric on-chip interconnects
 - Interconnects optimized for different application requirements
- Asymmetric resource management algorithms
 - E.g., network congestion control
- Interaction of multiprogrammed multithreaded workloads

Talk Outline

- Problem and Motivation
- How Do We Get There: Examples
- Accelerated Critical Sections (ACS)
- Bottleneck Identification and Scheduling (BIS)
- Staged Execution and Data Marshaling
- Thread Cluster Memory Scheduling (if time permits)
- Ongoing/Future Work
- Conclusions

Summary

- Applications and phases have varying performance requirements
- Designs evaluated on multiple metrics/constraints: energy, performance, reliability, fairness, ...
- **One-size-fits-all** design cannot satisfy all requirements and metrics: **cannot get the best of all worlds**
- **Asymmetry** in design enables tradeoffs: **can get the best of all worlds**
 - Asymmetry in core microarch. → **Accelerated Critical Sections, BIS, DM**
→ Good parallel performance + Good serialized performance
 - Asymmetry in memory scheduling → **Thread Cluster Memory Scheduling**
→ Good throughput + good fairness
- Simple asymmetric designs can be effective and low-cost

Thank You

Onur Mutlu

onur@cmu.edu

<http://www.ece.cmu.edu/~omutlu>

Email me with any questions and feedback!

Architecting and Exploiting Asymmetry in Multi-Core Architectures

Onur Mutlu

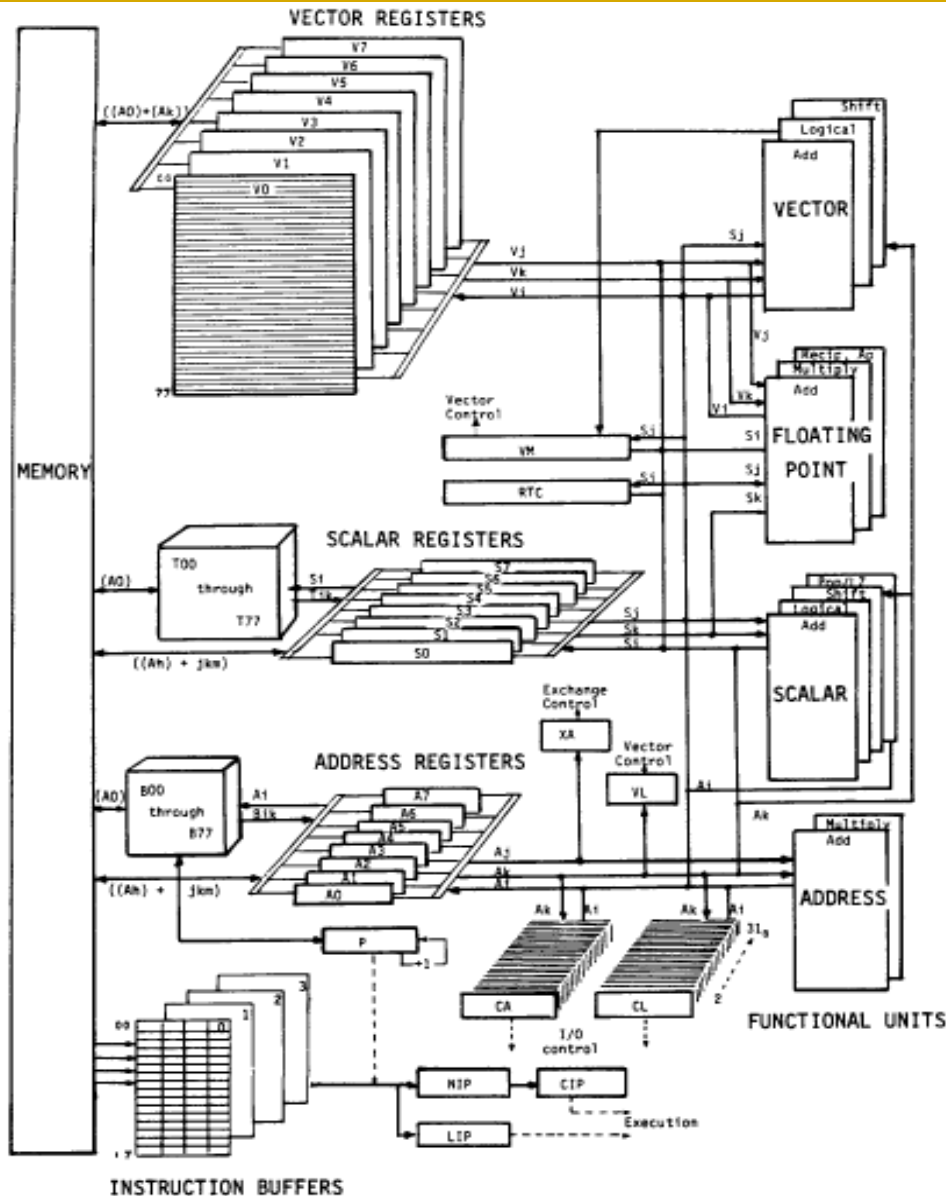
onur@cmu.edu

July 23, 2013

BSC/UPC

SAFARI Carnegie Mellon

Vector Machine Organization (CRAY-1)



- CRAY-1
- Russell, “The CRAY-1 computer system,” CACM 1978.
- Scalar and vector modes
- 8 64-element vector registers
- 64 bits per element
- 16 memory banks
- 8 64-bit scalar registers
- 8 24-bit address registers

Identifying and Accelerating Resource Contention Bottlenecks

Thread Serialization

- Three fundamental causes

1. Synchronization

2. Load imbalance

3. Resource contention

Memory Contention as a Bottleneck

■ Problem:

- ❑ Contended memory regions cause serialization of threads
- ❑ Threads accessing such regions can form the critical path
- ❑ Data-intensive workloads (MapReduce, GraphLab, Graph500) can be sped up by 1.5 to 4X by ideally removing contention

■ Idea:

- ❑ Identify contended regions dynamically
- ❑ Prioritize caching the data from threads which are slowed down the most by such regions in faster DRAM/eDRAM

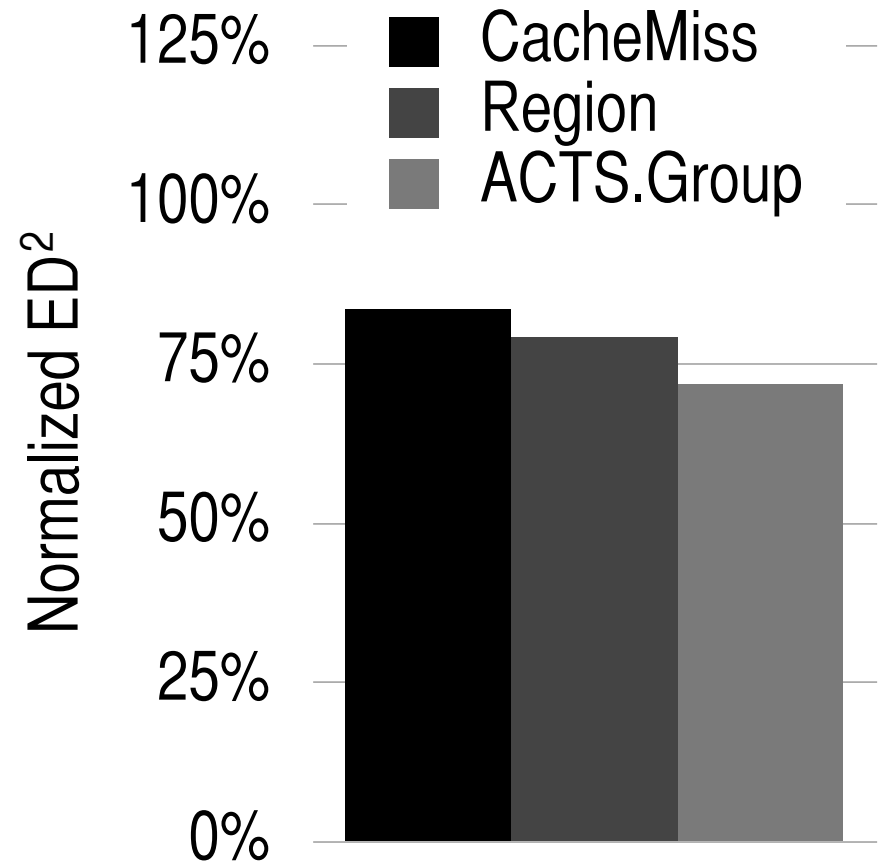
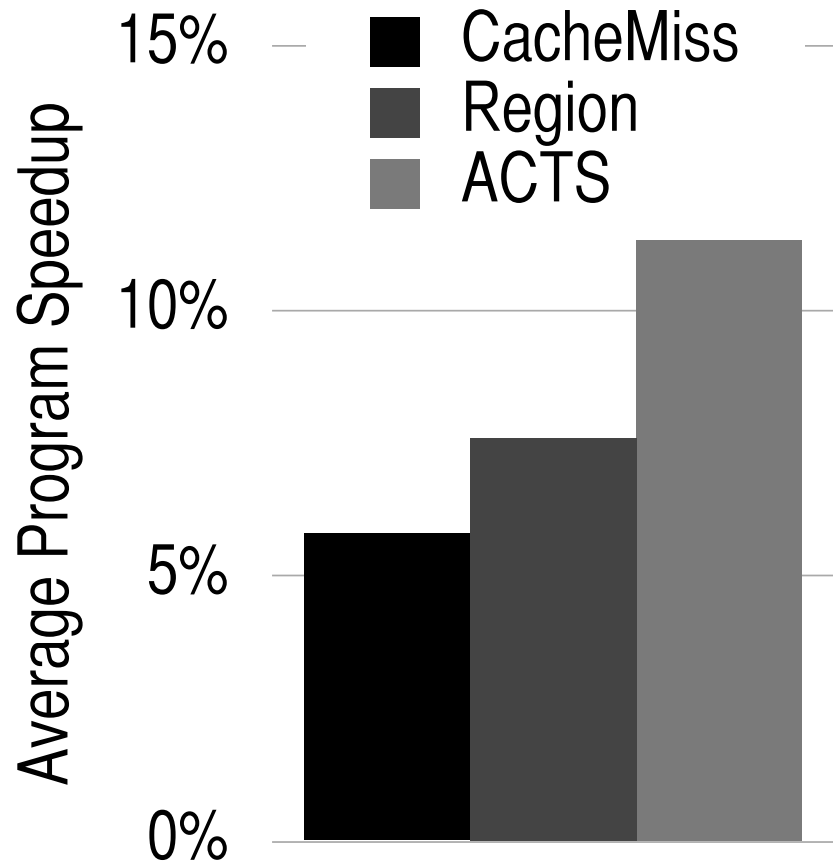
■ Benefits:

- ❑ Reduces contention, serialization, critical path

Evaluation

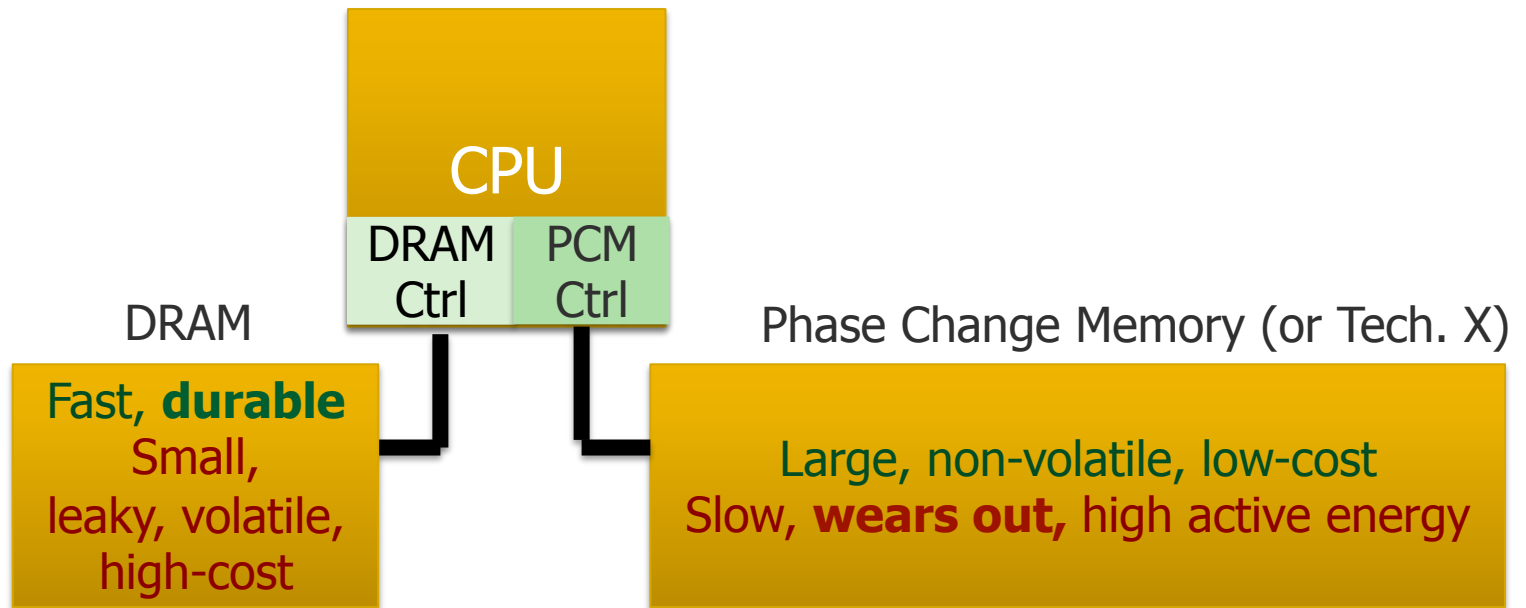
- Workloads: MapReduce, GraphLab, Graph500
- Cycle-level x86 platform simulator
 - **CPU**: 8 out-of-order cores, 32KB private L1, 512KB shared L2
 - **Hybrid Memory**: DDR3 1066 MT/s, 32MB DRAM, 8GB PCM
- Mechanisms
 - Baseline: DRAM as a conventional cache to PCM
 - CacheMiss: Prioritize caching data from threads with highest cache miss latency
 - Region: Cache data from most contended memory regions
 - **ACTS**: Prioritize caching data from threads most slowed down due to memory region contention

Caching Results



Heterogeneous Main Memory

Heterogeneous Memory Systems



Hardware/software manage data allocation and movement
to achieve the best of multiple technologies

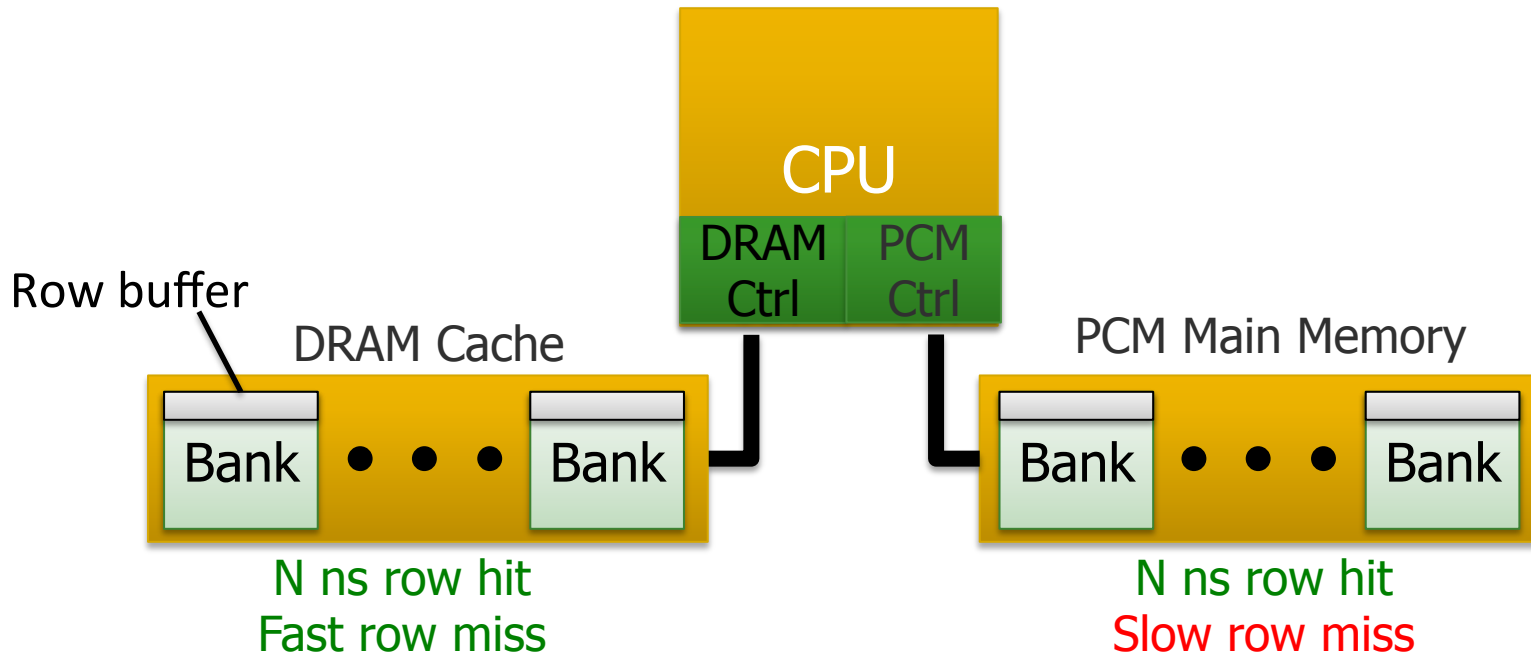
Meza, Chang, Yoon, Mutlu, Ranganathan, "Enabling Efficient and Scalable Hybrid Memories,"
IEEE Comp. Arch. Letters, 2012.

One Option: DRAM as a Cache for PCM

- PCM is main memory; DRAM caches memory rows/blocks
 - Benefits: Reduced latency on DRAM cache hit; write filtering
- Memory controller hardware manages the DRAM cache
 - Benefit: Eliminates system software overhead
- Three issues:
 - What data should be placed in DRAM versus kept in PCM?
 - What is the granularity of data movement?
 - How to design a low-cost hardware-managed DRAM cache?
- Two idea directions:
 - Locality-aware data placement [Yoon+ , CMU TR 2011]
 - Cheap tag stores and dynamic granularity [Meza+, IEEE CAL 2012]

DRAM vs. PCM: An Observation

- Row buffers are the same in DRAM and PCM
- Row buffer **hit** latency **same** in DRAM and PCM
- Row buffer **miss** latency **small** in DRAM, **large** in PCM



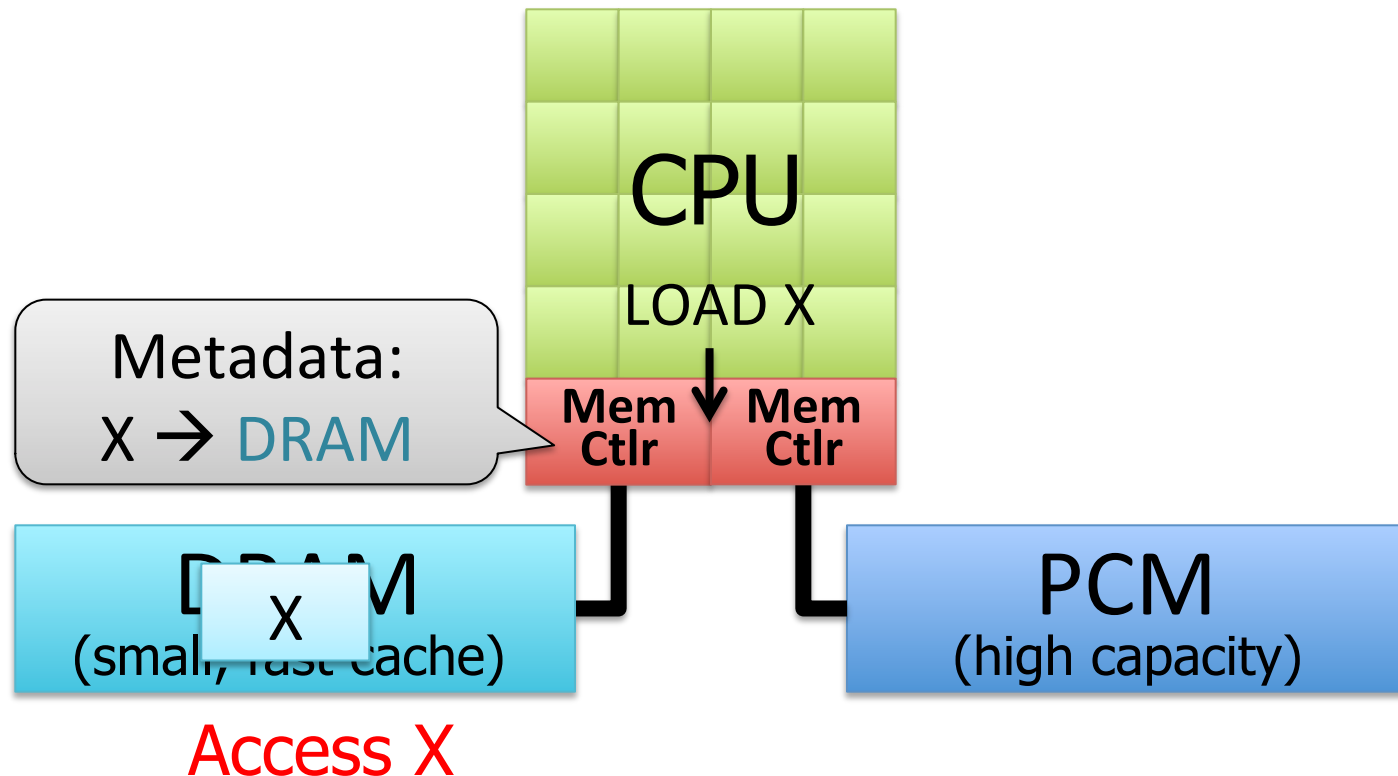
- Accessing the row buffer in PCM is fast
- What incurs high latency is the PCM array access → avoid this

Row-Locality-Aware Data Placement

- Idea: Cache in DRAM only those rows that
 - Frequently cause row buffer conflicts → because row-conflict latency is smaller in DRAM
 - Are reused many times → to reduce cache pollution and bandwidth waste
- Simplified rule of thumb:
 - Streaming accesses: Better to place in PCM
 - Other accesses (with some reuse): Better to place in DRAM
- Bridges half of the performance gap between all-DRAM and all-PCM memory on memory-intensive workloads
- Yoon et al., “Row Buffer Locality-Aware Data Placement in Hybrid Memories,” CMU SAFARI Technical Report, 2011.

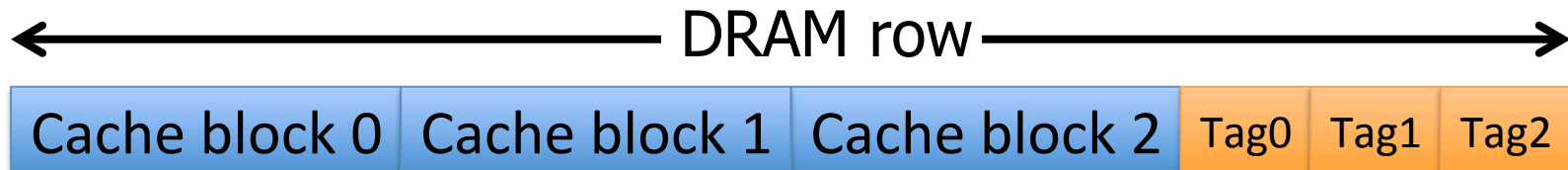
The Problem with Large DRAM Caches

- A large DRAM cache requires a large metadata (tag + block-based information) store
- How do we design an efficient DRAM cache?



Idea 1: Tags in Memory

- Store tags in the same row as data in DRAM
 - Store metadata in same row as their data
 - Data and metadata can be accessed together



- Benefit: No on-chip tag storage overhead
- Downsides:
 - Cache hit determined only after a DRAM access
 - Cache hit requires two DRAM accesses

Idea 2: Cache Tags in SRAM

- Recall Idea 1: Store all metadata in DRAM
 - To reduce metadata storage overhead
- Idea 2: Cache in on-chip SRAM frequently-accessed metadata
 - Cache only a small amount to keep SRAM size small

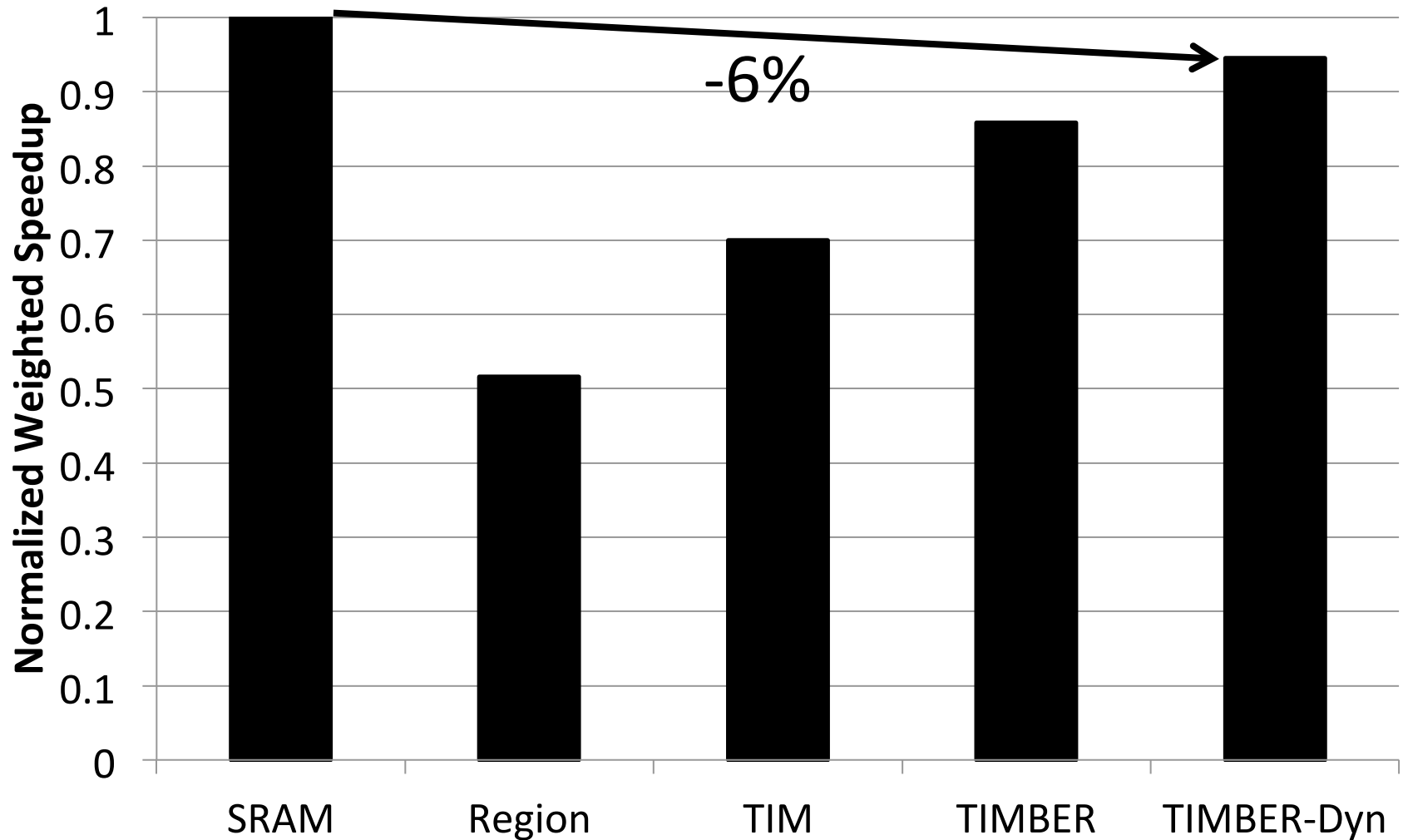
Idea 3: Dynamic Data Transfer Granularity

- Some applications benefit from caching more data
 - They have good spatial locality
- Others do not
 - Large granularity wastes bandwidth and reduces cache utilization
- Idea 3: Simple dynamic caching granularity policy
 - Cost-benefit analysis to determine best DRAM cache block size
 - Group main memory into sets of rows
 - Some row sets follow a fixed caching granularity
 - The rest of main memory follows the best granularity
 - Cost-benefit analysis: access latency versus number of cachings
 - Performed every quantum

Methodology

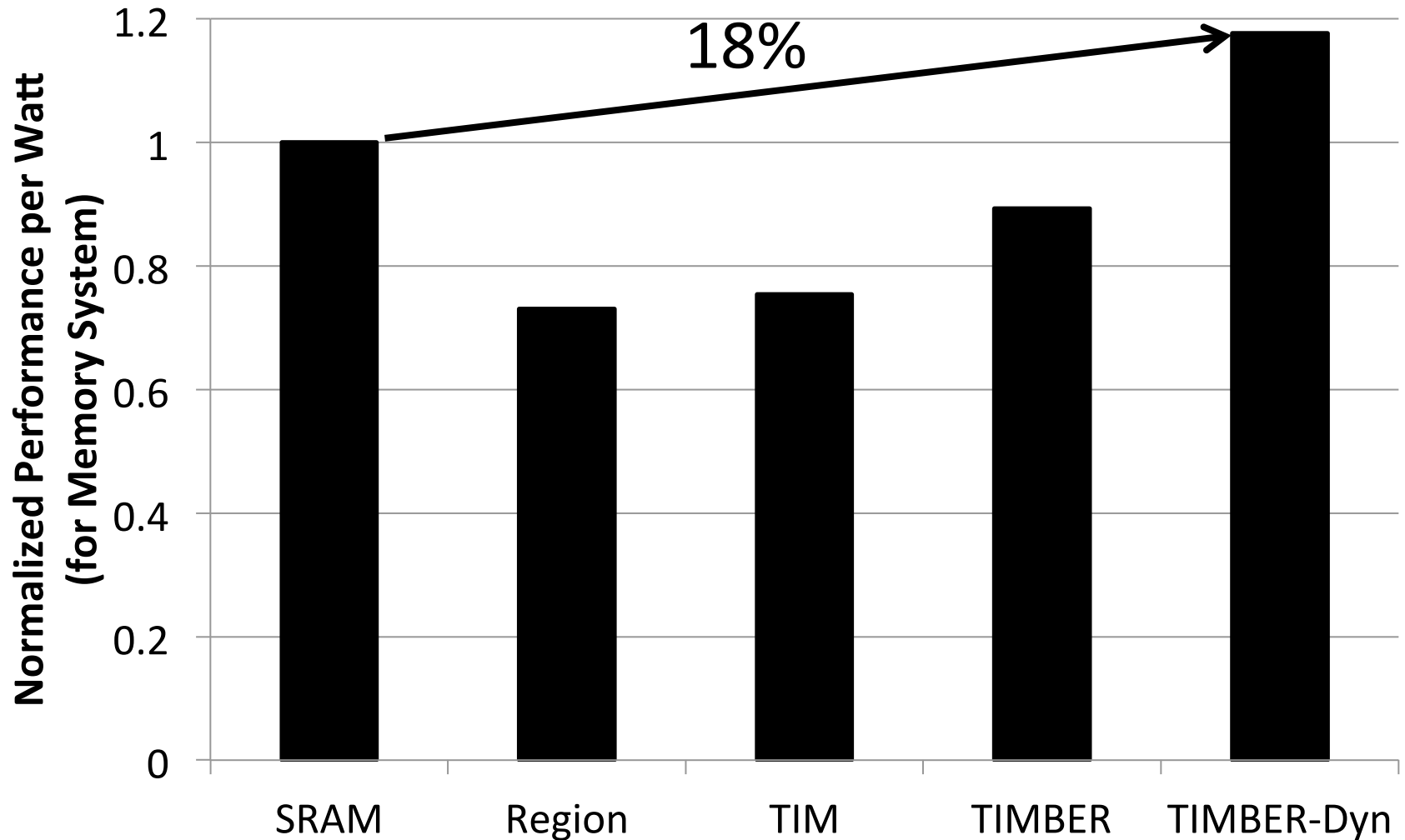
- System: 8 out-of-order cores at 4 GHz
- Memory: 512 MB direct-mapped DRAM, 8 GB PCM
 - 128B caching granularity
 - DRAM row hit (miss): 200 cycles (400 cycles)
 - PCM row hit (clean / dirty miss): 200 cycles (640 / 1840 cycles)
- Evaluated metadata storage techniques
 - All SRAM system (8MB of SRAM)
 - Region metadata storage
 - TIM metadata storage (same row as data)
 - TIMBER, 64-entry direct-mapped (8KB of SRAM)

TIMBER Performance



Meza, Chang, Yoon, Mutlu, Ranganathan, “[Enabling Efficient and Scalable Hybrid Memories](#),” IEEE Comp. Arch. Letters, 2012.

TIMBER Energy Efficiency



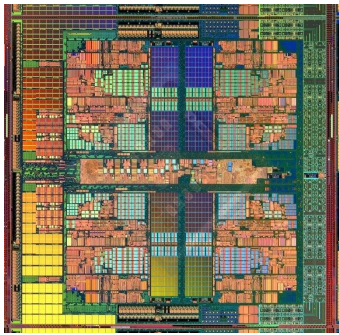
Summary

- Applications and phases have varying performance requirements
- Designs evaluated on multiple metrics/constraints: energy, performance, reliability, fairness, ...
- **One-size-fits-all** design cannot satisfy all requirements and metrics: **cannot get the best of all worlds**
- **Asymmetry** in design enables tradeoffs: **can get the best of all worlds**
 - Asymmetry in core microarch. → **Accelerated Critical Sections, BIS, DM**
→ Good parallel performance + Good serialized performance
 - Asymmetry in main memory → **Data Management for DRAM-PCM Hybrid Memory** → Good performance + good efficiency
- Simple asymmetric designs can be effective and low-cost

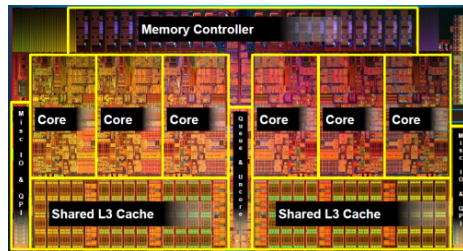
Memory QoS

Trend: Many Cores on Chip

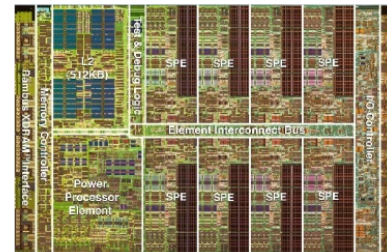
- Simpler and lower power than a single large core
- Large scale parallelism on chip



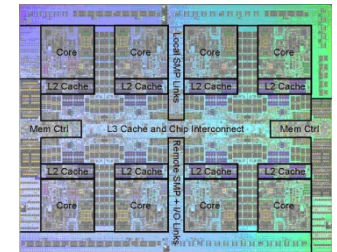
AMD Barcelona
4 cores



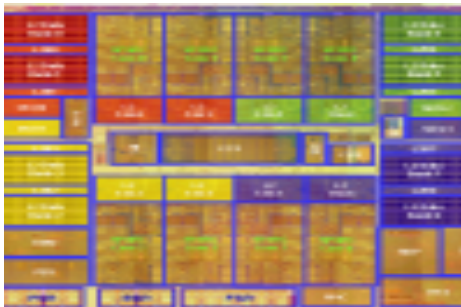
Intel Core i7
8 cores



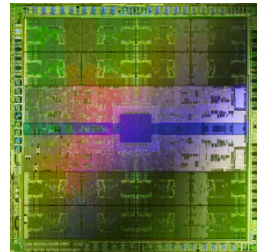
IBM Cell BE
8+1 cores



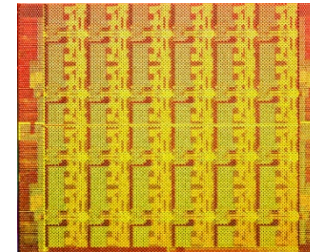
IBM POWER7
8 cores



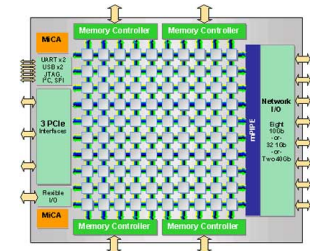
Sun Niagara II
8 cores



Nvidia Fermi
448 "cores"



Intel SCC
48 cores, networked

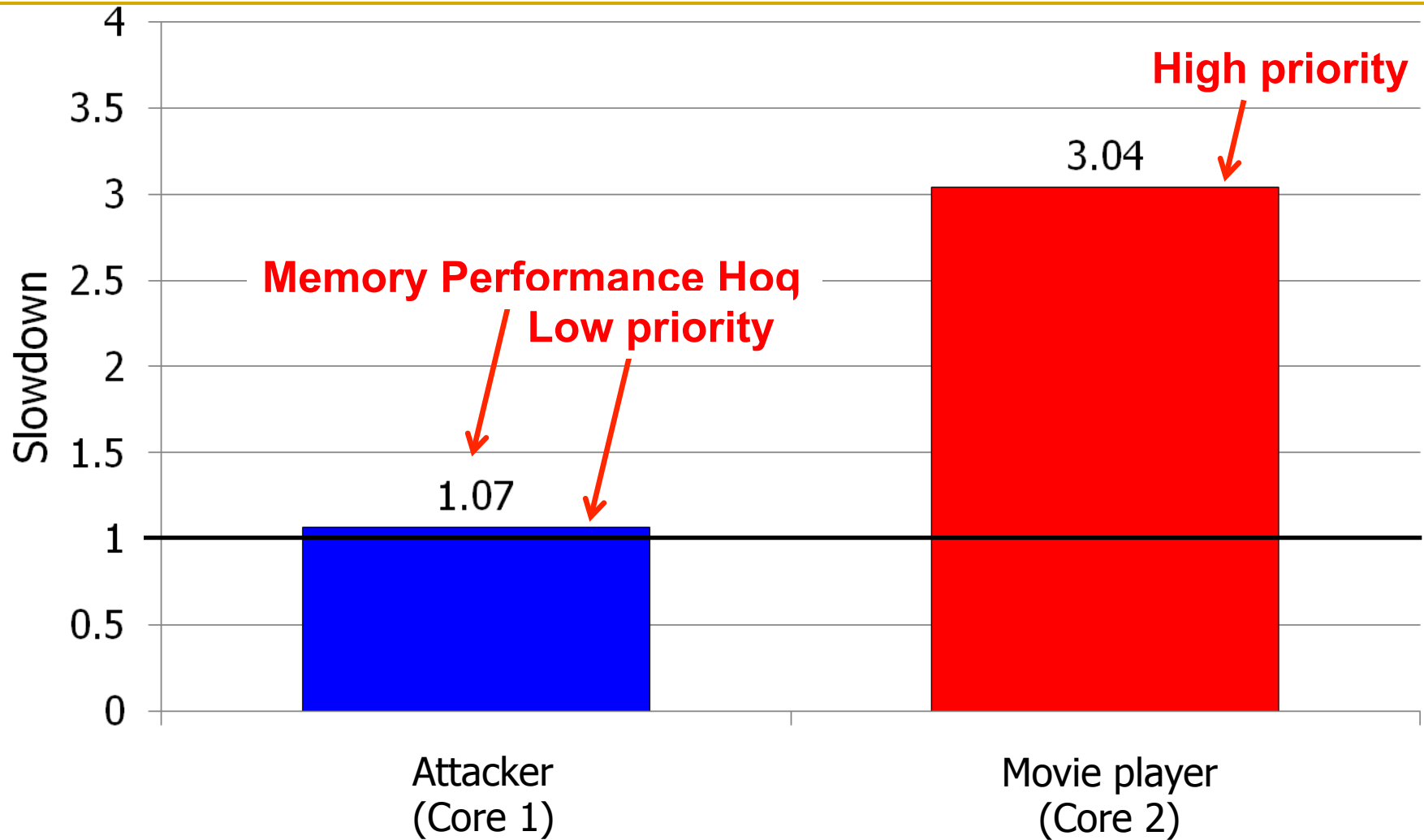


Tiler TILE Gx
100 cores, networked

Many Cores on Chip

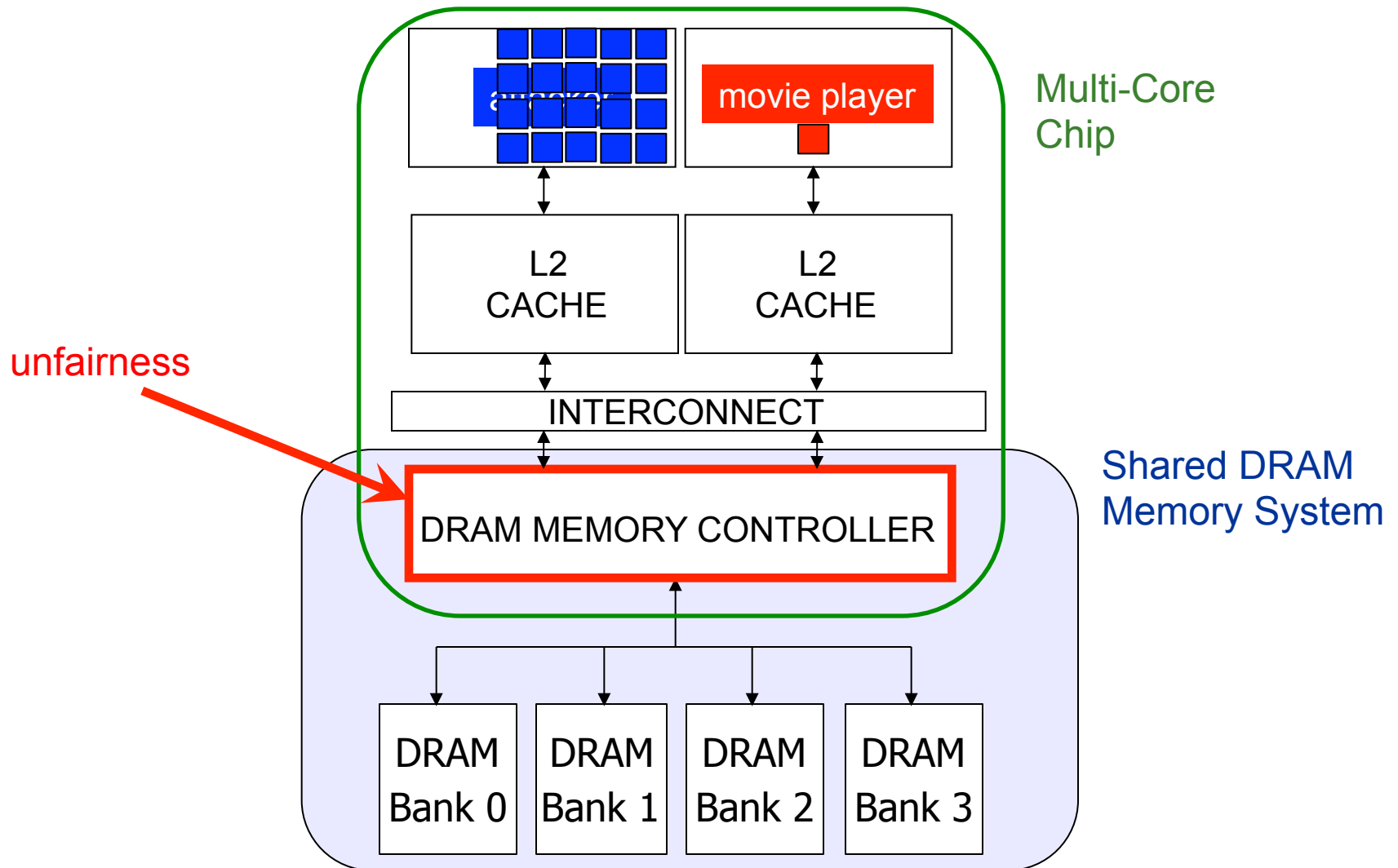
- What we want:
 - N times the system performance with N times the cores
- What do we get today?

(Un)expected Slowdowns



Moscibroda and Mutlu, “[Memory performance attacks: Denial of memory service in multi-core systems](#),” USENIX Security 2007.

Why? Uncontrolled Memory Interference



A Memory Performance Hog

```
// initialize large arrays A, B  
for (j=0; j<N; j++) {  
    index = j*linesize; streaming  
    A[index] = B[index];  
    ...  
}
```

STREAM

- Sequential memory access
- Very high row buffer locality (96% hit rate)
- Memory intensive

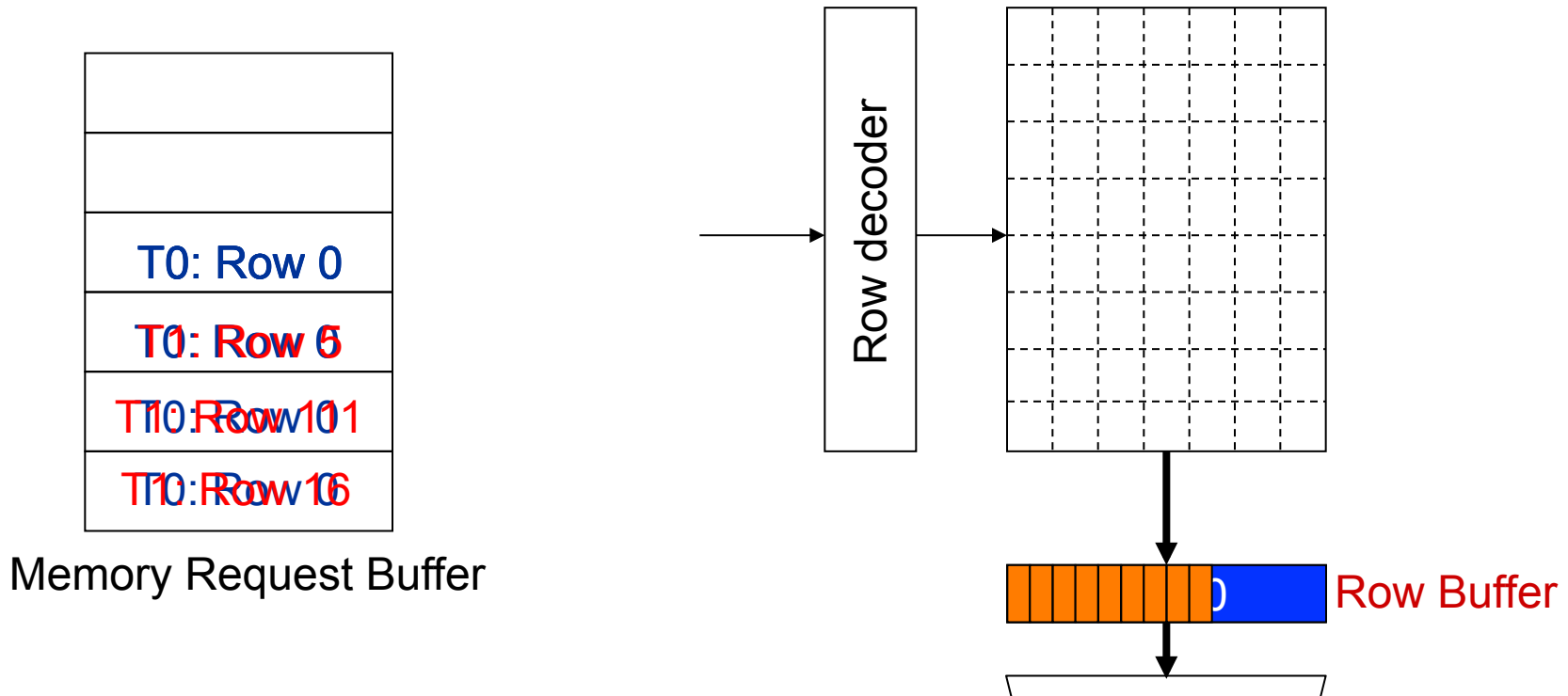
```
// initialize large arrays A, B  
for (j=0; j<N; j++) {  
    index = rand(); random  
    A[index] = B[index];  
    ...  
}
```

RANDOM

- Random memory access
- Very low row buffer locality (3% hit rate)
- Similarly memory intensive

Moscibroda and Mutlu, “[Memory Performance Attacks](#),” USENIX Security 2007.

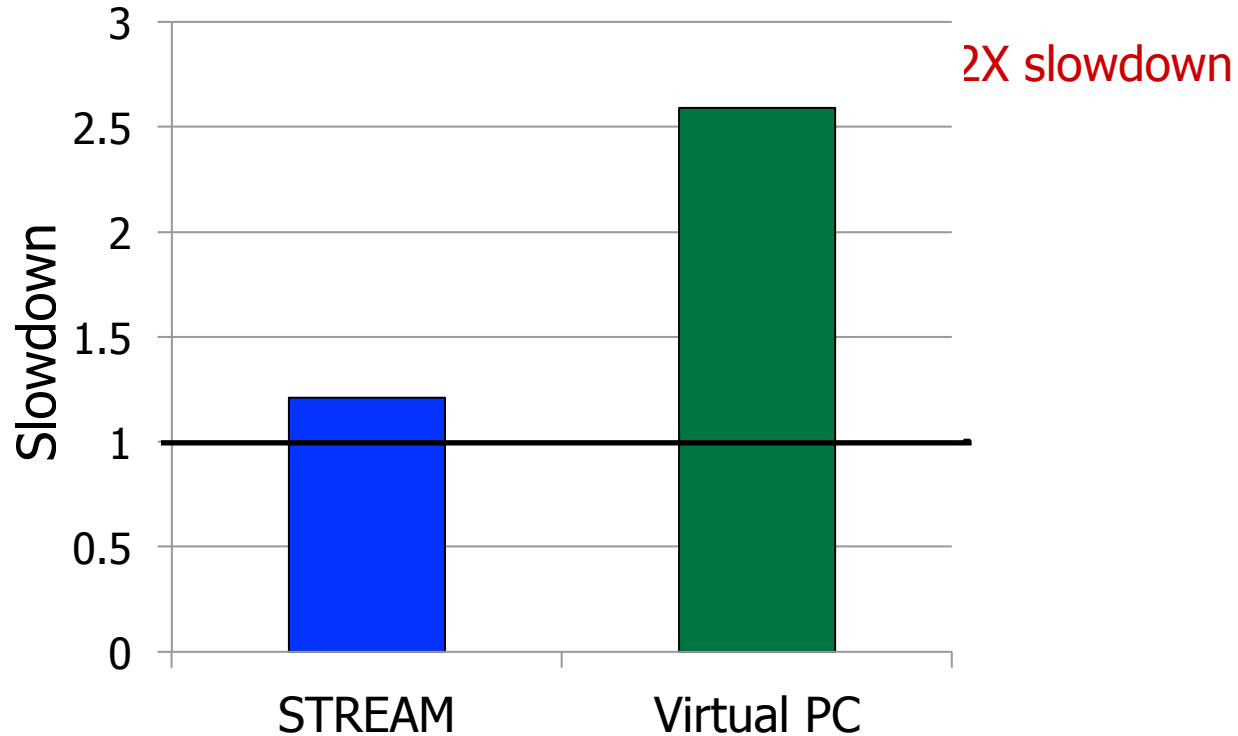
What Does the Memory Hog Do?



Row size: 8KB, cache block size: 64B
128 (8KB/64B) requests of T0 serviced before T1

Moscibroda and Mutlu, “[Memory Performance Attacks](#),” USENIX Security 2007.

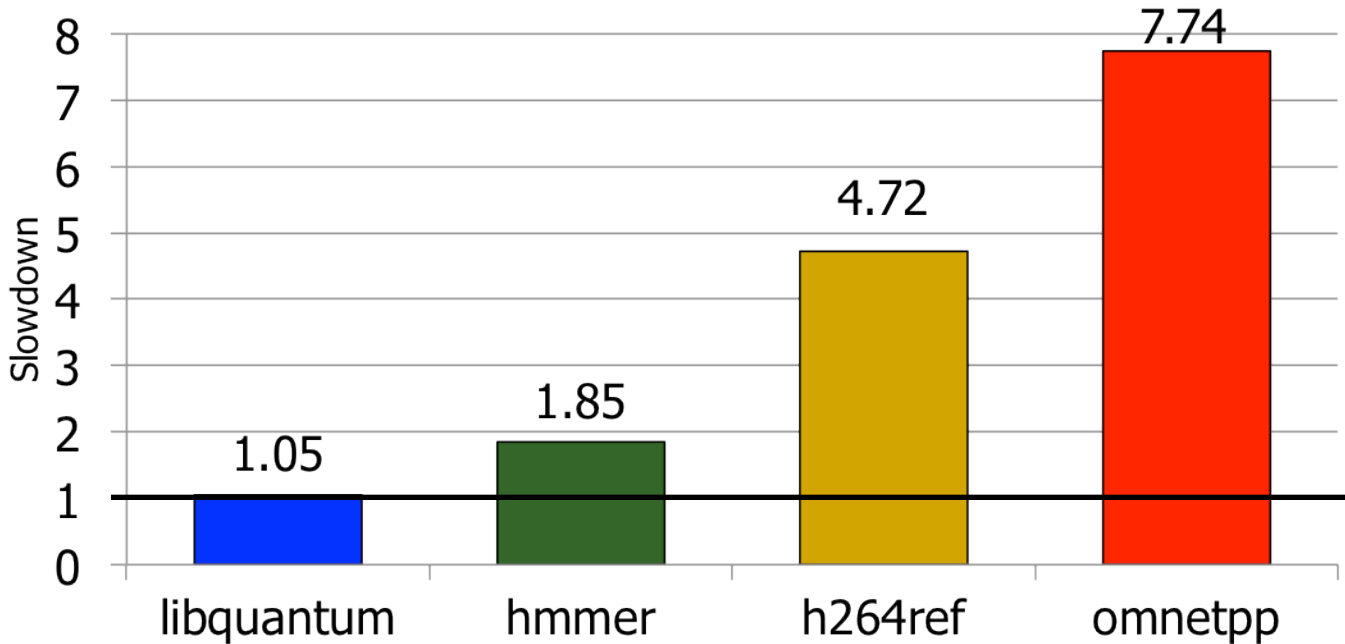
Effect of the Memory Performance Hog



Results on Intel Pentium D running Windows XP
(Similar results for Intel Core Duo and AMD Turion, and on Fedora Linux)

Moscibroda and Mutlu, “[Memory Performance Attacks](#),” USENIX Security 2007.

Greater Problem with More Cores



- Vulnerable to denial of service (DoS) [Usenix Security'07]
- Unable to enforce priorities or SLAs [MICRO'07,'10,'11, ISCA'08'11'12, ASPLOS'10]
- Low system performance [IEEE Micro Top Picks '09,'11a,'11b,'12]

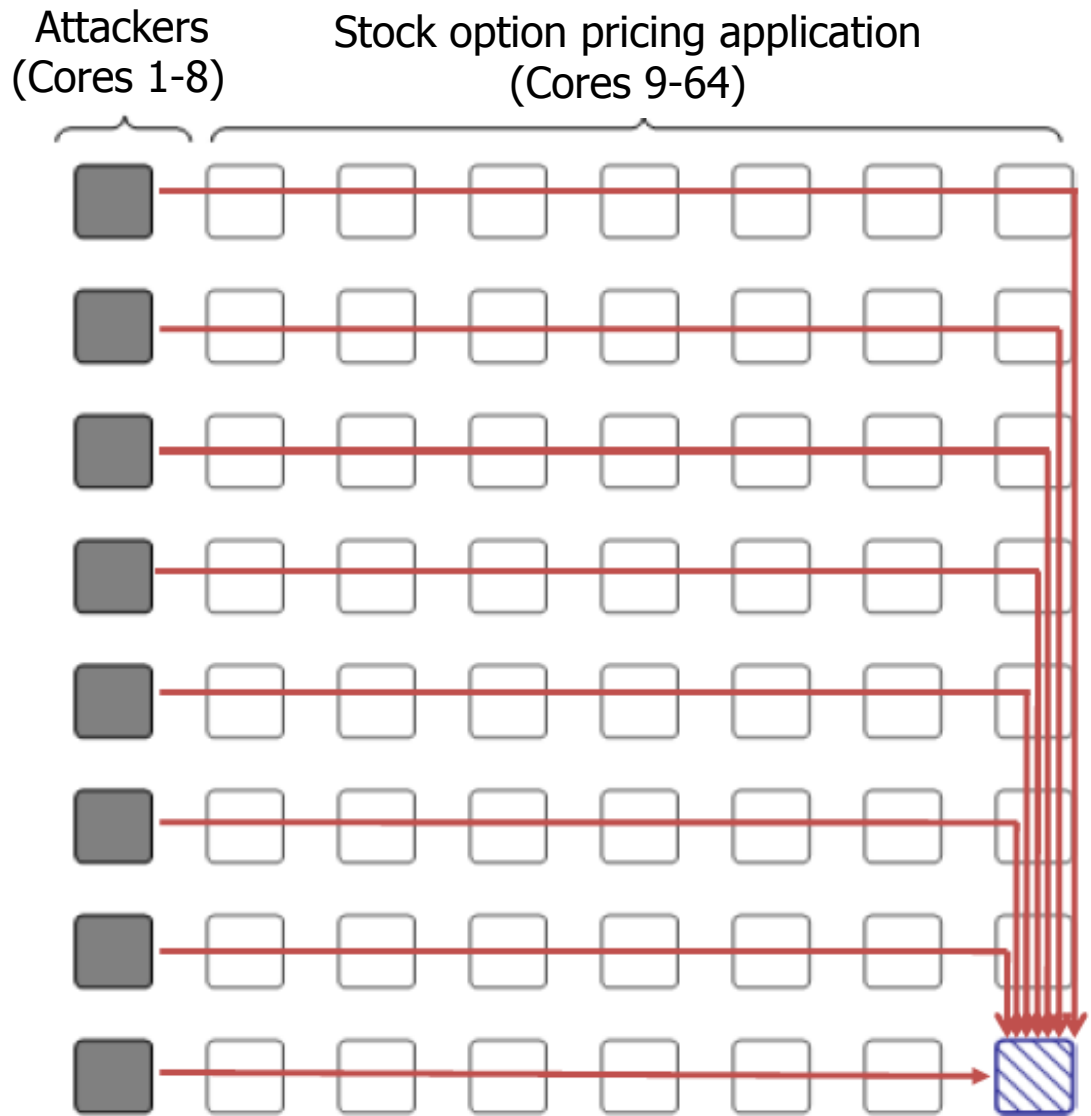
Uncontrollable, unpredictable system

Distributed DoS in Networked Multi-Core Systems

Cores connected via
packet-switched
routers on chip

~5000X slowdown

Grot, Hestness, Keckler, Mutlu,
“Preemptive virtual clock: A Flexible,
Efficient, and Cost-effective QOS
Scheme for Networks-on-Chip,”
MICRO 2009.



Solution: QoS-Aware, Predictable Memory

- Problem: Memory interference is uncontrolled → uncontrollable, unpredictable, vulnerable system
- Goal: We need to control it → Design a QoS-aware system
- Solution: Hardware/software cooperative memory QoS
 - Hardware designed to provide a configurable fairness substrate
 - Application-aware memory scheduling, partitioning, throttling
 - Software designed to configure the resources to satisfy different QoS goals
 - E.g., fair, programmable memory controllers and on-chip networks provide QoS and predictable performance
[2007-2012, Top Picks'09,'11a,'11b,'12]

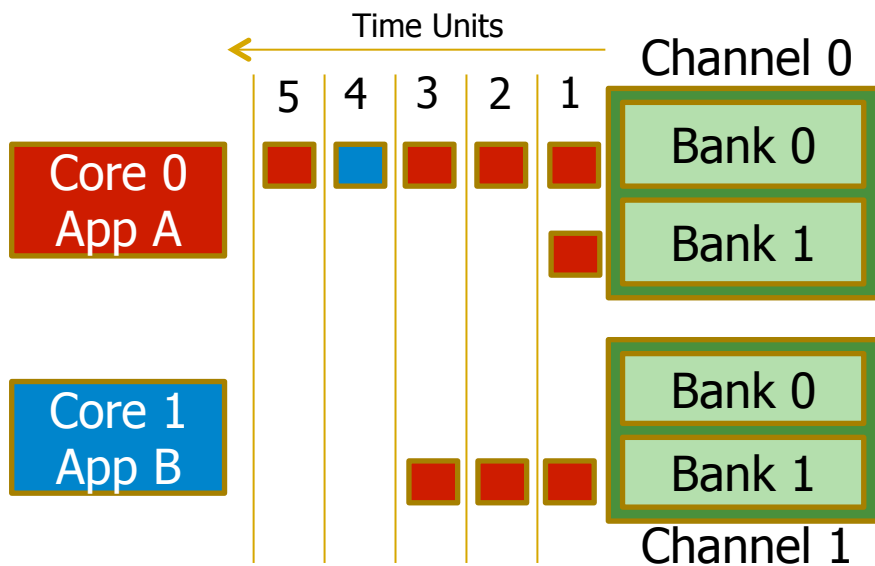
Designing QoS-Aware Memory Systems: Approaches

- **Smart resources:** Design each shared resource to have a configurable interference control/reduction mechanism
 - QoS-aware memory controllers [Mutlu+ MICRO'07] [Moscibroda+, Usenix Security'07] [Mutlu+ ISCA'08, Top Picks'09] [Kim+ HPCA'10] [Kim+ MICRO'10, Top Picks'11] [Ebrahimi+ ISCA'11, MICRO'11] [Ausavarungnirun+, ISCA'12]
 - QoS-aware interconnects [Das+ MICRO'09, ISCA'10, Top Picks '11] [Grot+ MICRO'09, ISCA'11, Top Picks '12]
 - QoS-aware caches
- **Dumb resources:** Keep each resource free-for-all, but reduce/control interference by injection control or data mapping
 - Source throttling to control access to memory system [Ebrahimi+ ASPLOS'10, ISCA'11, TOCS'12] [Ebrahimi+ MICRO'09] [Nychis+ HotNets'10]
 - QoS-aware data mapping to memory controllers [Muralidhara+ MICRO'11]
 - QoS-aware thread scheduling to cores

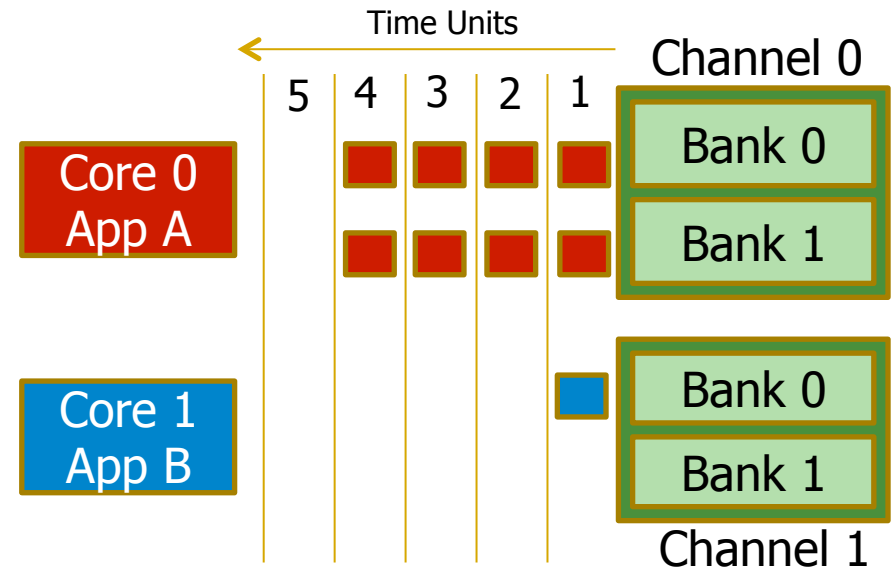
A Mechanism to Reduce Memory Interference

■ Memory Channel Partitioning

- ❑ Idea: System software maps badly-interfering applications' pages to different channels [Muralidhara+, MICRO'11]



Conventional Page Mapping



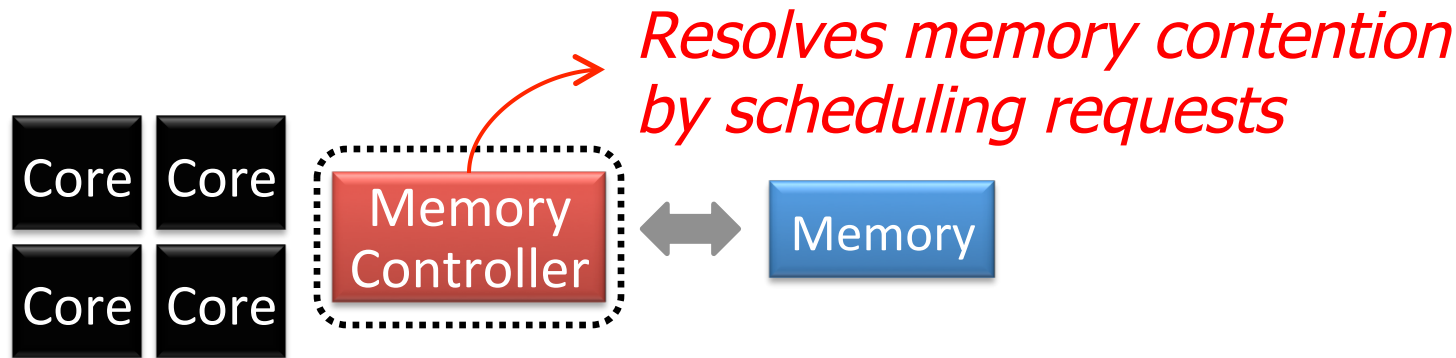
Channel Partitioning

- Separate data of low/high intensity and low/high row-locality applications
- Especially effective in reducing interference of threads with “medium” and “heavy” memory intensity
 - ❑ 11% higher performance over existing systems (200 workloads)

Designing QoS-Aware Memory Systems: Approaches

- **Smart resources:** Design each shared resource to have a configurable interference control/reduction mechanism
 - **QoS-aware memory controllers** [Mutlu+ MICRO'07] [Moscibroda+, Usenix Security'07] [Mutlu+ ISCA'08, Top Picks'09] [Kim+ HPCA'10] [Kim+ MICRO'10, Top Picks'11] [Ebrahimi+ ISCA'11, MICRO'11] [Ausavarungnirun+, ISCA'12]
 - QoS-aware interconnects [Das+ MICRO'09, ISCA'10, Top Picks '11] [Grot+ MICRO'09, ISCA'11, Top Picks '12]
 - QoS-aware caches
- **Dumb resources:** Keep each resource free-for-all, but reduce/control interference by injection control or data mapping
 - Source throttling to control access to memory system [Ebrahimi+ ASPLOS'10, ISCA'11, TOCS'12] [Ebrahimi+ MICRO'09] [Nychis+ HotNets'10]
 - QoS-aware data mapping to memory controllers [Muralidhara+ MICRO'11]
 - QoS-aware thread scheduling to cores

QoS-Aware Memory Scheduling



- How to schedule requests to provide
 - ❑ High system performance
 - ❑ High fairness to applications
 - ❑ Configurability to system software

- Memory controller needs to be aware of threads

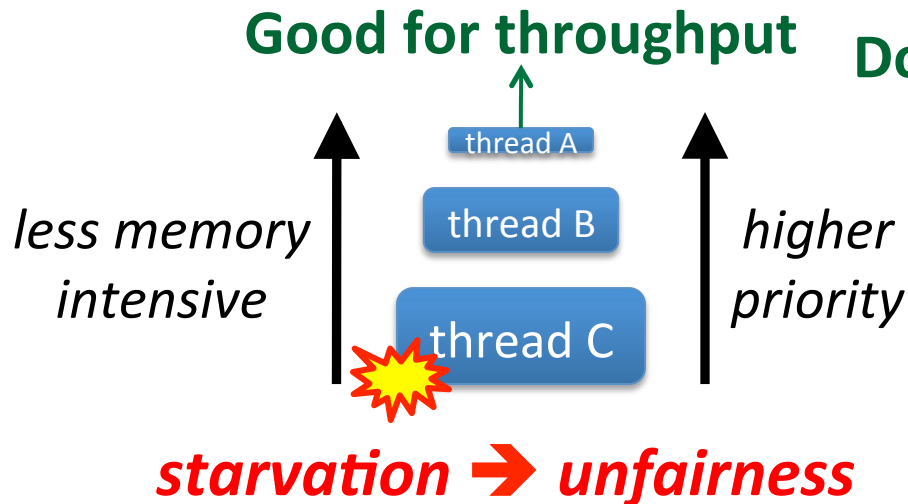
QoS-Aware Memory Scheduling: Evolution

- **Stall-time fair memory scheduling** [Mutlu+ MICRO'07]
 - Idea: Estimate and balance thread slowdowns
 - Takeaway: **Proportional thread progress improves performance, especially when threads are "heavy"** (memory intensive)
- **Parallelism-aware batch scheduling** [Mutlu+ ISCA'08, Top Picks'09]
 - Idea: Rank threads and service in rank order (to preserve bank parallelism); batch requests to prevent starvation
 - Takeaway: **Preserving within-thread bank-parallelism improves performance**; request batching improves fairness
- **ATLAS memory scheduler** [Kim+ HPCA'10]
 - Idea: Prioritize threads that have attained the least service from the memory scheduler
 - Takeaway: **Prioritizing "light" threads improves performance**

Throughput vs. Fairness

Throughput biased approach

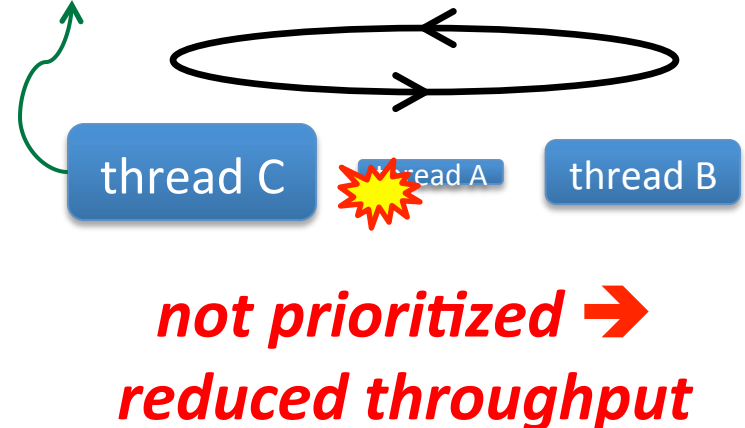
Prioritize less memory-intensive threads



Fairness biased approach

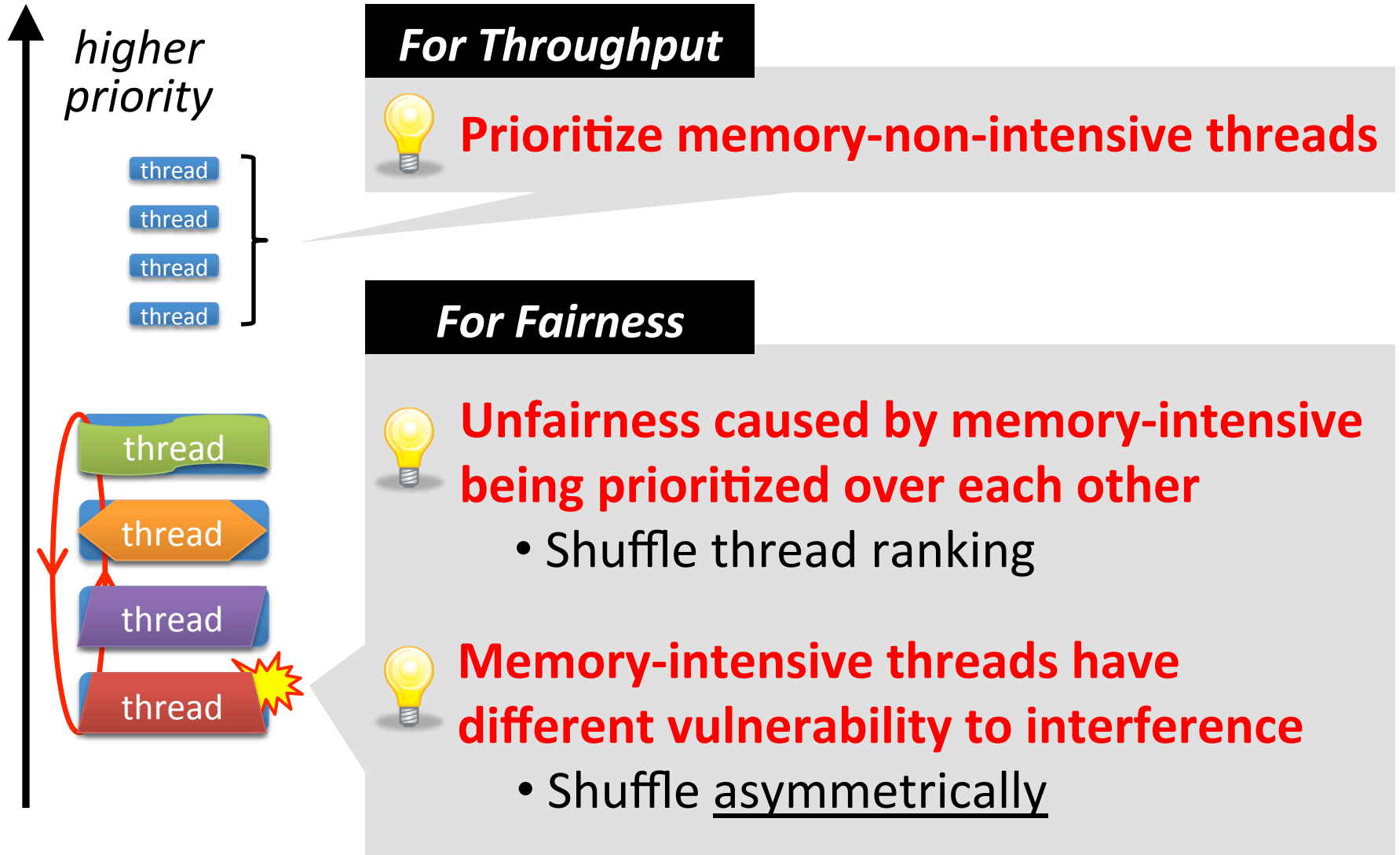
Take turns accessing memory

Does not starve



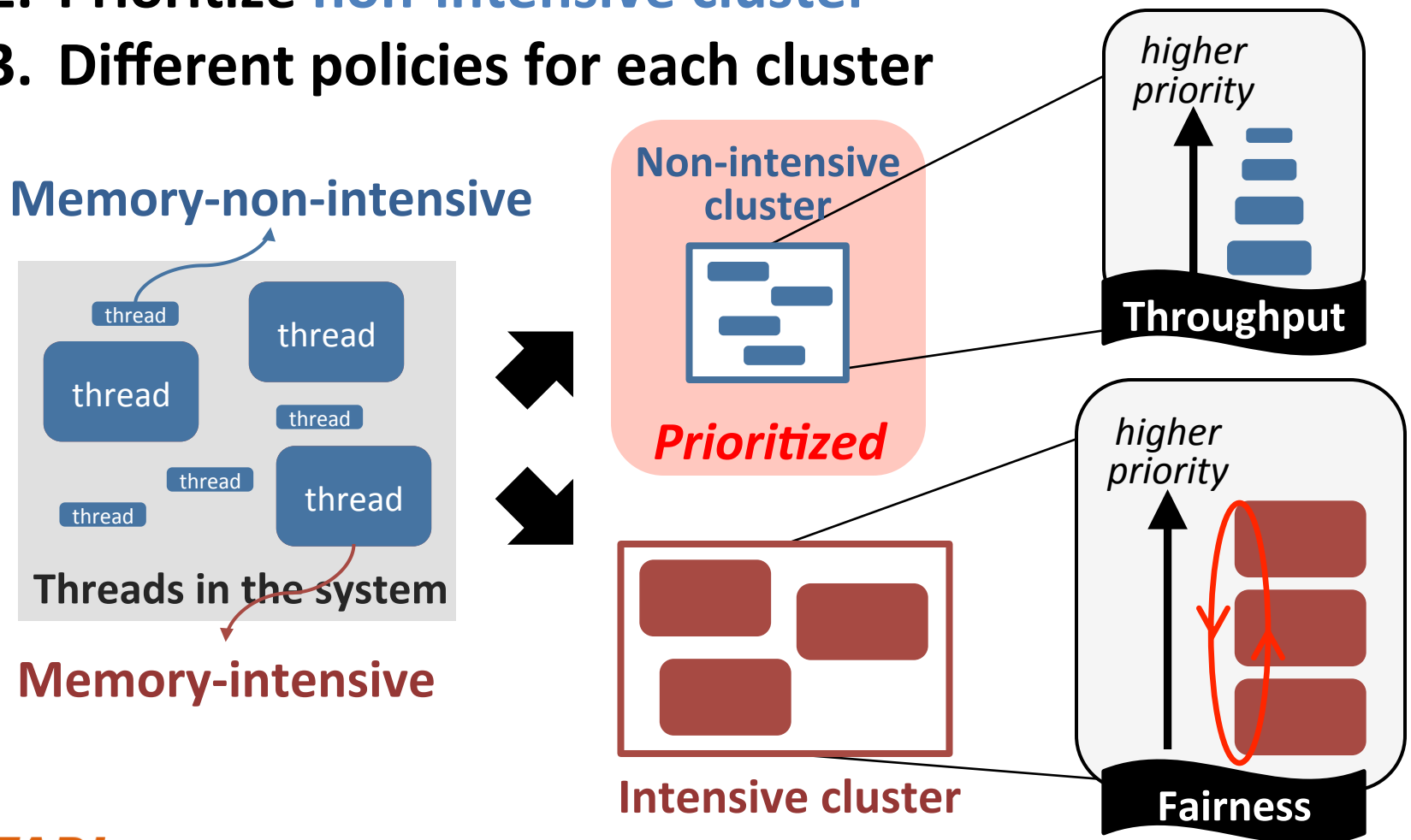
Single policy for all threads is insufficient

Achieving the Best of Both Worlds



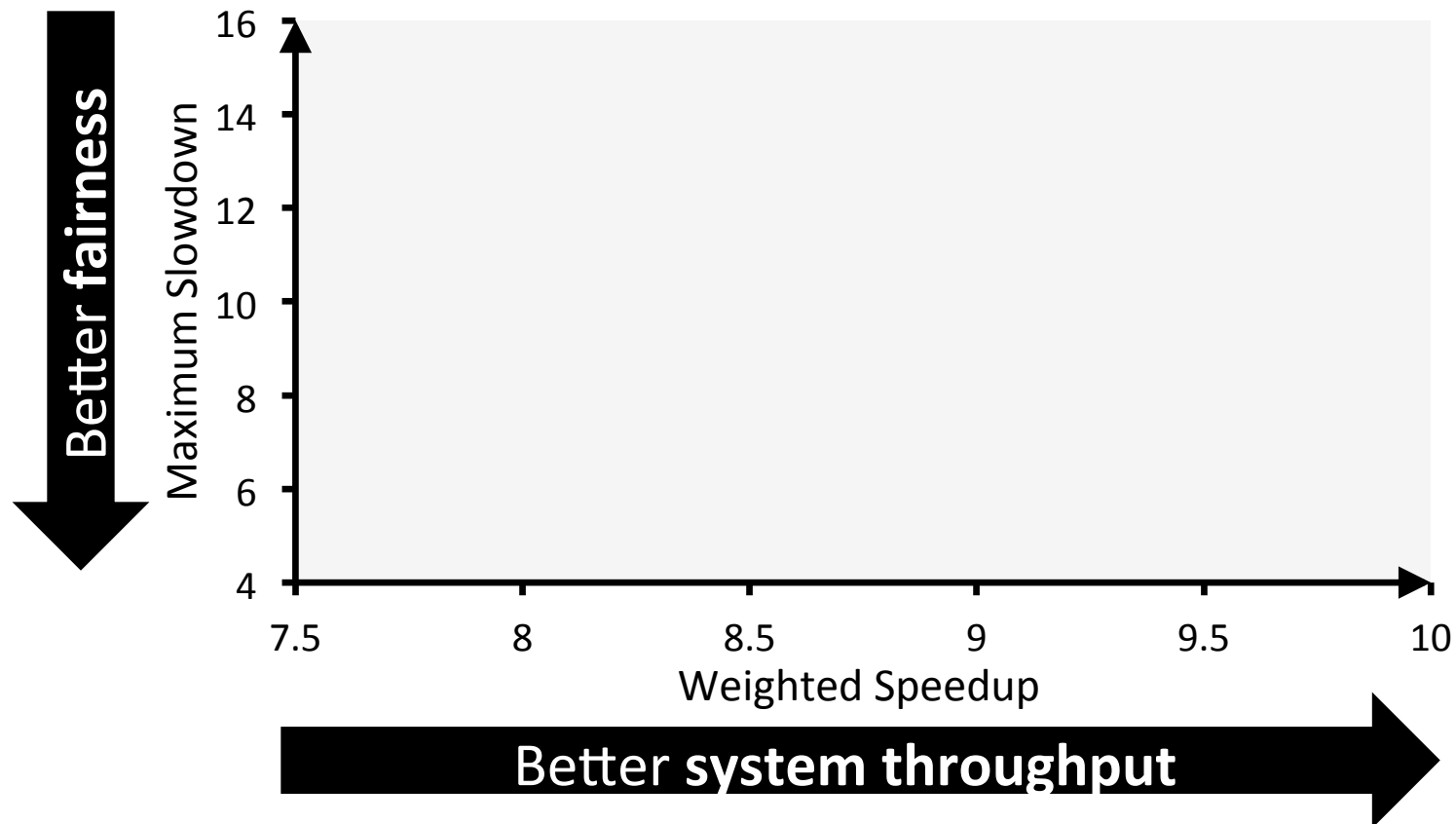
Thread Cluster Memory Scheduling [Kim+ MICRO'10]

1. Group threads into two **clusters**
2. Prioritize **non-intensive cluster**
3. Different policies for each cluster



TCM: Throughput and Fairness

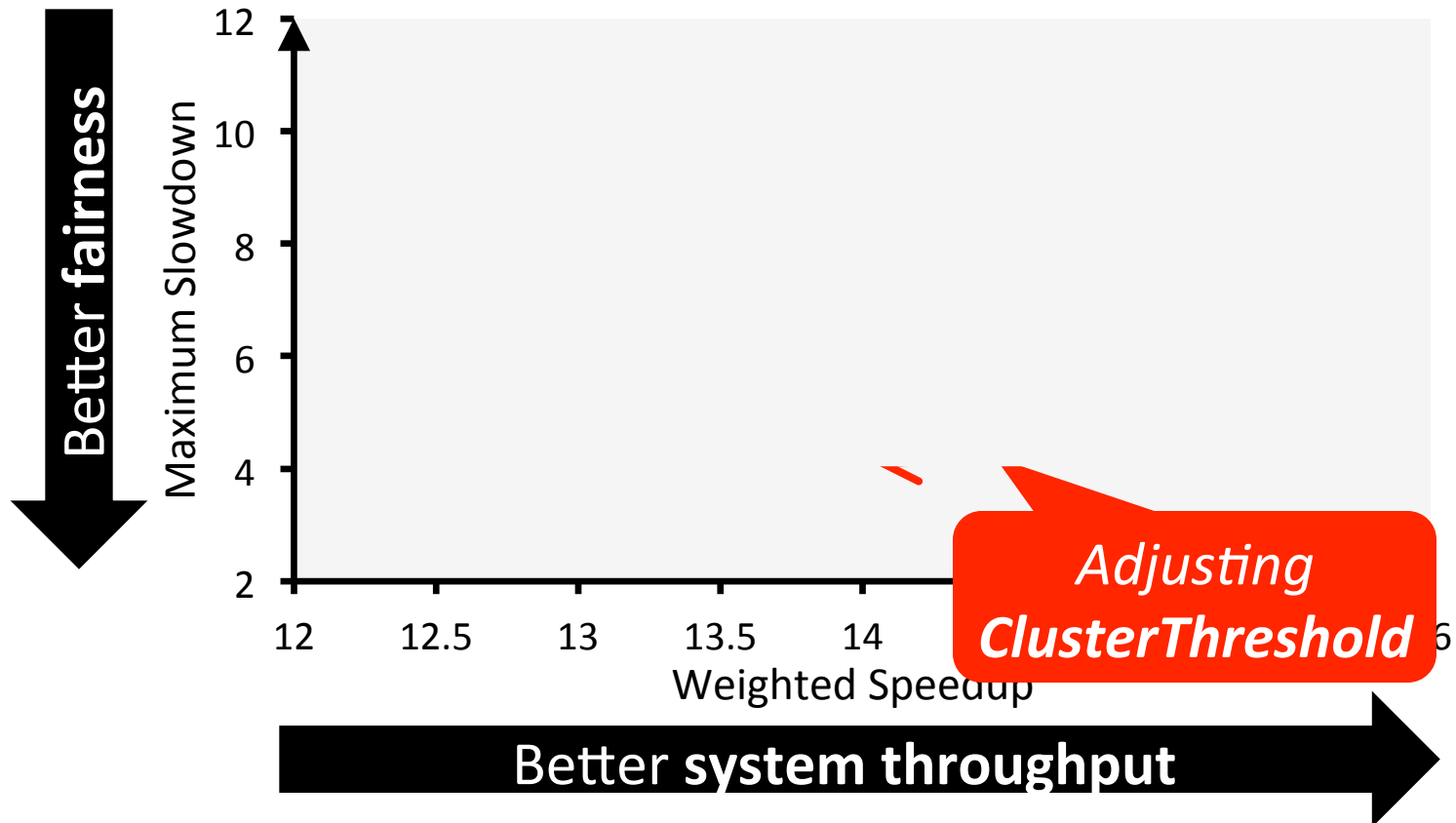
24 cores, 4 memory controllers, 96 workloads



*TCM, a heterogeneous scheduling policy,
provides best fairness and system throughput*

TCM: Fairness-Throughput Tradeoff

When configuration parameter is varied...



TCM allows robust fairness-throughput tradeoff

Memory Control in CPU-GPU Systems

- **Observation:** Heterogeneous CPU-GPU systems require memory schedulers with **large request buffers**
- **Problem:** Existing monolithic application-aware memory scheduler designs are **hard to scale** to large request buffer sizes
- **Solution:** Staged Memory Scheduling (SMS)
decomposes the memory controller into three simple stages:
 - 1) Batch formation: maintains row buffer locality
 - 2) Batch scheduler: reduces interference between applications
 - 3) DRAM command scheduler: issues requests to DRAM
- Compared to state-of-the-art memory schedulers:
 - ❑ SMS is significantly simpler and more scalable
 - ❑ SMS provides higher performance and fairness

Memory QoS in a Parallel Application

- Threads in a multithreaded application are inter-dependent
- Some threads can be on the critical path of execution due to synchronization; some threads are not
- How do we schedule requests of inter-dependent threads to maximize multithreaded application performance?
- Idea: **Estimate limiter threads** likely to be on the critical path and prioritize their requests; **shuffle priorities of non-limiter threads** to reduce memory interference among them [Ebrahimi+, MICRO'11]
- Hardware/software cooperative limiter thread estimation:
 - Thread executing the most contended critical section
 - Thread that is falling behind the most in a *parallel for* loop

Some Related Past Work

- That I could not cover...
- How to handle prefetch requests in a QoS-aware multi-core memory system?
 - ❑ Prefetch-aware shared resource management, ISCA'11. [ISCA 2011 Talk](#)
 - ❑ Prefetch-aware memory controllers, MICRO'08, IEEE-TC'11. [Micro 2008 Talk](#)
 - ❑ Coordinated control of multiple prefetchers, MICRO'09. [Micro 2009 Talk](#)
- How to design QoS mechanisms in the interconnect?
 - ❑ Topology-aware, scalable QoS, ISCA'11.
 - ❑ Slack-based packet scheduling, ISCA'10.
 - ❑ Efficient bandwidth guarantees, MICRO'09.
 - ❑ Application-aware request prioritization, MICRO'09.

Summary: Memory QoS Approaches and Techniques

- Approaches: Smart vs. dumb resources
 - ❑ Smart resources: QoS-aware memory scheduling
 - ❑ Dumb resources: Source throttling; channel partitioning
 - ❑ Both approaches are effective in reducing interference
 - ❑ No single best approach for all workloads

- Techniques: Request scheduling, source throttling, memory partitioning
 - ❑ All approaches are effective in reducing interference
 - ❑ Can be applied at different levels: hardware vs. software
 - ❑ No single best technique for all workloads

- Combined approaches and techniques are the most powerful
 - ❑ Integrated Memory Channel Partitioning and Scheduling [MICRO'11]

Partial List of Referenced/ Related Papers

Heterogeneous Cores

- M. Aater Suleman, Onur Mutlu, Moinuddin K. Qureshi, and Yale N. Patt, **"Accelerating Critical Section Execution with Asymmetric Multi-Core Architectures"**
Proceedings of the 14th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), pages 253-264, Washington, DC, March 2009. [Slides \(ppt\)](#)
- M. Aater Suleman, Onur Mutlu, Jose A. Joao, Khubaib, and Yale N. Patt, **"Data Marshaling for Multi-core Architectures"**
Proceedings of the 37th International Symposium on Computer Architecture (ISCA), pages 441-450, Saint-Malo, France, June 2010. [Slides \(ppt\)](#)
- Jose A. Joao, M. Aater Suleman, Onur Mutlu, and Yale N. Patt, **"Bottleneck Identification and Scheduling in Multithreaded Applications"**
Proceedings of the 17th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), London, UK, March 2012. [Slides \(ppt\)](#) [\(pdf\)](#)

QoS-Aware Memory Systems (I)

- Rachata Ausavarungnirun, Kevin Chang, Lavanya Subramanian, Gabriel Loh, and Onur Mutlu,
"Staged Memory Scheduling: Achieving High Performance and Scalability in Heterogeneous Systems"
Proceedings of the 39th International Symposium on Computer Architecture (ISCA), Portland, OR, June 2012.
- Sai Prashanth Muralidhara, Lavanya Subramanian, Onur Mutlu, Mahmut Kandemir, and Thomas Moscibroda,
"Reducing Memory Interference in Multicore Systems via Application-Aware Memory Channel Partitioning"
Proceedings of the 44th International Symposium on Microarchitecture (MICRO), Porto Alegre, Brazil, December 2011
- Yoongu Kim, Michael Papamichael, Onur Mutlu, and Mor Harchol-Balter,
"Thread Cluster Memory Scheduling: Exploiting Differences in Memory Access Behavior"
Proceedings of the 43rd International Symposium on Microarchitecture (MICRO), pages 65-76, Atlanta, GA, December 2010. [Slides \(pptx\)](#) [\(pdf\)](#)
- Eiman Ebrahimi, Chang Joo Lee, Onur Mutlu, and Yale N. Patt,
"Fairness via Source Throttling: A Configurable and High-Performance Fairness Substrate for Multi-Core Memory Systems"
ACM Transactions on Computer Systems (TOCS), April 2012.

QoS-Aware Memory Systems (II)

- Onur Mutlu and Thomas Moscibroda,
"Parallelism-Aware Batch Scheduling: Enabling High-Performance and Fair Memory Controllers"
*IEEE Micro, Special Issue: Micro's Top Picks from 2008 Computer Architecture Conferences (**MICRO TOP PICKS**), Vol. 29, No. 1, pages 22-32, January/February 2009.*
- Onur Mutlu and Thomas Moscibroda,
"Stall-Time Fair Memory Access Scheduling for Chip Multiprocessors"
*Proceedings of the 40th International Symposium on Microarchitecture (**MICRO**), pages 146-158, Chicago, IL, December 2007. [Slides \(ppt\)](#)*
- Thomas Moscibroda and Onur Mutlu,
"Memory Performance Attacks: Denial of Memory Service in Multi-Core Systems"
*Proceedings of the 16th USENIX Security Symposium (**USENIX SECURITY**), pages 257-274, Boston, MA, August 2007. [Slides \(ppt\)](#)*

QoS-Aware Memory Systems (III)

- Eiman Ebrahimi, Rustam Miftakhutdinov, Chris Fallin, Chang Joo Lee, Onur Mutlu, and Yale N. Patt,
"Parallel Application Memory Scheduling"
*Proceedings of the 44th International Symposium on Microarchitecture (**MICRO**), Porto Alegre, Brazil, December 2011. Slides (pptx)*
- Boris Grot, Joel Hestness, Stephen W. Keckler, and Onur Mutlu,
"Kilo-NOC: A Heterogeneous Network-on-Chip Architecture for Scalability and Service Guarantees"
*Proceedings of the 38th International Symposium on Computer Architecture (**ISCA**), San Jose, CA, June 2011. Slides (pptx)*
- Reetuparna Das, Onur Mutlu, Thomas Moscibroda, and Chita R. Das,
"Application-Aware Prioritization Mechanisms for On-Chip Networks"
*Proceedings of the 42nd International Symposium on Microarchitecture (**MICRO**), pages 280-291, New York, NY, December 2009. Slides (pptx)*

Heterogeneous Memory

- Justin Meza, Jichuan Chang, HanBin Yoon, Onur Mutlu, and Parthasarathy Ranganathan, **"Enabling Efficient and Scalable Hybrid Memories Using Fine-Granularity DRAM Cache Management"**
IEEE Computer Architecture Letters (CAL), May 2012.
- HanBin Yoon, Justin Meza, Rachata Ausavarungnirun, Rachael Harding, and Onur Mutlu, **"Row Buffer Locality-Aware Data Placement in Hybrid Memories"**
SAFARI Technical Report, TR-SAFARI-2011-005, Carnegie Mellon University, September 2011.
- Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger, **"Architecting Phase Change Memory as a Scalable DRAM Alternative"**
Proceedings of the 36th International Symposium on Computer Architecture (ISCA), pages 2-13, Austin, TX, June 2009. [Slides \(pdf\)](#)
- Benjamin C. Lee, Ping Zhou, Jun Yang, Youtao Zhang, Bo Zhao, Engin Ipek, Onur Mutlu, and Doug Burger, **"Phase Change Technology and the Future of Main Memory"**
IEEE Micro, Special Issue: Micro's Top Picks from 2009 Computer Architecture Conferences (MICRO TOP PICKS), Vol. 30, No. 1, pages 60-70, January/February 2010.

Flash Memory

- Yu Cai, Eric F. Haratsch, Onur Mutlu, and Ken Mai,
"Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis"
Proceedings of the Design, Automation, and Test in Europe Conference (DATE), Dresden, Germany, March 2012. Slides (ppt)

Latency Tolerance

- Onur Mutlu, Jared Stark, Chris Wilkerson, and Yale N. Patt,
"Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-order Processors"
Proceedings of the
9th International Symposium on High-Performance Computer Architecture
(HPCA), pages 129-140, Anaheim, CA, February 2003. [Slides \(pdf\)](#)
- Onur Mutlu, Hyesoon Kim, and Yale N. Patt,
"Techniques for Efficient Processing in Runahead Execution Engines"
Proceedings of the 32nd International Symposium on Computer Architecture
(ISCA), pages 370-381, Madison, WI, June 2005. [Slides \(ppt\)](#) [Slides \(pdf\)](#)
- Onur Mutlu, Hyesoon Kim, and Yale N. Patt,
"Address-Value Delta (AVD) Prediction: Increasing the Effectiveness of Runahead Execution by Exploiting Regular Memory Allocation Patterns"
Proceedings of the 38th International Symposium on Microarchitecture
(MICRO), pages 233-244, Barcelona, Spain, November 2005. [Slides \(ppt\)](#)
[Slides \(pdf\)](#)

Scaling DRAM: Refresh and Parallelism

- Jamie Liu, Ben Jaiyen, Richard Veras, and Onur Mutlu,
"RAIDR: Retention-Aware Intelligent DRAM Refresh"
Proceedings of the 39th International Symposium on Computer Architecture (ISCA), Portland, OR, June 2012.
- Yoongu Kim, Vivek Seshadri, Donghyuk Lee, Jamie Liu, and Onur Mutlu,
"A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM"
Proceedings of the 39th International Symposium on Computer Architecture (ISCA), Portland, OR, June 2012.