

Fairness via Source Throttling:

A Configurable and High-Performance Fairness Substrate for
Multi-Core Memory Systems

Eiman Ebrahimi | Chang Joo Lee | Onur Mutlu | Yale N. Patt

ASPLOS 2010

HPS Research Group
The University of Texas at Austin

Computer Architecture Laboratory
Carnegie Mellon University

Executive Summary

- **Motivation:** Cores in a chip-multiprocessor system share multiple hardware resources in the memory subsystem
 - Interference in the shared resources can lead to unfair slowdown for some applications
- **Problem:** Existing fairness mechanisms focus on a single resource
 - Multiple independently implemented mechanisms can make contradictory decisions, leading to low fairness and loss of performance
- **Goal:** provides fairness in the entire shared memory system without degrading performance
- **Key Contributions:** Fairness via Source Throttling(FST) provides two major mechanisms
 - 1) Runtime fairness evaluation
 - 2) Dynamic request throttling
- **Major Results:** improve performance by 25.6%/14.5% and reduce unfairness by 44%/36%

Outline

- Background, Problem & Goal
- Novelty
- Key Approach & Ideas
- Mechanism
- Methodology & Evaluation
- Conclusion

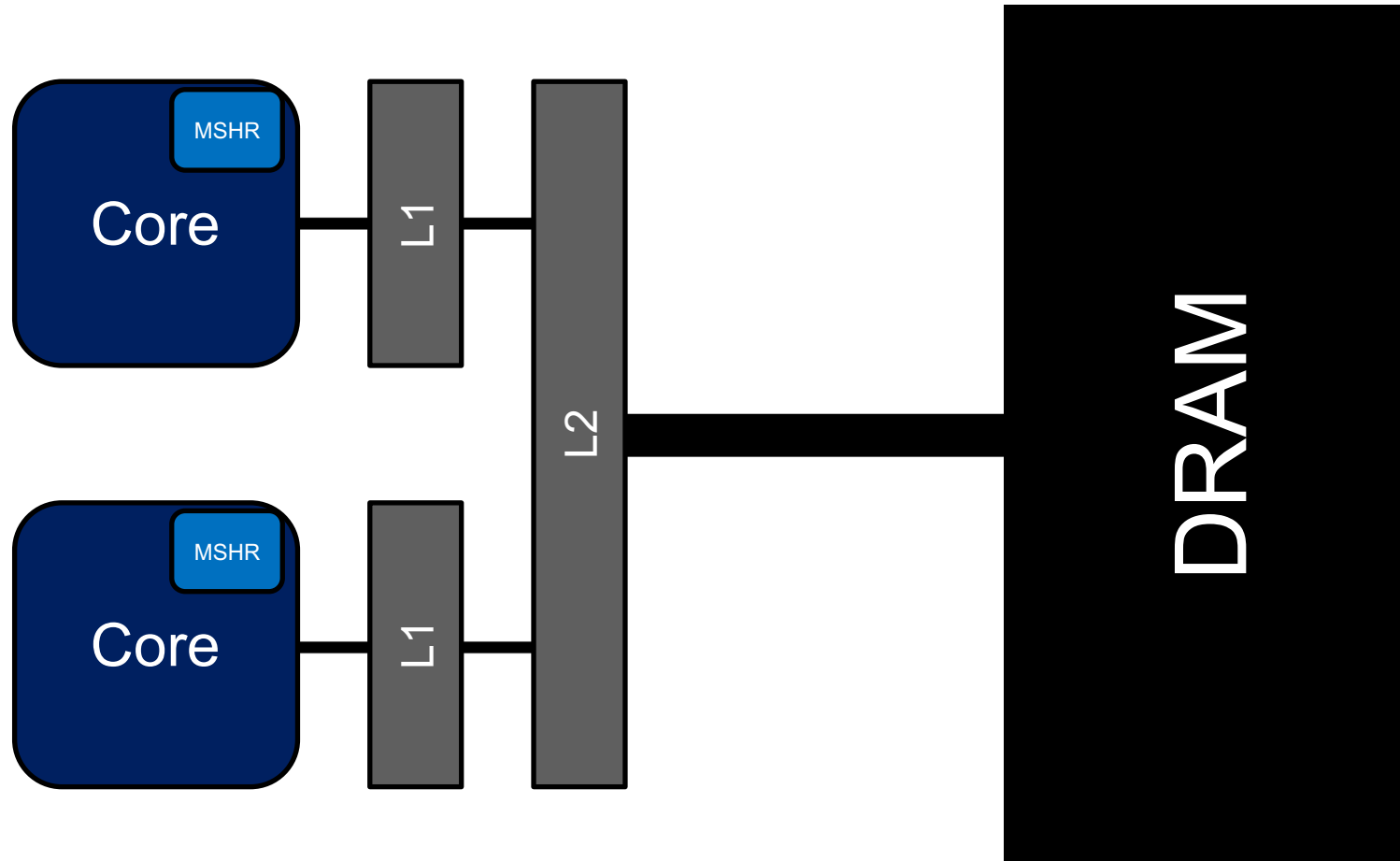
- Takeaways
- Strengths and Weaknesses
- Thoughts & Ideas
- Open Discussion

Outline

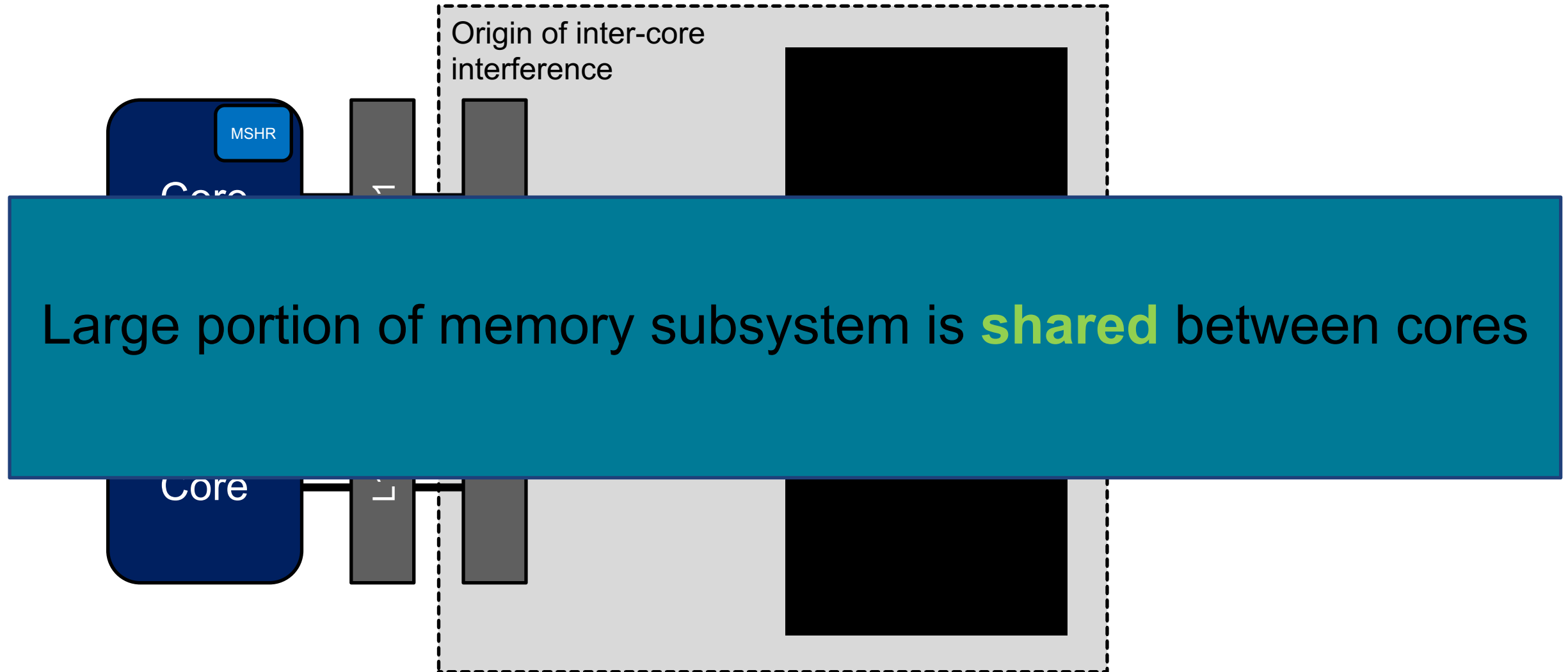
- Background, Problem & Goal
- Novelty
- Key Approach & Ideas
- Mechanism
- Methodology & Evaluation
- Conclusion

- Takeaways
- Strengths and Weaknesses
- Thoughts & Ideas
- Open Discussion

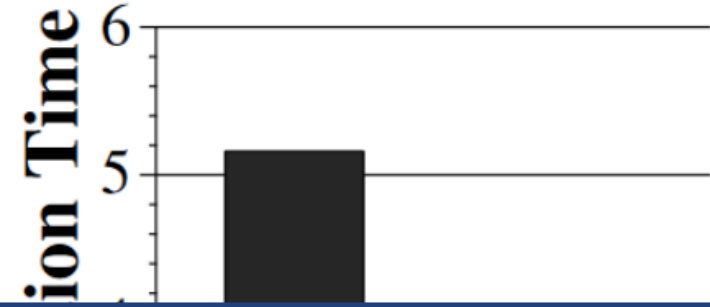
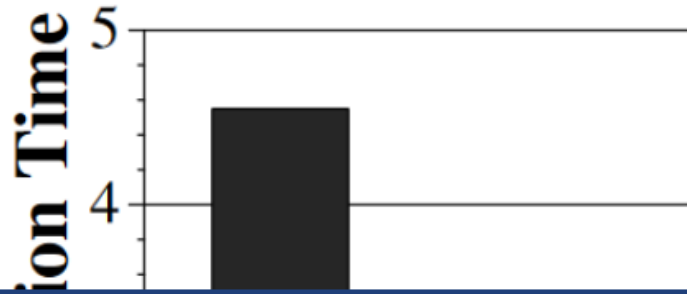
System Layout



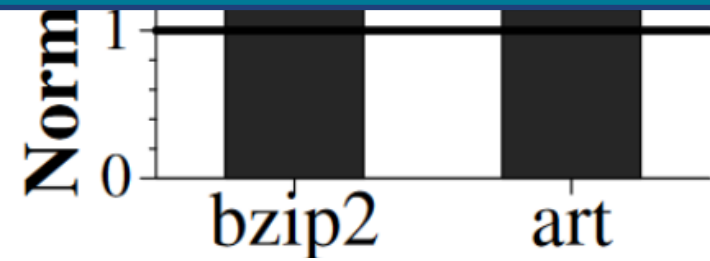
System Layout



Interference and delays lead to slowdown



Goal: all applications of **equal priority**
experience the **same slowdown**



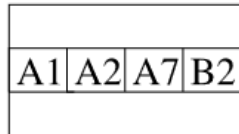
Disparity in slowdowns due to unfairness

Previous Approach

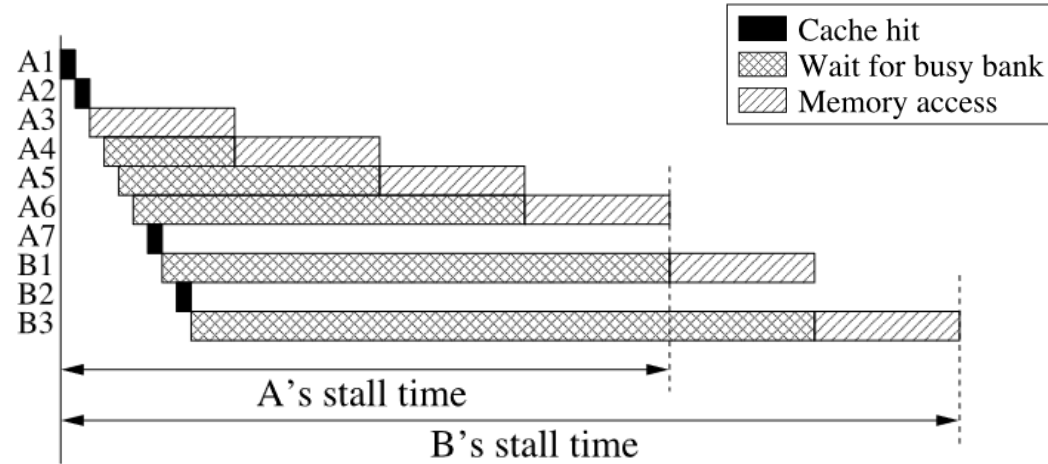
- Previous research focused on individual resources

Previous Approach

Access order:
A1, A2, A3, A4, A5, A6, A7, B1, B2, B3

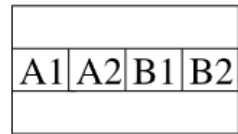


Shared L2 cache

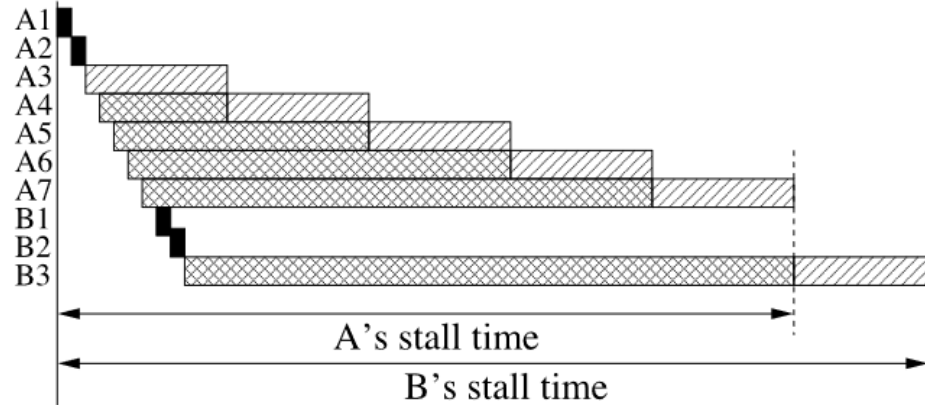


(e) Initial state for no fairness control (f) Memory-related stall time of no fairness control

Access order:
A1, A2, A3, A4, A5, A6, A7, B1, B2, B3



Shared L2 cache



(g) Initial state for fair cache

(h) Memory-related time of fair cache

Previous Approach

- Previous research focused on **individual** resources
- It is **challenging** to properly **coordinate multiple** fairness mechanisms
 - Partitioning one resource may change demands on another shared resource

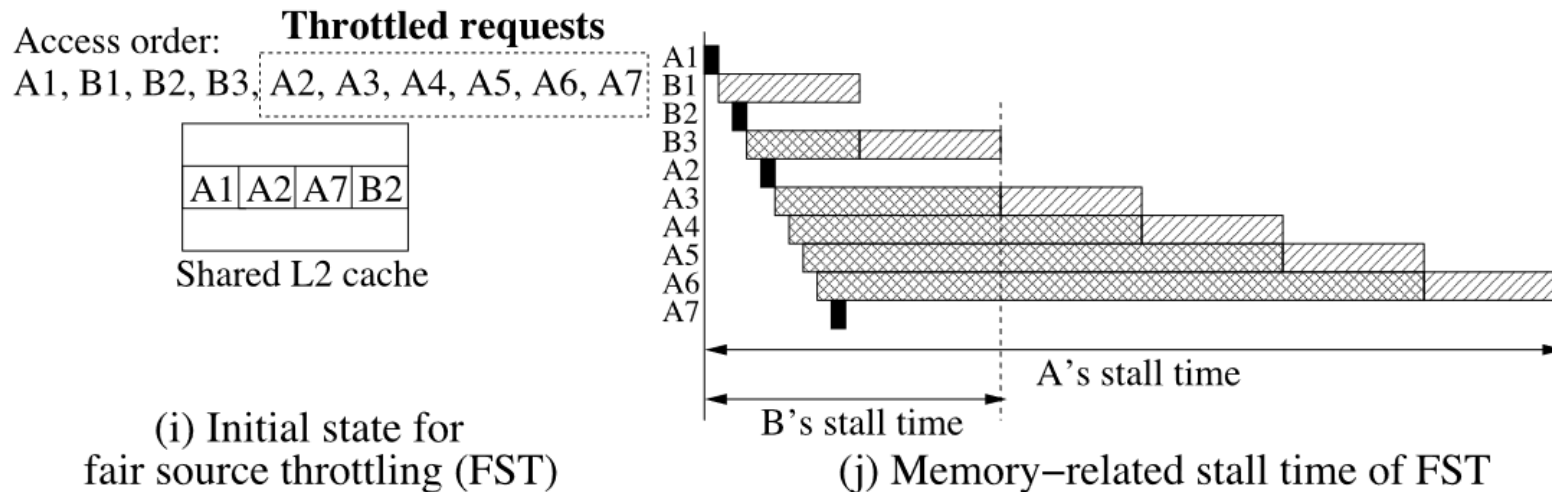
Outline

- Background, Problem & Goal
- **Novelty**
- Key Approach & Ideas
- Mechanism
- Methodology & Evaluation
- Conclusion

- Takeaways
- Strengths and Weaknesses
- Thoughts & Ideas
- Open Discussion

System-Wide Fairness

- Tackle unfairness in the **entire shared memory system**
 - Eliminate the need for multiple fairness mechanisms
- Control fairness by **orchestrating memory requests**
 - Rate of memory request injections
 - Number of memory request injections

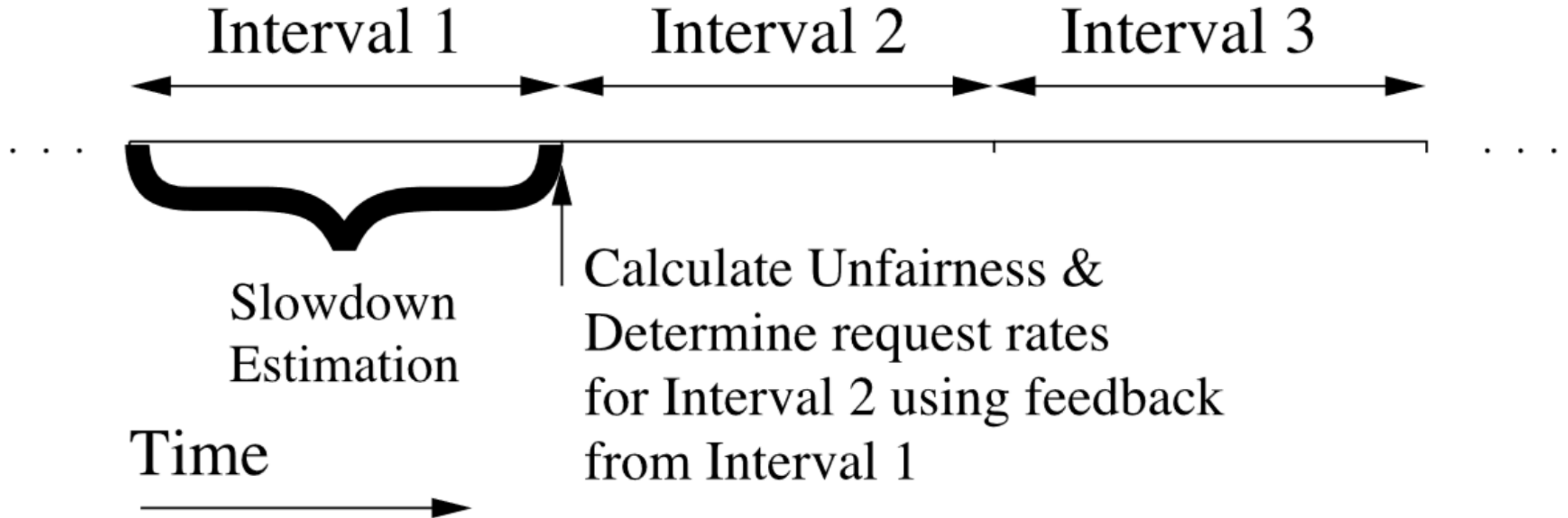


Outline

- Background, Problem & Goal
- Novelty
- **Key Approach & Ideas**
- Mechanism
- Methodology & Evaluation
- Conclusion

- Takeaways
- Strengths and Weaknesses
- Thoughts & Ideas
- Open Discussion

Interval based Estimation and Throttling



Outline

- Background, Problem & Goal
- Novelty
- Key Approach & Ideas
- **Mechanism**
- Methodology & Evaluation
- Conclusion

- Takeaways
- Strengths and Weaknesses
- Thoughts & Ideas
- Open Discussion

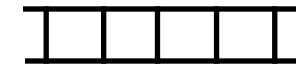
Runtime Fairness Evaluation

- Goal: **dynamically estimate** system unfairness
- Slowdown: $T_{\text{shared}}/T_{\text{alone}}$
 - T_{shared} : number of cycles to execute **simultaneously** with other applications
 - T_{alone} : number of cycles to execute **alone**
- Estimating T_{alone} while running multiple applications
 - T_{excess} : number of excess execution cycles induced by **inter-core interference**
 - $T_{\text{alone}} = T_{\text{shared}} - T_{\text{excess}}$

$$IS_i = \frac{T_i^{\text{shared}}}{T_i^{\text{alone}}}, \quad \text{Unfairness} = \frac{\text{MAX}\{IS_0, IS_1, \dots, IS_{N-1}\}}{\text{MIN}\{IS_0, IS_1, \dots, IS_{N-1}\}}$$

Tracking Inter-Core Interference

- Three sources of inter-core interference:
 - Shared cache
 - DRAM bus and bank
 - DRAM row-buffer
- InterferencePerCore bit-vector
 - Indicate whether a core is **delayed** due to inter-core interference
- Bit-vector for each source
 - Update main copy by taking union of the source bit-vectors



Cache
Bus&Bank
Row-Buffer



Cache Interference

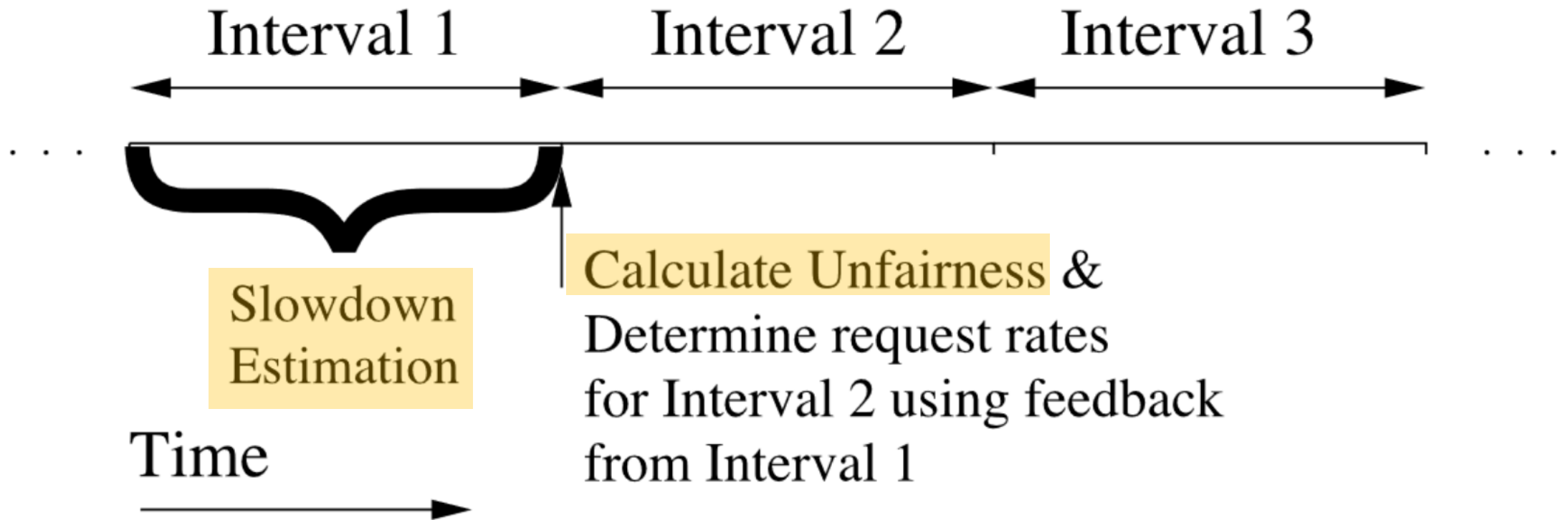
- Goal: Estimating inter-core interference on the cache by tracking **cache misses** caused by **another core**
- Pollution filter for each core
 - Bit-vector is indexed by the lower order bits of the accessed cache line address
 - A set entry in the bit-vector indicates that a cache line belonging to this core was **evicted** by **another core**
- Three steps in case of cache miss:
 - 1) on cache miss access pollution filter with the missing address and check whether **bit is set**
 - 2) set the bit in the **InterferencePerCore** vector and reset the bit in the **pollution filter**
 - 3) when the interfered-with memory **request is serviced** reset the InterferencePerCore bit

DRAM Bus & Bank Interference

- Goal: Estimate inter-core interference caused by an **inability** to access DRAM due to another core using the bus or requesting service from the bank
- This situation is easily detectable
 - If detected the corresponding InterferencePerCore bit is set
- The InterferencePerCore bit is reset when **no request** from this core is being **prevented** access to DRAM by another cores requests

DRAM Row-Buffer Interference

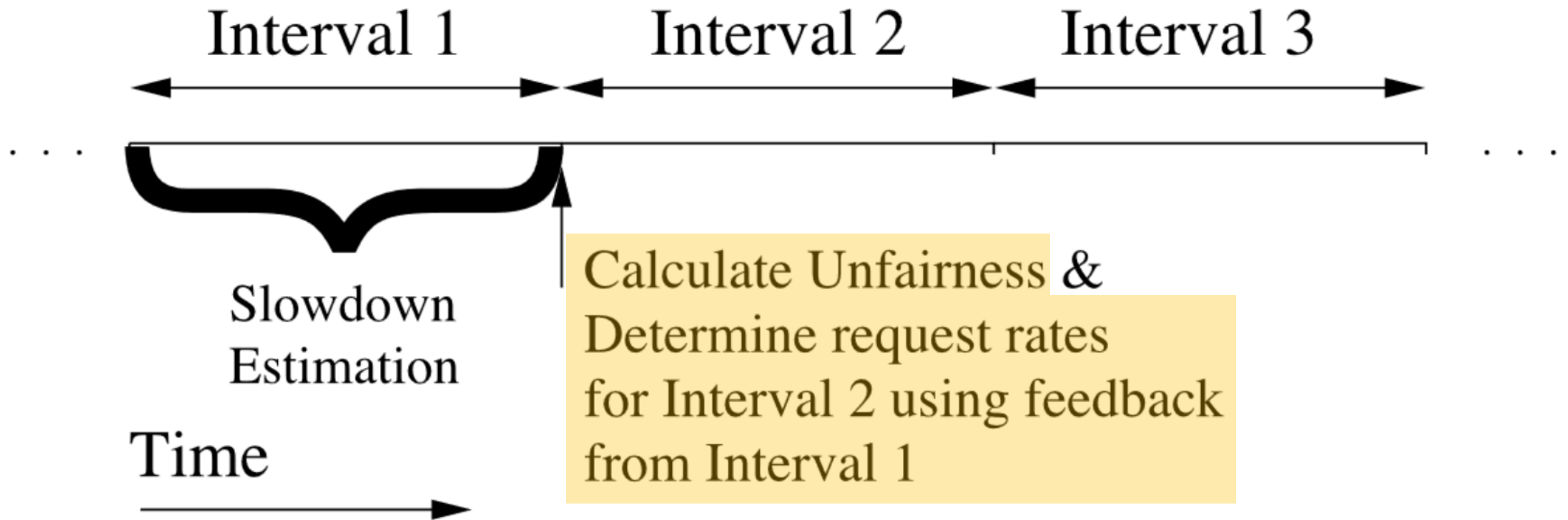
- Goal: Estimate interference caused by the **conversion** of row-buffer hits to a **miss/conflict** due to another cores memory request
- Shadow Row-Buffer Address Register for each core and for each bank
 - Whenever memory request accesses some row X, the SRAR is updated to X
- Three Steps in case of Row-Buffer miss:
 - 1) on row-buffer miss consult SRAR
 - 2) if the **SRAR** bit is set, interference is present, hence InterferencePerCore bit is set
 - 3) once the memory **request is serviced** the InterferencePerCore bit is reset



Estimation of T_{excess} for Core i

- Every cycle:
 - Check whether core i experiences interference
 - Increment T_{excess} by 1
- $T_{\text{alone}} = T_{\text{shared}} - T_{\text{excess}}$

$$IS_i = \frac{T_i^{\text{shared}}}{T_i^{\text{alone}}}, \quad \text{Unfairness} = \frac{\text{MAX}\{IS_0, IS_1, \dots, IS_{N-1}\}}{\text{MIN}\{IS_0, IS_1, \dots, IS_{N-1}\}}$$



Dynamic Request Throttling

- Check whether the estimated unfairness is **bigger** than a certain unfairness threshold
 - Throttle **down** application with the **smallest** slowdown
 - Throttle **up** application with the **largest** slowdown
- After fairness was achieved for a certain number of successive intervals:
 - Throttle up **all** applications

Throttling Mechanisms

- 1) Adjust MSHR quota
 - **MSHR quota** determines the **max. number of outstanding misses** for each core
 - Reduce the pressure by decreasing the number of concurrent request contending for service
- 2) Adjust the rate of issuing requests to the shared cache
 - Reduce number of memory **requests per unit time**
 - This allows requests from other applications to be prioritized

System Software Support

- Different Fairness Objectives:
 - The goal to be achieved by FST can be **configured** by system software (trigger condition)
- Thread Weights:
 - **Adjust priority** of different applications by applying weights
- Thread Migration and Context Switches:
 - On context switch or thread migration the corresponding interference state is cleared
 - On restart of execution, the thread starts with max. throttle and then FST dynamically adapts

Scalability to more Cores

- Each core maintains a set of $N-1$ counters, with N being the number of cores, which keep track of the inter-core interference caused by each other core
- This can be used to **identify** which core experiences the most slowdown (App_{slow}) and who of the other cores is the main contributor ($App_{interfering}$)
- Once identified, the **main contributor** will be throttled down and App_{slow} will be throttled up
- Cores other than the App_{slow} and $App_{interfering}$ are throttled up every *threshold* intervals to optimize performance

Preventing Bank Service Denial due to FR-FCFS

- FR-FCFS has the potential to **starve** application with no row-buffer locality
 - Even if the interfering application gets throttled down the problem can still exist
 - This denial of service can happen continuously
- Stop prioritizing row-buffer hits

Outline

- Background, Problem & Goal
- Novelty
- Key Approach & Ideas
- Mechanism
- **Methodology & Evaluation**
- Conclusion

- Takeaways
- Strengths and Weaknesses
- Thoughts & Ideas
- Open Discussion

System Specification

- In-house cycle-accurate x86 CMP simulator
- Faithfully model all port contention, queuing effects, bank conflicts, and other major DDR3 DRAM system constraints

Execution Core	6.6 GHz out of order processor, 15 stages, Decode/retire up to 4 instructions Issue/execute up to 8 micro instructions 256-entry reorder buffer
Front End	Fetch up to 2 branches; 4K-entry BTB 64K-entry Hybrid branch predictor
On-chip Caches	L1 I-cache: 32KB, 4-way, 2-cycle, 64B line L1 D-cache: 32KB, 4-way, 2-cycle, 64B line Shared unified L2: 1MB (2MB for 4-core), 8-way (16-way for 4-core), 16-bank, 15-cycle (20-cycle for 4-core), 1 port, 64B line size
DRAM Controller	On-chip, FR-FCFS scheduling policy [27] 128-entry MSHR and memory request buffer
DRAM and Bus	667MHz bus cycle, DDR3 1333MHz [20] 8B-wide data bus Latency: 15-15-15ns (t_{RP} - t_{RCD} - t_{CL}) 8 DRAM banks, 16KB row buffer per bank Round-trip L2 miss latency: Row-buffer hit: 36ns, conflict: 66ns

Workloads

- 18 two-application workloads from the SPEC CPU 2000/2006 benchmark
 - Two-application workloads were chosen such that **at least one** of them is highly memory-intensive

- 10 four-application workloads from the SPEC CPU 2000/2006 benchmark
 - Four-applications workloads were chosen such that **at least one** of them has high intensity and **one** has at least medium or high intensity

Metrics

- Weighted Speedup(Wspeedup):
 - IPC^{alone} is the IPC(instructions per cycle) measured when running **alone**
 - IPC^{shared} is measured while running in **tandem** with other applications
- Harmonic mean of Speedups(Hspeedup):
 - Balanced measure between fairness and system throughput

$$Wspeedup = \sum_{i=0}^{N-1} \frac{IPC_i^{shared}}{IPC_i^{alone}} \qquad Hspeedup = \frac{N}{\sum_{i=0}^{N-1} \frac{IPC_i^{alone}}{IPC_i^{shared}}}$$

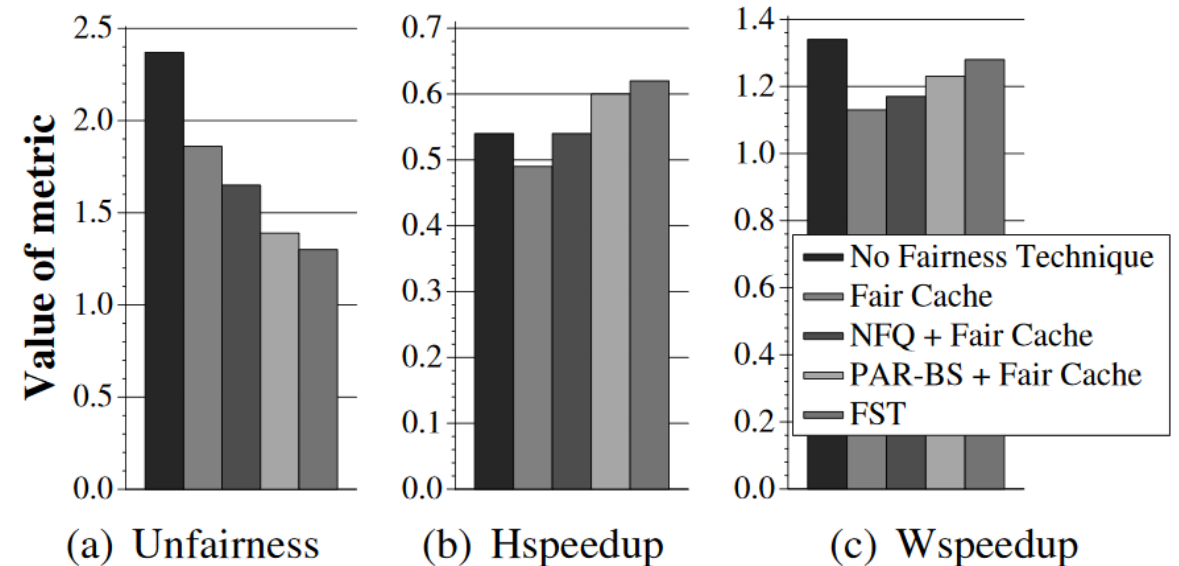
Methodology

- NoFairness:
 - Employs **no fairness techniques** in the shared memory subsystem
 - Uses LRU cache replacement and FR-FCFS
- FairCache:
 - Uses Virtual private caches technique for **fair capacity management**
- NFQ+FairCache:
 - Uses a **network fair queuing**(NFQ) fair memory scheduler combined with FairCache
- PAR-BS+FairCache:
 - Use **parallelism-aware batch scheduling** fair memory scheduler and FairCache

2-Core System Results

- All Fairness techniques degrade Wspeedup to some extent
- Unsophisticated fairness mechanisms can have a **negative** effect on system performance
- FST provides a significantly better **balance** between system fairness and performance

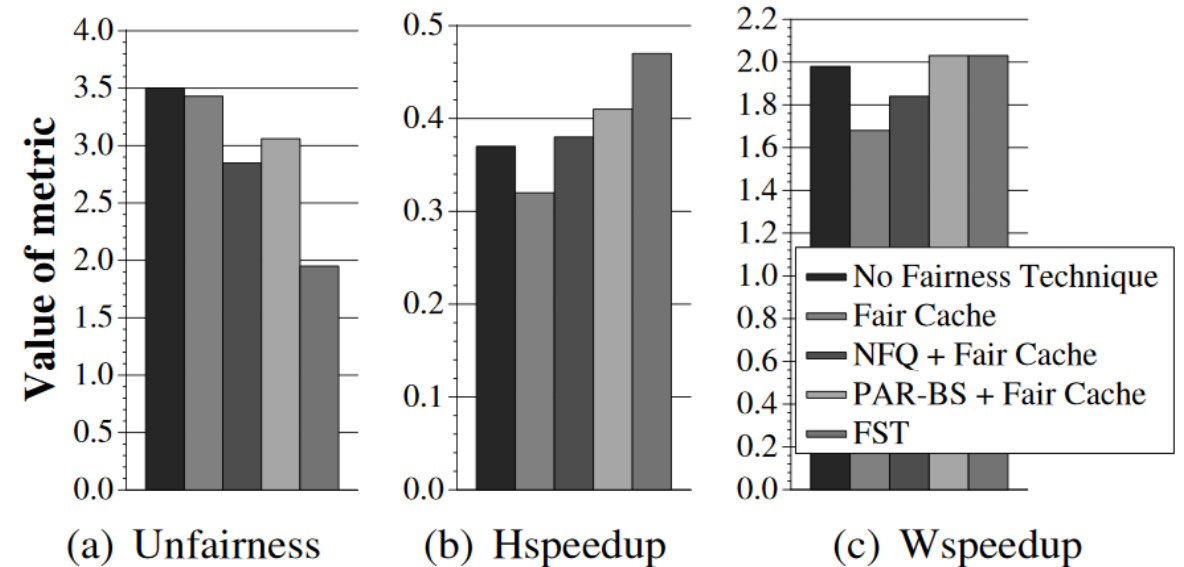
Average performance on the 2 core system



4-Core System Results

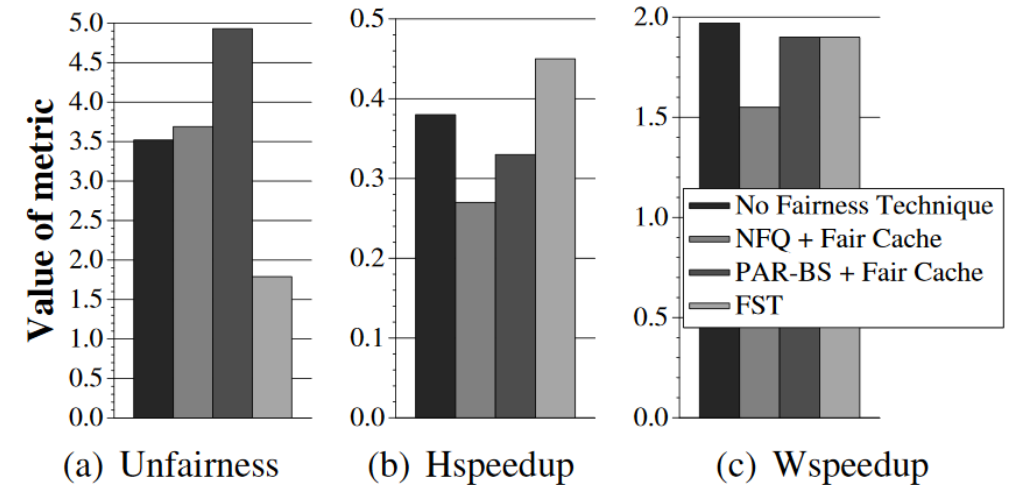
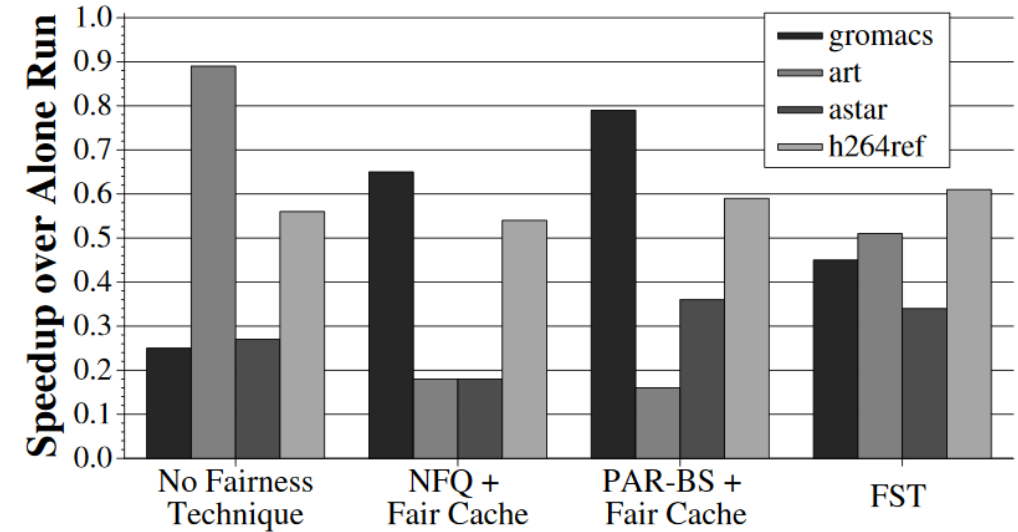
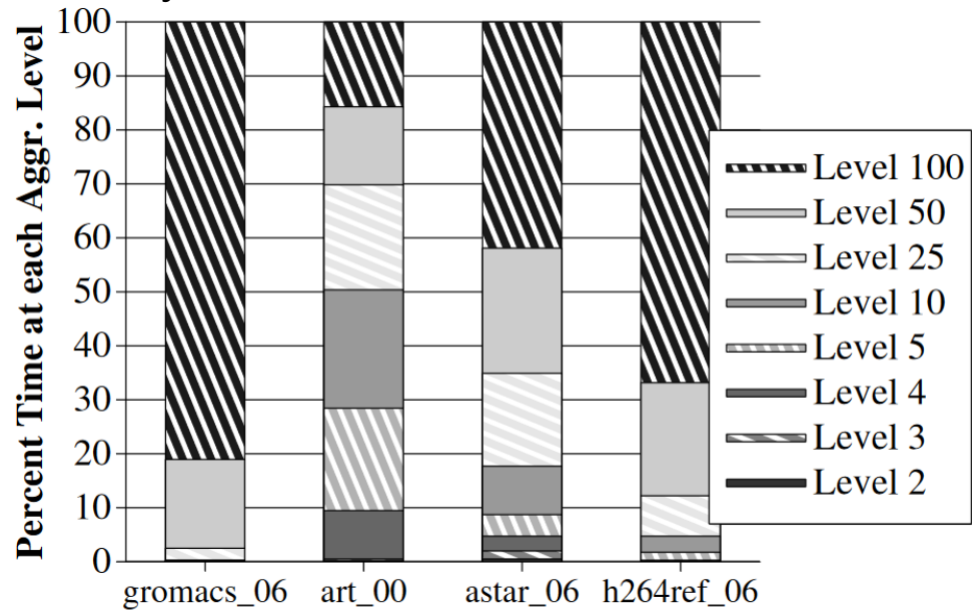
- Previous fairness mechanisms fail to improve system fairness significantly
 - Prioritize nonintensive applications regardless of whether or not those experience slowdown
- FST is the best technique for system fairness and Hspeedup, while not falling behind in Wspeedup

Average performance on the 4 core system



Case Study

- Art and Astar are memory intensive:
 - These are slowed down **too much** by NFQ+FairCache and PAR-BS+FairCache, causing high unfairness
 - Inability to detect when slowdown is caused



Hardware Cost

	Cost for N cores	Cost for N = 4
<i>ExcessCycles</i> counters	$N \times N \times 16$ bits/counter	256 bits
Interference pollution filter per core	$2048 \text{ entries} \times N \times$ $(1 \text{ pollution bit} + (\log_2 N) \text{ bit processor id})/\text{entry}$	24,576 bits
<i>InterferingCoreId</i> per MSHR entry	$32 \text{ entries/core} \times N \times 2 \text{ interference sources} \times (\log_2 N) \text{ bits/entry}$	512 bits
<i>InterferencePerCore</i> bit-vector	$3 \text{ interference sources} \times N \times N \times 1 \text{ bit}$	48 bits
Shadow row-buffer address register	$N \times \# \text{ of DRAM banks (B)} \times 32 \text{ bits/address}$	1024 bits
<i>Successive Fairness Achieved Intervals</i> counter <i>Intervals To Wait To Throttle Up</i> counter per core <i>Inst Count Each Interval</i> per core	$(2 \times N + 1) \times 16$ bits/counter	144 bits
Core id per tag store entry in K MB L2 cache	$16384 \text{ blocks/Megabyte} \times K \times (\log_2 N) \text{ bit/block}$	65,536 bits
Total hardware cost for N-core system	Sum of the above	92092 = 11.24 KB
Percentage area overhead (as fraction of the baseline K MB L2 cache)	$\text{Sum (KB)} \times 100 / (K \times 1024)$	11.24KB/2048KB = 0.55%

Hardware cost of FST on a 4-core CMP system

Outline

- Background, Problem & Goal
- Novelty
- Key Approach & Ideas
- Mechanism
- Methodology & Evaluation
- **Conclusion**

- Takeaways
- Strengths and Weaknesses
- Thoughts & Ideas
- Open Discussion

Executive Summary & Conclusion

- **Motivation:** Cores in a chip-multiprocessor system share multiple hardware resources in the memory subsystem
 - Interference in the shared resources can lead to unfair slowdown for some applications
- **Problem:** Existing fairness mechanisms focus on a single resource
 - Multiple independently implemented mechanisms can make contradictory decisions, leading to low fairness and loss of performance
- **Goal:** provides fairness in the entire shared memory system without degrading performance
- **Key Contributions:** Fairness via Source Throttling(FST) provides two major mechanisms
 - 1) Runtime fairness evaluation
 - 2) Dynamic request throttling
- **Major Results:** improve performance by 25.6%/14.5% and reduce unfairness by 44%/36%

Outline

- Background, Problem & Goal
- Novelty
- Key Approach & Ideas
- Mechanism
- Methodology & Evaluation
- Conclusion

- **Takeaways**
- Strengths and Weaknesses
- Thoughts & Ideas
- Open Discussion

Takeaways

- In order to ensure good performance for multiple applications in a shared system, controlling system-wide fairness is necessary
- By implementing FST one can decrease system complexity, due to the fact that no more coordination between multiple fairness techniques is needed

Outline

- Background, Problem & Goal
- Novelty
- Key Approach & Ideas
- Mechanism
- Methodology & Evaluation
- Conclusion

- Takeaways
- **Strengths and Weaknesses**
- Thoughts & Ideas
- Open Discussion

Strengths

- A new approach to an old problem, which will only **get worse** with rising core counts
- In addition to improving **system-wide fairness** it also provides comparable or superior **performance** compared to prior fairness mechanisms
- Reduce system **complexity** by replacing multiple resource-based mechanisms with FST
- FST can accomplish multiple **different** fairness objectives
- The evaluation provides a good overview, while the case study provides more insight
- It is well written

Weaknesses/Limitations

- **False** positive and negative in the pollution filter
- Implementation cost of FST may **scale** poorly since the number of cores directly determines the cost
- Diminishing returns on a system with a lot of **thread migration** and **context switches**
- The optimal unfairness threshold mentioned in the paper might be hard to find

Outline

- Background, Problem & Goal
- Novelty
- Key Approach & Ideas
- Mechanism
- Methodology & Evaluation
- Conclusion

- Takeaways
- Strengths and Weaknesses
- **Thoughts & Ideas**
- Open Discussion

Thoughts and Ideas

- Interval-based estimation and throttling
 - What impact will an application with rapidly and randomly changing memory intensity have?
- Aggressiveness levels
 - Would it make sense to have throttle function based on slowdown instead of fixed levels?
- Security aspects are not evaluated
 - Could a single or a group of bad actors attack FST?

Outline

- Background, Problem & Goal
- Novelty
- Key Approach & Ideas
- Mechanism
- Methodology & Evaluation
- Conclusion

- Takeaways
- Strengths and Weaknesses
- Thoughts & Ideas
- **Open Discussion**

Discussion starters

- Will the problem become more important over time?
- Are there situations where FST might not work?
- Do you think the increase in cost due to higher bank and core counts will be overshadowed by the increase in performance?
- Can you think of some disadvantages that I missed or even some way of improving FST?

Backup Slides

Algorithm 2 Estimation of T_{excess} for core i **Every cycle**

if *inter-core cache or DRAM bus or DRAM bank or DRAM row-buffer interference* **then**
 set *InterferencePerCore* bit i
 set *InterferingCoreId* in delayed memory request
end if
if *InterferencePerCore* bit i is set **then**
 Increment *ExcessCycles* for core i
end if

1: check whether some sort of inter-core interference is present
 2: if increment the ExcessCycles counter

Every L2 cache fill for a miss due to interference OR**Every time a memory request which is a row-buffer miss due to interference is serviced**

reset *InterferencePerCore* bit of core i
InterferingCoreId of core $i = i$ (no interference)

1: whenever a interfered with memory request is serviced reset *InterferencePerCore* bit and set *InterferingCoreId* of core i to i
 2: whenever a memory request is scheduled and also has no other request waiting on any bank busy servicing another core

Every time a memory request is scheduled to DRAM

if Core i has no requests waiting on any bank which is busy servicing another core j ($j \neq i$) **then**
 reset *InterferencePerCore* bit of core i
end if

Dynamic request throttling

Algorithm 1 Dynamic Request Throttling

```
if Estimated Unfairness > Unfairness Threshold then  
  Throttle down application with the smallest slowdown  
  Throttle up application with the largest slowdown  
  Reset Successive Fairness Achieved Intervals  
  
else  
  if Successive Fairness Achieved Intervals = threshold then  
    Throttle all applications up  
    Reset Successive Fairness Achieved Intervals  
  else  
    Increment Successive Fairness Achieved Intervals  
  end if  
end if
```

- This is a simplified version for dual cores
- After a certain number of fair intervals both cores are allowed to throttle up

Algorithm 3 Dynamic Request Throttling - General Form

```

if Estimated Unfairness > Unfairness Threshold then
  Throttle down application that causes most interference
  (Appinterfering) for application with largest slowdown
  Throttle up application with the largest slowdown (Appslow)
  Reset Successive Fairness Achieved Intervals
  Reset Intervals To Wait To Throttle Up for Appinterfering.

```

1) Responsible for throttling down the most interfering application

```

// Preventing bank service denial

```

```

if Appinterfering throttled lower than Switchthr AND causes
greater than Interferencethr amount of Appslow's total interference
then

```

2) Solving bank service denial due to FR-FCFS

```

  Temporarily stop prioritizing Appinterfering due to row hits in
  memory controller

```

```

end if

```

```

if AppRowHitNotPrioritized has not been Appinterfering for
SwitchBackthr intervals then

```

```

  Allow it to be prioritized in memory controller based on row-buffer
  hit status of its requests

```

```

end if

```

```

for all applications except Appinterfering and Appslow do

```

3) Throttling up all applications that are neither App_{slow} nor $App_{interfering}$ every *threshold1* intervals

```

  if Intervals To Wait To Throttle Up = threshold1 then

```

```

    throttle up

```

```

    Reset Intervals To Wait To Throttle Up for this app.

```

```

  else

```

```

    Increment Intervals To Wait To Throttle Up for this app.

```

```

  end if

```

```

end for

```

```

else

```

```

  if Successive Fairness Achieved Intervals = threshold2 then

```

```

    Throttle up application with the smallest slowdown

```

```

    Reset Successive Fairness Achieved Intervals

```

```

  else

```

```

    Increment Successive Fairness Achieved Intervals

```

```

  end if

```

```

end if

```

4) Throttling up ??? application after number of *threshold2* intervals

FST Parameter used in the evaluation

- We use 8 different aggressiveness levels:
 - 2%, 3%, 4%, 5%, 10%, 25%, 50% and 100%

<i>Fairness Threshold</i>	<i>Successive Fairness Achieved Intervals Threshold</i>	<i>Intervals Wait To Throttle Up</i>	<i>Interval Length</i>
1.4	4	2	25Kinsts
<i>Switch_{thr}</i>		<i>Interference_{thr}</i>	<i>SwitchBack_{thr}</i>
5%		70%	3 intervals