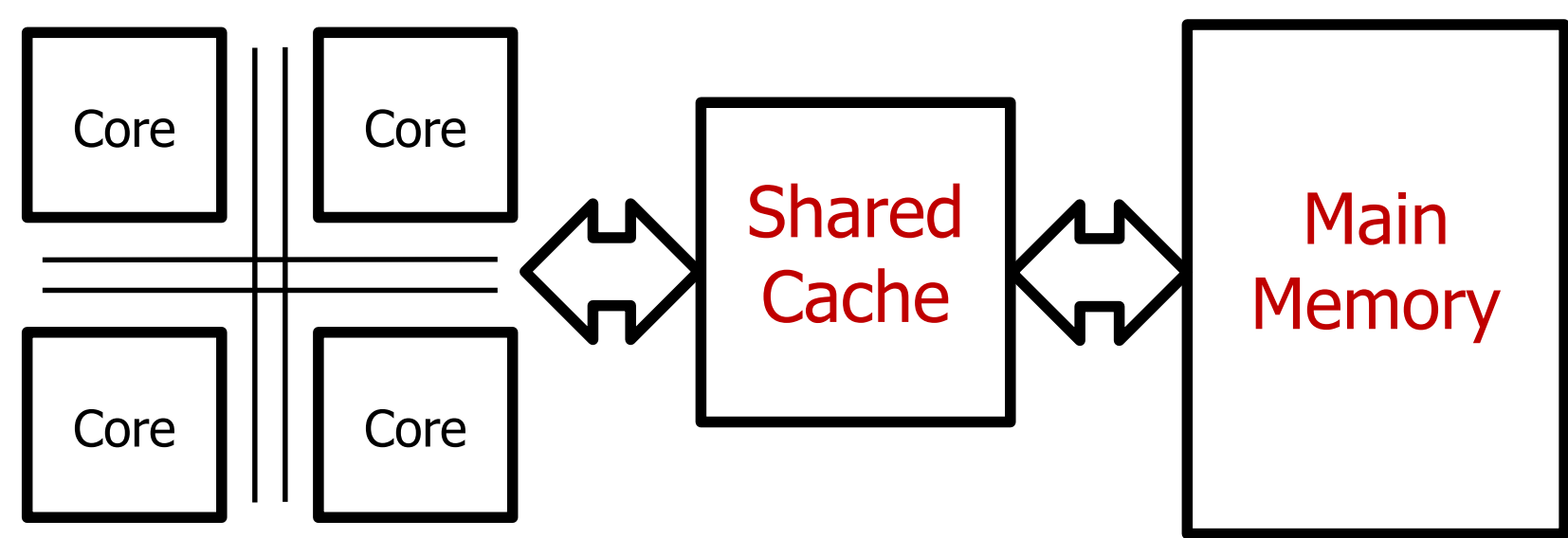


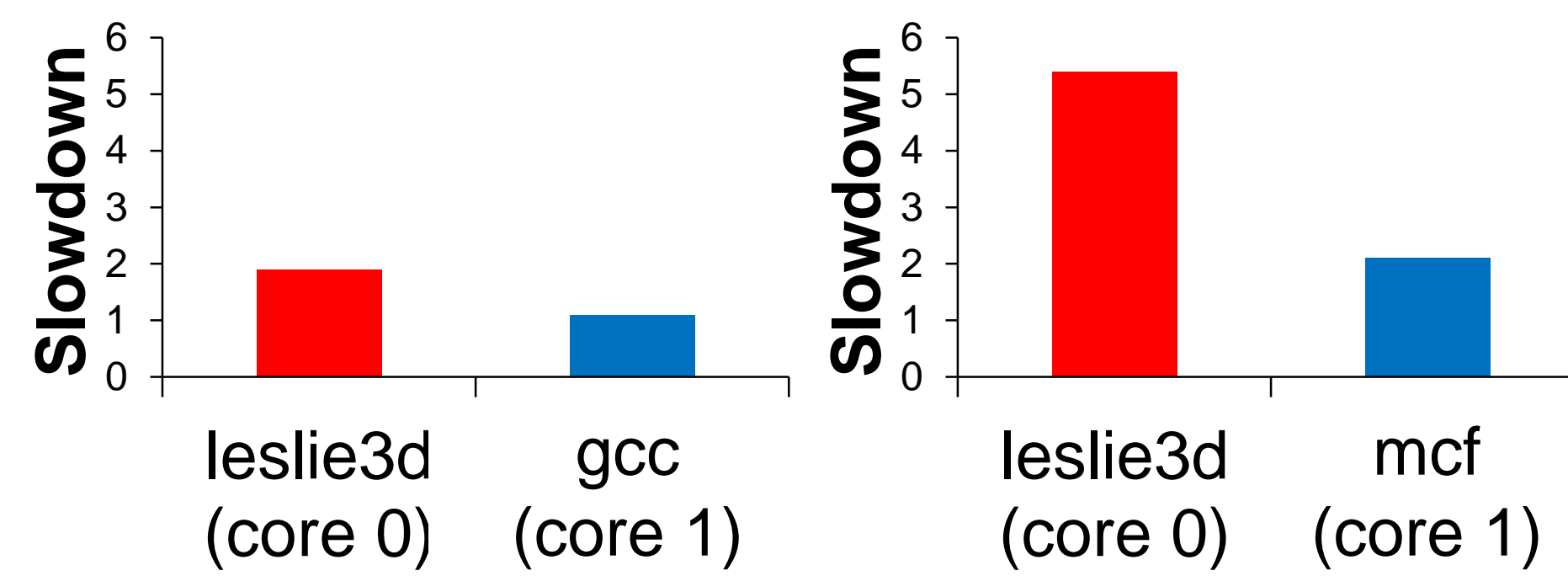
Application Slowdown Model: Quantifying and Controlling Impact of Interference at Shared Caches and Main Memory

Lavanya Subramanian, Vivek Seshadri, Arnab Ghosh, Samira Khan, Onur Mutlu
Carnegie Mellon University, Intel Labs, University of Virginia

Problem:
Shared Resource Interference



Impact of
Shared Resource Interference



1. High application slowdowns
2. Unpredictable application slowdowns

Our Goal

Provide high and predictable performance in the presence of shared resource interference

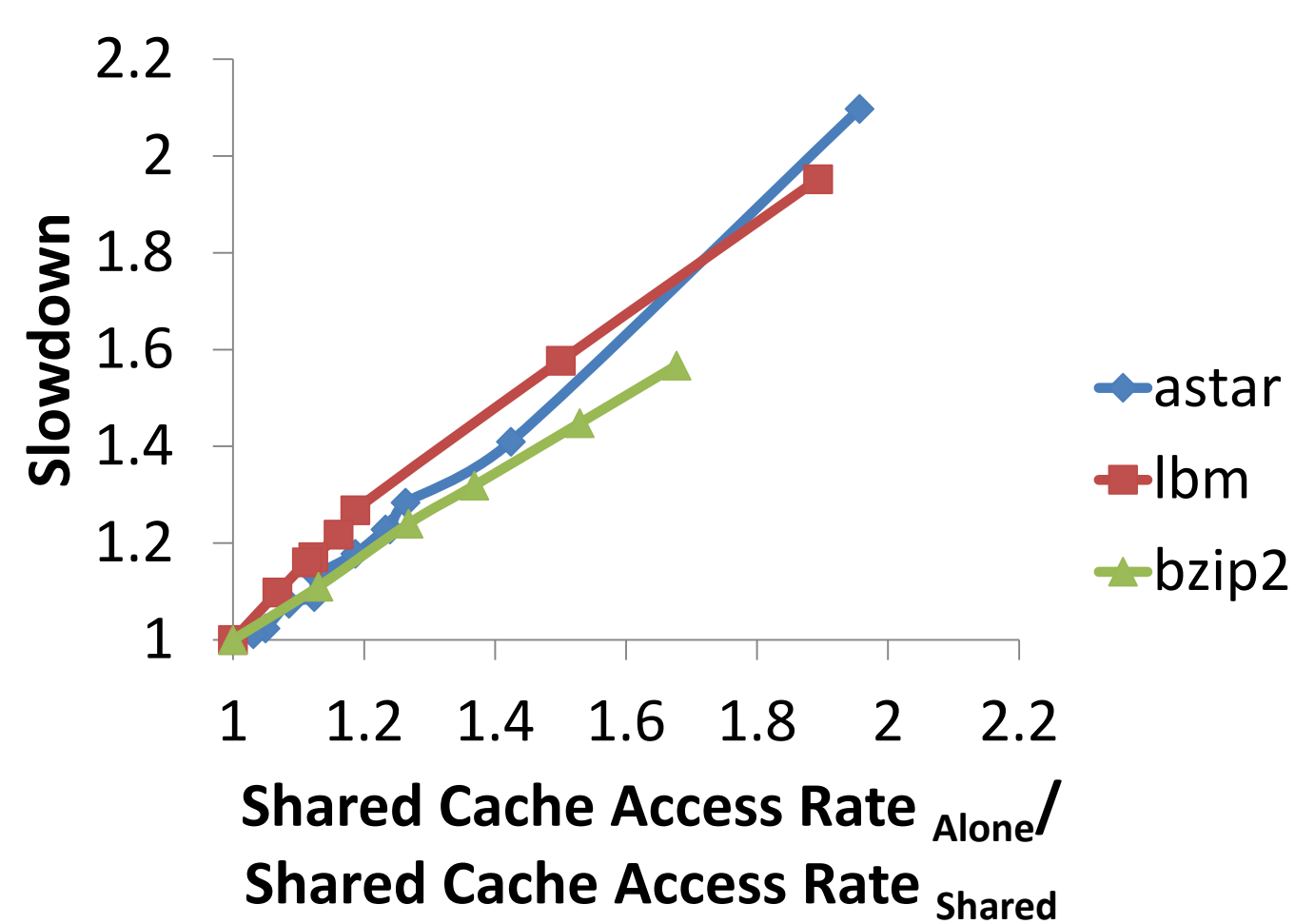
Our Approach

1. Build a model to estimate slowdowns
2. Leverage our model for slowdown-aware resource management

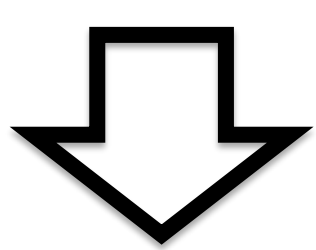
Application Slowdown Model (ASM)

Observation: Proxy for Performance

For a memory bound application,
Performance \propto Cache access rate



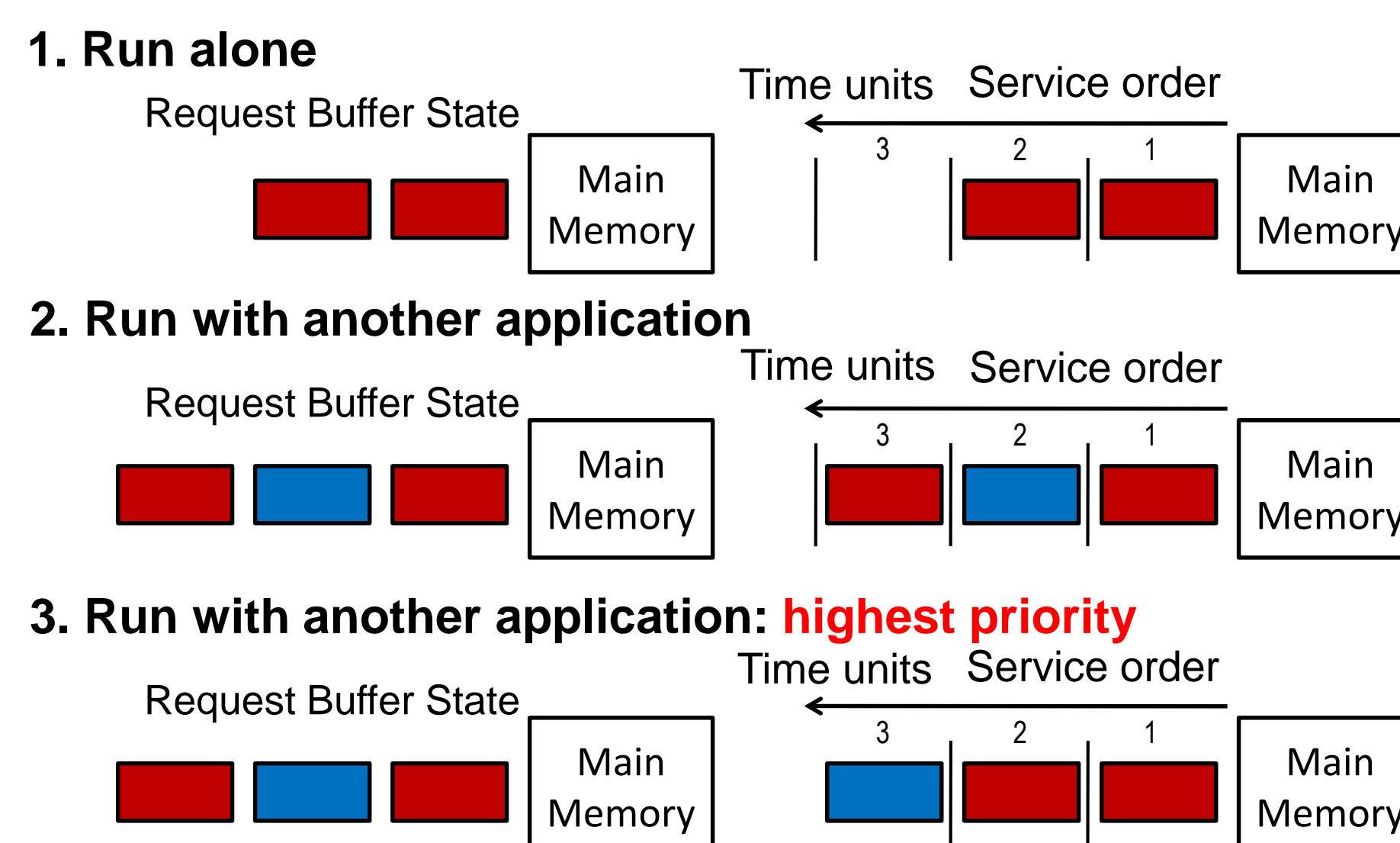
$$\text{Slowdown} = \frac{\text{Performance}_{\text{Alone}}}{\text{Performance}_{\text{Shared}}}$$



$$\text{Slowdown} = \frac{\text{Cache Access Rate}_{\text{Alone}} (\text{CAR}_{\text{Alone}})}{\text{Cache Access Rate}_{\text{Shared}} (\text{CAR}_{\text{Shared}})}$$

Challenge: Estimating Cache Access Rate Alone

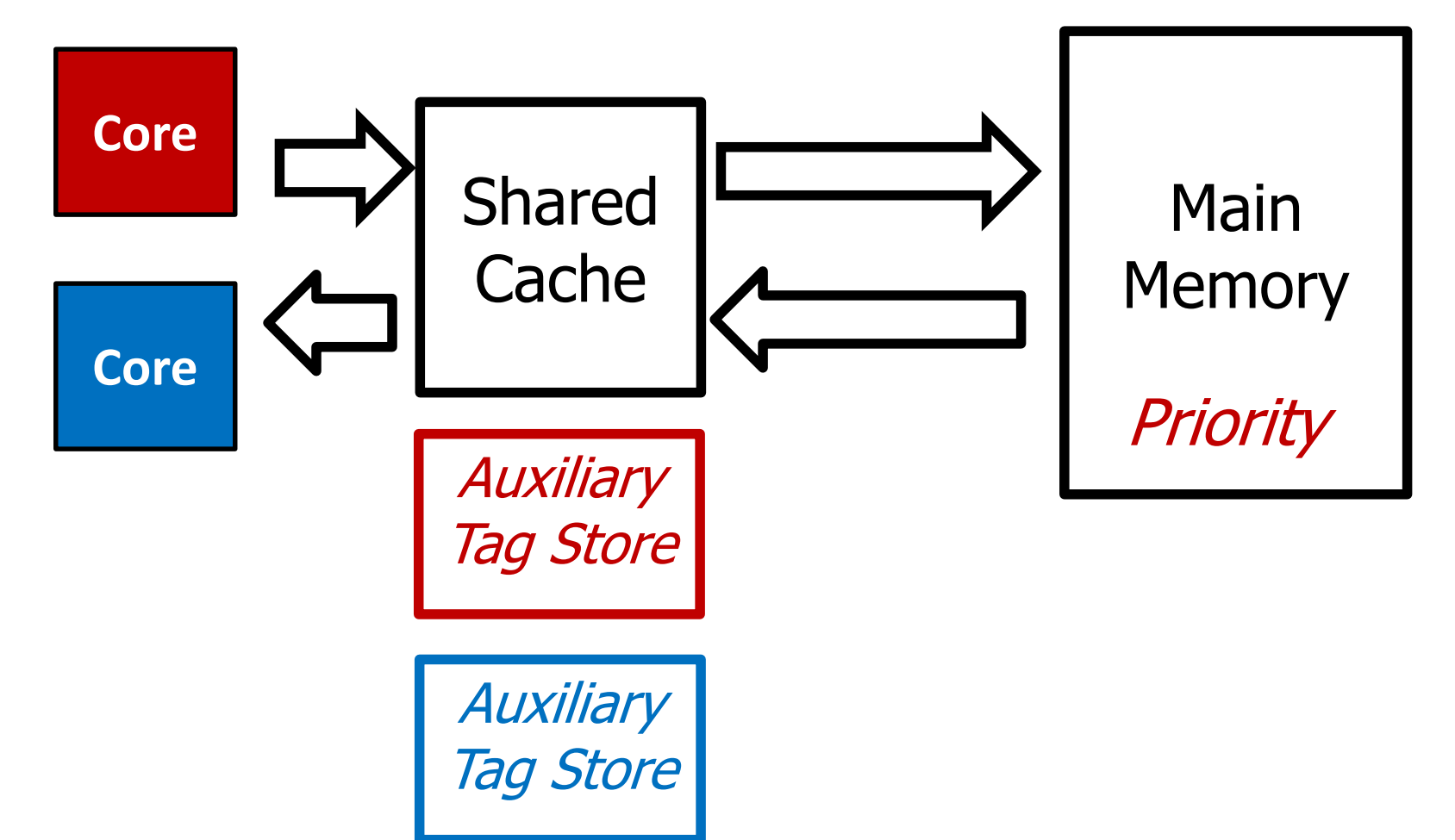
Minimize memory bandwidth contention:
Using priority (Subramanian et al., HPCA 2013)



Highest priority \rightarrow Little interference
(almost as if application were run alone)

Enables estimation of **miss service time**

Quantify shared cache capacity contention:
Using auxiliary tag stores (Pomerene et al., 1989)



Auxiliary tag store counts **#contention misses**

$$\text{Cache Contention Cycles} = \# \text{Contention Misses} \times \text{Average Miss Service Time}$$

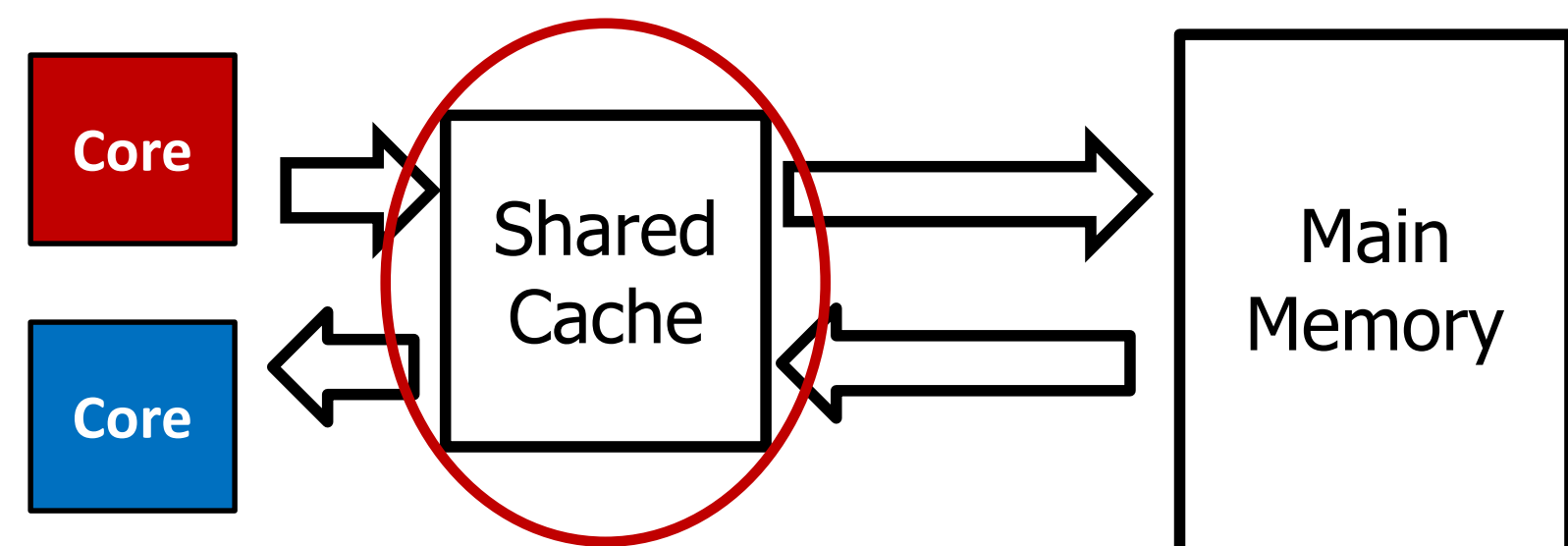
From auxiliary tag store when given high priority
Measured when given high priority

Remove contention cycles when estimating $\text{CAR}_{\text{Alone}}$

Average error of ASM: 10%; Average error of previous models: 30%

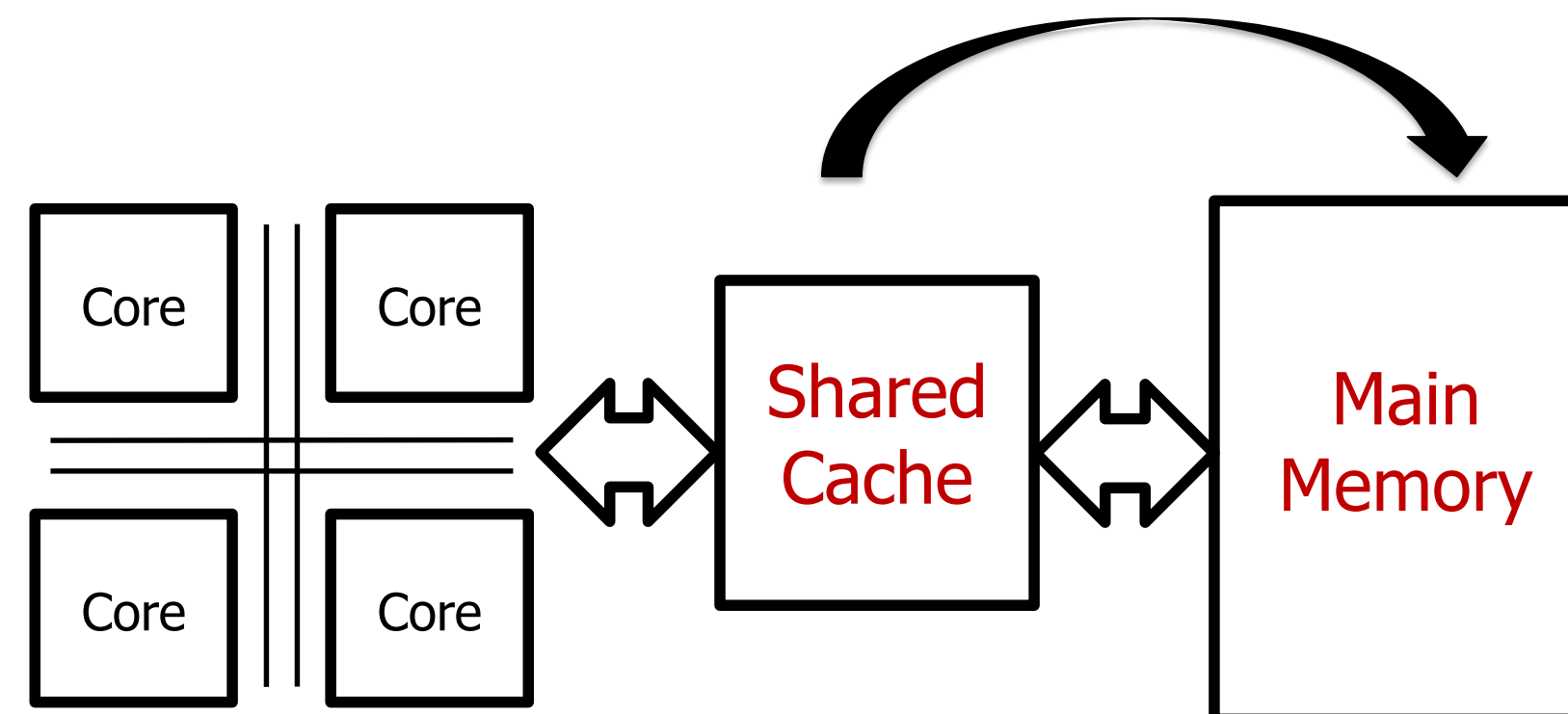
Leveraging the Application Slowdown Model

Slowdown-aware cache capacity partitioning



Previous work: Reduce miss counts;
Our proposal: Reduce slowdowns

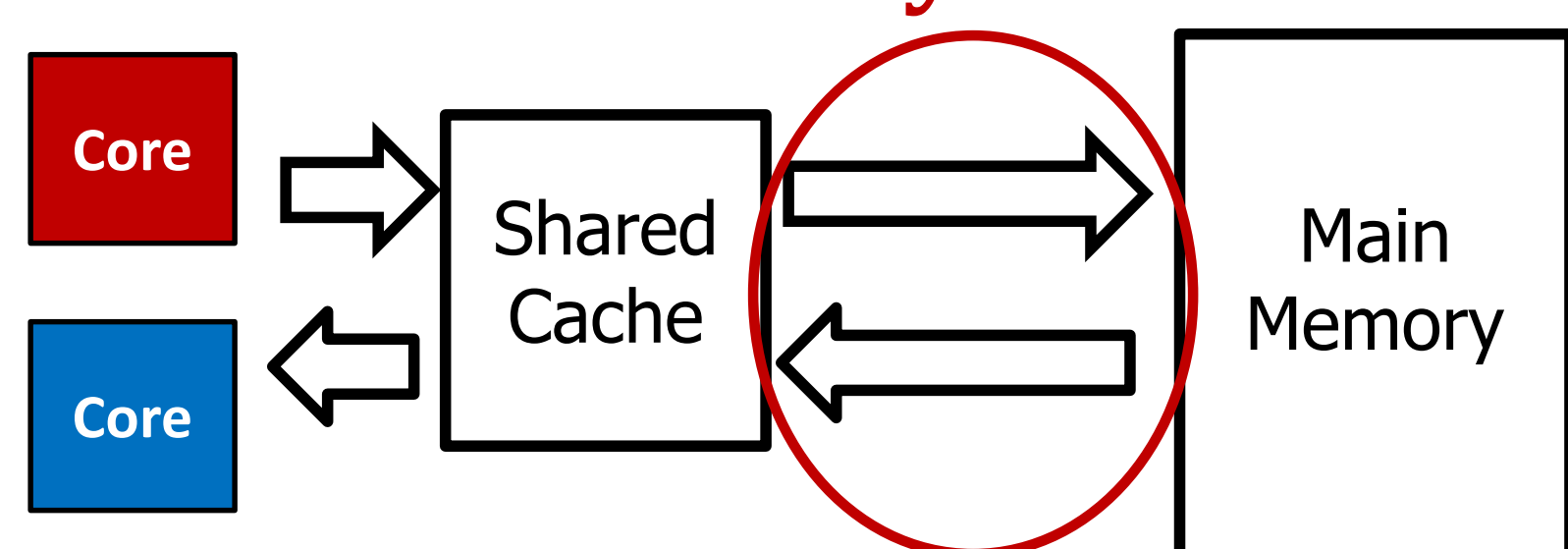
Coordinated Resource Allocation Schemes



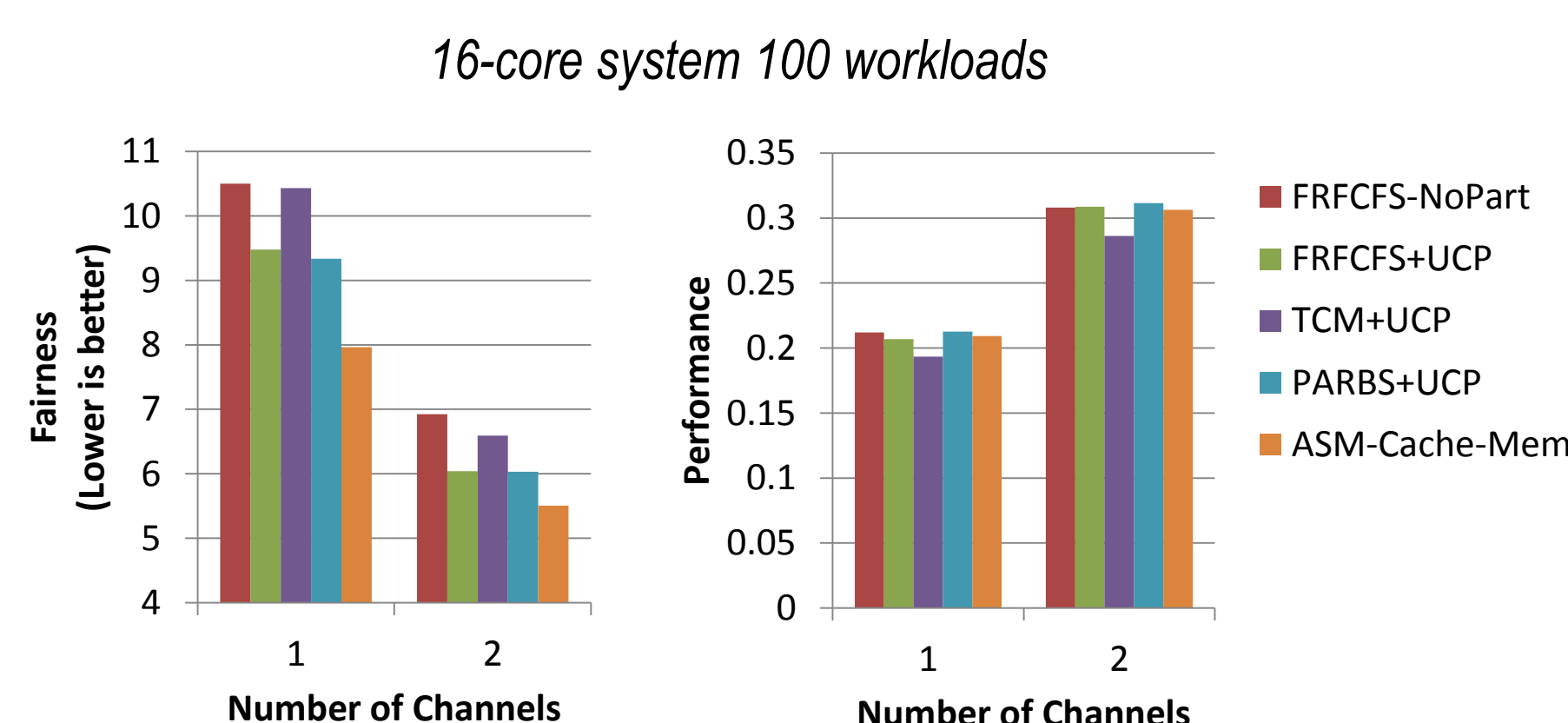
Providing Slowdown Guarantees

- Cache allocation with the goal of meeting a slowdown bound
- Allocate just enough cache space to critical application
- Allocate remaining cache space to other applications

Slowdown-aware memory bandwidth partitioning



Allocation memory bandwidth proportional to slowdowns



Significant fairness benefits across different channel counts

