CAGE: A Contention-Aware Game-theoretic Model for Heterogeneous Resource Assignment

Diman Zad Tootaghaj Farshid Farhat

The Pennsylvania State University (US)



Why distributed resource management approach?



Distributed

Centralized

Distributed Resource Management

Observations:

Centralized Approaches:

As the **number of applications increase**, centralized approaches are not **scalable**.

The central decision maker does not know about the **applications' need**. Applications have **different** resource needs. Applications' needs **changes** during time.

<u>Key Idea</u>: Using a game theoretic resource allocations, where applications compete for the shared resources.

<u>Results:</u>

We show that **CAGE** is **strategy proof**:

• No application can get more utilization by bidding more or less than the true value of the resource

Spec 2006 Cache Demand



Applications' phases

Consider hmmer and mcf applications:



[1] M. K. Qureshi et al. "Utility-based cache partitioning: A low-overhead, highperformance, runtime mechanism to partition shared caches." *in MICRO-39.* IEEE, 2006.

Distributed Resource Management

Two Case Studies:

- 1) A Cache Congestion game
- 2) Main Processor and Co-processor Congestion game.

Problem Definition:

Given N Players (applications) and a set of M resources, we want to find an

- allocation scheme that maximizes the total utility of all players.
- Game theoretic model: Each player *i* chooses a subset of resources (from a given
- family of subsets), each resource m has a utility gain u_m which depends on the
- number of players using this resource.
- Each player wants to maximize its own utility.

CAGE Model



Cache Congestion Game

Cache Partitioning:

We have a hierarchical cache where larger chunks are more congested and smaller chunks are more private. The applications can decide to run their code on the congested part with more cache space or less congested part with less cache capacity.

Utility-based cache partitioning:

Shows that LRU-based cache partitioning gives more cache space to applications that have higher demand and lower to those who have lower demand. But higher demand doesn't mean higher performance (streaming applications).

Evaluation

Cache Partitioning:

Comparison of CAGE with solo, shared on different mixes of applications:



Conclusion

Observations:

Centralized Approaches:

As the **number of applications increase**, centralized approaches are not **scalable**. The central decision maker does not know about the **applications' need**. Applications have **different** resource needs. Applications' needs **changes** during time.

Results:

We show that **CAGE** is **strategy proof**:

No application can get more utilization by bidding more or less than the true

value of the resource

Questions?



Hidden Slides

Spec 2006 Cache Demand

IPC with respect to the size of LLC:



Example (hmmer, mcf)



Static approaches: 27.8% Improvement CAGE: 36.84/2+27.8/2>27.8

Parallel

• Consider the following bipartite graph and utility functions for each resource:



Sequential Auction:

• Parallel actions of each agent: Each user bids for the most valuable resource.



Sequential Auction:

• Parallel actions of each agent:



Sequential Auction:

• Parallel actions of each agent:

