

SNIPER: EXPLORING THE LEVEL OF ABSTRACTION FOR SCALABLE AND ACCURATE PARALLEL MULTI-CORE SIMULATION

TREVOR E. CARLSON, WIM HEIRMAN, LIEVEN EECKHOUT

TREVOR.CARLSON@ELIS.UGENT.BE HTTP://WWW.ELIS.UGENT.BE/~TCARLSON WEDNESDAY, NOVEMBER 16TH, 2011 SC11, SEATTLE, WA



NEED FOR HIGH-ABSTRACTION MODELS



DEMANDS ON SIMULATION ARE INCREASING

- Number of cores per node is increasing (100s)
 Requiring parallel simulation to keep pace
- Cache sizes are increasing
 - Requires
 longer
 benchmark
 runtimes



LLC Cache Sizes

TIME TO MARKET MATTERS

Industry

- Tend to have enough compute bandwidth
- Simulation latency limited
- Results should be ready when you arrive at work in the morning

Academia

- Not enough compute power
- Both latency and bandwidth limited
- Exploration and investigation work

• How many scenarios can I run?

- N = total number of simulation scenarios
- d = days until SC deadline
- t = average time per simulation
- B = number of benchmarks
- A = number of architectures

minimize t to maximize N

FAST AND ACCURATE SIMULATION IS NEEDED

Sniper Simulator

- Interval core model
- Accurate structures (caches, branch predictors, etc.)
- Parallel simulator scales with the number of simulated cores

Key Questions

- What is the right level of abstraction?
- When to use these abstraction models?



OUTLINE

- Motivation
- Research Background
- Experimental Setup
- Results
- Conclusions

NEEDED DETAIL DEPENDS ON FOCUS



ONE-IPC MODELING - TOO SIMPLE?

- Simple high-abstraction model
 - Widely used
- One-IPC modeling
 - Scalar, in-order issue
 - Account for non-unit instruction exec latencies
 - Perfect branch prediction
 - L1 D-cache accesses are assumed to be hidden
 - All other cache accesses incur penalty

INTERVAL SIMULATION

 Out-of-order core performance model with inorder simulation speed



D. Genbrugge et al., HPCA'10 S. Eyerman et al., ACM TOCS, May 2009 T. Karkhanis and J. E. Smith, ISCA'04, ISCA'07 9

KEY BENEFITS OF THE INTERVAL MODEL

- Models superscalar OOO execution
- Models impact of ILP
- Models second-order effects: MLP

• Allows for constructing CPI stacks

MULTI-CORE INTERVAL SIMULATION





Instantaneous dispatch rate is determined by the longest critical path in the old window:

Instantaneous dispatch rate = min (W / L, D) Little's law

Assumes a balanced architecture

L = longest critical path length in cycles W = instructions in the old window (max = ROB length) D = maximum dispatch rate (processor width)

LONG BACK-END MISS EVENTS Isolated long-latency load



S. Eyerman et al., ACM TOCS, May 2009

LONG BACK-END MISS EVENTS OVERLAPPING LONG-LATENCY LOADS



S. Eyerman et al., ACM TOCS, May 2009

CORE TIMING LONG-LATENCY LOAD



If long-latency load (LLC miss):

core sim time += miss-latency

AND walk the window to issue independent miss events: these are hidden under the long-latency load – second-order effects

CYCLE STACKS

- Where did my cycles go?
- CPI stack
 - Cycles per instruction
 - Broken up in components
- Normalize by either
 - Number of instructions (CPI stack)
 - Execution time (time stack)
- Different from miss rates: cycle stacks directly quantify the effect on performance

CPL

CONSTRUCTING CPI STACKS

- Interval simulation: track why time is advanced
 - No miss events
 - Issue instructions at base CPI
 - Increment base component
 - Miss event
 - Fast-forward time by X cycles
 - Increment component by X



CPI

EXPERIMENTAL SETUP

Benchmarks

- Complete SPLASH-2 suite
 - 1 to 16 threads
 - Linux pthreads API
- Graphite simulation infrastructure
 - Fast simulation
 - Parallel and modular base
 - Pin-based user-level functional-first simulation

EXPERIMENTAL SETUP: ARCHITECTURE



SNIPER SIMULATION ENVIRONMENT

- User-level, x86-64, parallel
- Branched from MIT Graphite in August 2010
- Adds interval core model, CPI stacks, modern branch predictor, shared cache models, DVFS, OpenMP support, etc.
- Hardware-validated against a 16-core Intel Xeon X7460 Dunnington machine



SNIPER SIMULATION ENVIRONMENT

- Code will be available this week
- Open source (MIT license, interval model with academic license) available at

http://snipersim.org



SIMULATION CONSIDERATIONS

Is a one-IPC core model accurate enough?

• What are the available abstraction levels?

• What is the speed/accuracy trade-off for each level?

INTERVAL PROVIDES NEEDED ACCURACY



INTERVAL: GOOD OVERALL ACCURACY



16 cores

SIMULATION PERFORMANCE



SYNCHRONIZATION VARIABILITY



Execution time variability

Variability due to relaxed synchronization is application specific

FLEXIBILITY TO CHOOSE NEEDED FIDELITY



CONCLUSIONS

- Faster simulation methodologies are needed to bridge the gap to 100s and 1000s of cores per socket
- Many performance/accuracy trade-offs, and application and optimization-target specific
- The interval model is a good option for fast and accurate core simulation
 - Simulated vs. HW error of only 25% for 16-thread applications
- Download the Sniper simulator soon at snipersim.org





SNIPER: EXPLORING THE LEVEL OF ABSTRACTION FOR SCALABLE AND ACCURATE PARALLEL MULTI-CORE SIMULATION



TREVOR.CARLSON@ELIS.UGENT.BE HTTP://WWW.ELIS.UGENT.BE/~TCARLSON WEDNESDAY, NOVEMBER 16TH, 2011 SC11, SEATTLE, WA



