# Dynamic Performance Tuning for Speculative Threads

Yangchun Luo, Venkatesan Packirisamy, Nikhil Mungre<sup>†</sup>, Ankit Tarkas<sup>†</sup>, Wei-Chung Hsu, and Antonia Zhai

> Dept. of Computer Science and Engineering <sup>†</sup>Dept. of Electronic Computer Engineering University of Minnesota – Twin Cities



### Motivation

Multicore processors brought forth large potentials of computing power



Intel Core i7 die photo

Exploiting these potentials demands thread-level parallelism

# **Exploiting Thread-Level Parallelism**



# But Unpredictable

ersity of Minnesota

# Addressing Unpredictability

State-of-the-art approach:

- Profiling-directed optimization
- Compiler analysis

#### A typical compilation framework







#### Hard to model TLS overhead statically, even with profiling





#### Profile-based decision may not be appropriate for actual input





Behavior of the same code changes over time





### Impact on Cache Performance





Cache impact is difficult to model statically

# Our proposal



# Parallelization Decisions at Runtime



New execution model facilitates dynamic tuning

# **Evaluating Performance Impact**





Simple metric: cost of speculative failure

Comprehensive evaluation: sequential performance prediction

# **Sequential Performance Prediction**



📇 University of Minnesota

#### Cache Behavior: Scenario #1



#### **Count** the miss in squashed execution



#### Cache Behavior: Scenario #2



### Tune the performance





# How to prioritize the search



# **Evaluation Infrastructure**

#### **Benchmarks**

• SPEC2000 written in C, -O3 optimization

#### **Underlying architecture**

- 4-core, chip-multiprocessor (CMP)
- speculation supported by coherence

#### Simulator

- Superscalar with detailed memory model
- simulates communication latency
- models bandwidth and contention

#### > Detailed, cycle-accurate simulation



# **Comparing Tuning Policies**



Parallel Code Overhead



# Profile-Based vs. Dynamic



Dynamic outperformed static by  $\approx 10\%$  (1.37x/1.25x)

Potentially go up to  $\approx 15\%$  (1.46x/1.27x)



### **Related works: Region Selection**



Dynamic Performance Tuning [Luo ISCA'09] (this work)



## Conclusions

Deciding how to extract speculative threads at runtime **Quantitatively** summarize performance profile **Compiler hints** avoid suboptimal decisions

Compiler and hardware work together. □ 37% speedup over sequential execution □ ≈10% speedup over profile-based decisions

> **Configurable HPMs enable efficient dynamic optimizations**

Dynamic performance tuning can improve TLS efficiency

