

Adaptive Write-Update and Write-Invalidate Cache Coherence Protocols for Producer-Consumer Sharing

Bangjie Liu

Language Technologies Institute
School of Computer Science
Carnegie Mellon University
bangjiel@andrew.cmu.edu

Hao Li

Language Technologies Institute
School of Computer Science
Carnegie Mellon University
haol2@andrew.cmu.edu

Abstract—Shared memory multicore systems play crucial roles in scientific and enterprise applications. They are efficient in general applications but perform poorly in applications that establish producer-consumer sharing patterns under write-invalidate cache coherence protocol. This is because write operations trigger the invalidation of copies reside in the cache of other cores, and thus introduces a large amount of memory reads that slow the system down. This project proposes an adaptive cache coherence protocol that can eliminate unnecessary memory accesses by speculatively pushing data modified by producer to potential consumers when producer-consumer sharing patterns are detected. We evaluate the proposed adaptive protocol on ... (results)

I. INTRODUCTION

Shared memory multicore systems play increasingly important roles in both scientific world and the industry. Cache coherence protocol has great influence on performance of shared memory multicore systems. Producer-consumer sharing refers to situations in multi-process synchronization wherein multiple processes share a common, fixed-size buffer and some processes, as known as producers, keep writing new data to that shared buffer while some readers keep reading data from it [1]. The most popular cache coherence protocol used in modern multiprocessor architecture is directory-based write-invalidate protocol, which is inefficient for producer-consumer sharing due to extensive invalidation traffics and expensive remote misses. This project focuses on evaluating the performance of producer-consumer application under write-invalidate and write-update protocols and proposes an adaptive protocol optimized for producer-consumer sharing patterns.

Researches have been done in this area, but they focus more on eliminating unnecessary hops by using additional hardware for cache directory and the proposed pattern detector is not sophisticated enough [2]. Instead of mitigating remote misses, this project focuses on mitigate the slowest part — memory accesses, and proposes a sophisticated detector that investigates producer-consumer sharing patterns at a fine-grained level. More specifically, this project works on single-producer-multiple-consumer patterns and assumes that communications among caches on different cores are physically feasible.

The remainder of this report is organized as follows. We start in Section 2 with an overview of our goals and what we have accomplished. Section 3 and 4 describe the simulator in our study and our analysis of the performance of write-invalidate and write-update protocols on a representative

producer-consumer application. In Section 5, we present our adaptive protocol, followed by a thorough evaluation on several applications in Section 6. Finally, we summarize our works and discuss future directions in Section 7 and 8.

II. GOALS

[75%] Implement cache simulators and testing tools (stack trace, logger, etc.) to fully evaluate and analyze the performance of directory-based write-invalidate and write-update protocols on representative multi-threaded producer-consumer applications.

[100%] Implement proposed adaptive cache coherence protocols and producer-consumer sharing pattern detector. Evaluate and analyze the its performance on representative producer-consumer, as well as general-purpose applications.

[125%] Equip the cache simulator with functionalities that can synchronize application threads so as to observe more accurate cache access events (mitigate the effect of pintool instrumentation).

As we write this report, we have accomplished both 75% and 100% goals. And due to time limits, we decide to work on a more thorough analysis of the proposed protocol on multiple benchmarks instead of the 125% goal.

III. CACHE SIMULATOR

We implement a cache simulator from scratch... [TODO: introduce cache simulator]

IV. WRITE-INVALIDATE AND WRITE-UPDATE PROTOCOLS EVALUATION

In this section, we evaluate the performance of write-invalidate and write-update protocols on a representative application that establishes producer-consumer sharing patterns. In the following evaluations, there are 1 producer and 2 consumers running on different cores.

A. Representative Application

In the application, there is a producer thread that keeps writing to a shared memory while there are multiple consumer threads reading for it. It works well for evaluation purposes because it offers overall performance evaluation as well as fine-grained investigation of the shared memory location.

B. Write-Invalidate Performance

As Figure 1 indicates, the entire programs performance is acceptable because other variables amortize the hit rates. But if we only look at the behavior of shared data the load hit rates are low readers because of a large amount of invalidation.

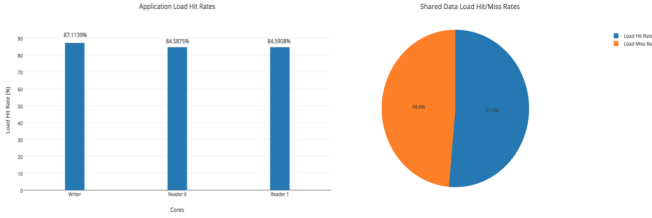


Fig. 1: Write-Invalidate Performance

C. Write-Update Performance

As Figure 2 indicates, the entire programs performance is good and for the shared data alone the load hit rate is perfect. This is because consumers get the data transferred from the produce and thus invalidation and memory reads are no longer needed.

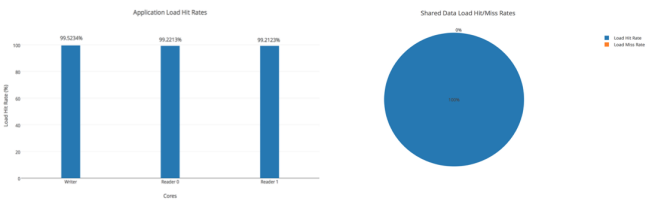


Fig. 2: Write-Invalidate Performance

V. ADAPTIVE CACHE COHERENCE PROTOCOL

In this section, we introduce an adaptive cache coherence protocol that speculatively pushes data to consumers once producer-consumer patterns are detected. Extra bits are needed to track access history for each cache line. They are associated with directory lines which are extended to the structure shown in Figure 3.

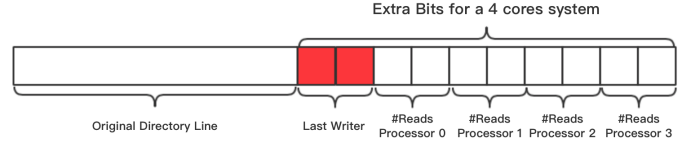


Fig. 3: Extended Directory Line Structure

The extra bits introduced for every cache line is

$$\log N + 2 * N$$

where N is the number of processors, *last_writer* tracks the last one to write this cache line, and there are 2 saturating bits for every processor to track the number of read operations on this cache line. Cache accesses will trigger updates on the extra bits as shown in Algorithm 1 and 2. Note that the extra bits will be discarded to save space when its associated directory line is evicted.

Algorithm 1 On Write Operations

```

1: Let writer_id be the processor performing writes
2: if writer_id == last_writer then
3:   for all pid ≠ writer_id do                                ▷ exclude itself
4:     if CountReads(pid) ≥ 1 then
5:       Cache2CacheDataPush(pid)                                ▷ write-update
6:     else
7:       InvalidateDataCopy(pid)                                ▷ write-invalidate
8:       SaturatingDecrease(pid)                                ▷ decrease by 1
9: else
10:  for all pid ≠ writer_id do                                ▷ exclude itself
11:    InvalidateDataCopy(pid)                                ▷ write-invalidate
12:    ClearCount(pid)                                          ▷ reset read counts
13:  UpdateLastWriter(writer_id)

```

Algorithm 2 On Read Operations

```

1: Let pid refer to a processor
2: for all pid do
3:   SaturatingIncrease(pid)                                ▷ increase by 1

```

VI. ADAPTIVE PROTOCOL EVALUATION

VII. SUMMARY

VIII. FUTURE WORKS

REFERENCES

- [1] https://en.wikipedia.org/wiki/Producer-consumer_problem
- [2] Cheng, Liqun, John B. Carter, and Donglai Dai. "An adaptive cache coherence protocol optimized for producer-consumer sharing." *High Performance Computer Architecture*, 2007. HPCA 2007. IEEE 13th International Symposium on. IEEE, 2007.