



# Understanding Application Contentiousness and Sensitivity on Modern Multicores

**Jason Mars and Lingjia Tang**

Department of Electrical Engineering and Computer Science, University of Michigan, 2260 Hayward Street, Ann Arbor, MI 48109-2121, USA

## Contents

1. Introduction	60
2. Contentiousness vs. Sensitivity	62
2.1 Definition	62
2.2 Contentiousness and Sensitivity	63
2.3 Experiment Design, Results, and Insights	64
2.3.1 Contentiousness	64
2.3.2 Sensitivity	65
2.3.3 Contentiousness vs. Sensitivity	65
3. LLC Misses as an Indicator?	67
4. Predicting Contention Characteristics	70
4.1 Modeling Contention Characteristics	70
4.2 Approximation Using PMUs	72
4.2.1 PMUs for Memory Resource Usage	72
4.2.2 Regression to Determine Coefficients	74
5. Evaluation	78
6. Related Work	81
7. Summary	81
References	82

## Abstract

Runtime systems to mitigate memory resource contention problems on multicore processors have recently attracted much research attention. One critical component of these runtimes is the *indicators* to rank and classify applications based on their contention characteristics. However, although there has been significant research effort, application contention characteristics remain not well understood and indicators have not been thoroughly evaluated.

In this chapter, we performed a thorough study of applications' contention characteristics to develop better indicators to improve contention-aware runtime systems.

The *contention characteristics* are composed of an application's *contentiousness*, and its *sensitivity* to contention. We show that contentiousness and sensitivity are not strongly correlated, and contrary to prior wisdom, a single indicator is not adequate to predict both. Also, while prior wisdom has relied on last level cache miss rate as one of the best indicators to predict an application's contention characteristics, we show that depending on the workloads, it can often be misleading. We then present prediction models that consider contention in various memory resources. Our regression analysis establishes an accurate model to predict application contentiousness. The analysis also demonstrates that performance counters alone may not be sufficient to accurately predict application sensitivity to contention. In this chapter, we also present an evaluation using SPEC CPU2006 benchmarks showing that, when predicting an application's contentiousness, the linear correlation coefficient  $R^2$  of our predictor and the real measured contentiousness is 0.834, as opposed to 0.224 when using last level cache miss rate.



## 1. INTRODUCTION

Multicore processors have become pervasive and can be found in a variety of computing domains, from the most basic desktop computers to the most sophisticated high performance datacenters. With each new generation of architectures, more cores are being added to a single die. In currently available multicore designs, much of the memory subsystem is shared. These shared components include on-chip caches, the memory bus, memory controllers, underlying interconnect, and even on-chip data prefetchers. For such architectures equipped with multiple processing cores, contention for shared resources significantly aggravates the existing memory wall problem and restricts the performance benefit of multicore processors.

There has been a significant research effort to mitigate the effects of contention using software runtime solutions. Ourselves and others have developed techniques that perform runtime contention detection and execution control [21, 29–31, 35] and online job scheduling [2, 6, 12, 14, 17–19, 34, 37]. To most effectively design these runtime systems, there are two important underlying research challenges.

**[Challenge 1]** It is important to have an in-depth understanding of application contention characteristics, including an application's *contentiousness*, which is the potential performance degradation it can cause to its co-runners, and an application's *sensitivity* to contention, which is the potential degradation it can suffer from its co-runners. In prior work, there have been conflicting conclusions about the relationship between an application's contentiousness and its sensitivity to contention. Some prior works [33, 37] argue that there is a clear distinction between an application's contentiousness

and contention sensitivity, while other works [12, 16] conclude that an application's contentiousness and sensitivity are strongly correlated for most applications and thus can be represented and estimated using a unified model. To address this disagreement, we perform thorough investigation of the contentiousness and sensitivity of general purpose applications on current systems.

**[Challenge 2]** Contention-aware runtimes use *indicators* for application contention characteristics to predict the potential performance degradation that may occur due to contention or detect contention as it occurs. Prior works use an application's last level cache (LLC) miss rate as an indicator to detect contention or to predict applications' contention characteristics in order to classify the applications [14, 21, 37]. In fact, LLC miss rate is argued to be one of the most precise indicators for contention-aware scheduling [37]. However, to the best of our knowledge, no prior work has thoroughly investigated how to use microarchitectural events to best construct indicators for an application's contention characteristics. It remains unclear that LLC miss rate is the best performance monitor-based indicator for all workloads. In particular, its accuracy for memory intensive workloads has not been thoroughly evaluated.

This chapter accomplishes five key objectives:

1. We investigate application contention characteristics through systematic experiments on latest multicore hardware and show that although contentiousness and contention sensitivity are consistent characteristics of an application on a given platform, they are not strongly correlated.
2. We explore the effectiveness of using LLC miss rate as an indicator for contentiousness and contention sensitivity and find that it can sometimes be misleading for both. One key insight of our work is that since contentiousness and contention sensitivity are not strongly correlated, no single indicator can accurately predict both.
3. We construct two models that combine usages of multiple memory resources including LLC, memory bandwidth and prefetchers to indicate an application's contention characteristics. Our insights are firstly, understanding contention characteristics require a holistic view of the entire memory subsystem. Secondly, a good indicator for an application's contentiousness must capture the *pressure* an application puts on the shared resources; meanwhile, a good indicator for an application's sensitivity must capture its *reliance* on the shared memory resources. And for many memory resources, pressure, and reliance are very different.

4. We select appropriate performance counters that can capture the usage of various memory resources. We then use regression analysis on a synthetic benchmark suite to establish an accurate model to predict an application's contentiousness. Regression also demonstrates that performance counters alone may not be sufficient to accurately predict an application's sensitivity to contention.
5. We present an evaluation using SPEC CPU2006 benchmarks that shows that when predicting an application's contentiousness, our predictor is much more accurate. The linear correlation coefficient  $R^2$  of our predictor and the real measured contentiousness is 0.834, as opposed to 0.224 when using last level cache miss rate.



## 2. CONTENTIOUSNESS vs. SENSITIVITY

In this section we present formal definitions of both *contentiousness* and *contention sensitivity*, and then investigate key questions about the nature of each and how they relate.

### 2.1 Definition

On multicore processors, an application's *contentiousness* is defined as the potential performance degradation it can cause to co-running application(s) due to its heavy demand on shared resources. On the other hand, an application's *sensitivity* to contention is defined by its potential to suffer performance degradation from the interference caused by its *contentious* co-runners.

As demonstrated in previous work [12], an application  $A$ 's sensitivity is formally defined using the following formula,

$$Sensitivity_A = \frac{IPC_{A(solo)} - \overline{IPC_{A(co-run)}}}{IPC_{A(solo)}}, \quad (1)$$

where  $IPC_{A(solo)}$  is  $A$ 's IPC when it is running alone and  $\overline{IPC_{A(co-run)}}$  is the statistical expectation of the  $A$ 's IPC when it co-runs with random co-runners. We extend this definition to include  $A$ 's contentiousness as,

$$Contentiousness_A = \frac{\overline{IPC_{B_i(solo)}} - \overline{IPC_{B_i(co-run_A)}}}{\overline{IPC_{B_i(solo)}}}, \quad (2)$$

where  $A$ 's contentiousness is quantified as the statistical expectation of the IPC degradation  $A$  causes to its random co-runner.

We can estimate  $Sensitivity_A$  and  $Contentiousness_A$  by co-locating  $A$  with various co-runners  $B_i$ , and take the average of  $A$ 's measured contentiousness and contention sensitivity.  $A$ 's sensitivity to co-runner  $B_i$  can be defined as,

$$Sensitivity_{A(\text{co-run}_{B_i})} = \frac{IPC_{A(\text{solo})} - IPC_{A(\text{co-run}_{B_i})}}{IPC_{A(\text{solo})}} \quad (3)$$

and the  $A$ 's average measured sensitivity is,

$$Sensitivity_{A(\text{avg})} = \frac{\sum_i^n Sensitivity_{A(\text{co-run}_{B_i})}}{n}. \quad (4)$$

Similarly, we can define  $A$ 's contentiousness when it is co-running with  $B_i$  and its average contentiousness as,

$$Contentiousness_{A(\text{co-run}_{B_i})} = \frac{IPC_{B_i(\text{solo})} - IPC_{i(\text{co-run}_A)}}{IPC_{B_i(\text{solo})}}, \quad (5)$$

$$Contentiousness_{A(\text{avg})} = \frac{\sum_i^n Contentiousness_{A(\text{co-run}_{B_i})}}{n}. \quad (6)$$

In this work we use Eq. (4) to estimate  $sensitivity_A$ , and Eq. (6) to estimate  $contentiousness_A$ .

## 2.2 Contentiousness and Sensitivity

In this section we address two important questions about an application's *contentiousness* and *sensitivity* to contention. We first investigated whether contention characteristics (both contentiousness and sensitivity to contention) are *consistent* characteristics of an application. We define consistent as, for a given machine, the relative ordering between all applications' contentiousness and sensitivity in general does not change across different co-runners.

Secondly, we investigated the correlation between an application's contentiousness and its sensitivity to contention. An important observation is that both an application's contentiousness, and its sensitivity to contention, involve the usage of shared resources. One intuition is that contentious applications may also be sensitive to contention and vice versa. Prior work has had conflicting conclusions about the relations between an application's contentiousness and contention sensitivity. There are four possible outcomes. An application can be (1) contentious and sensitive; (2) not contentious and insensitive; (3) contentious but not sensitive; and (4) not contentious but sensitive. Among these four outcomes, Jiang et al. [12, 16] conclude that typical applications' contentiousness and sensitivity are strongly correlated and

should be classified as either contentious and sensitive, or not contentious and insensitive. Xie and Loh [33] on the other hand, argue the existence of applications that are not contentious but sensitive. Meanwhile, other recent works [14, 37] argue that a contentious application that has high cache misses is likely to be very sensitive as well.

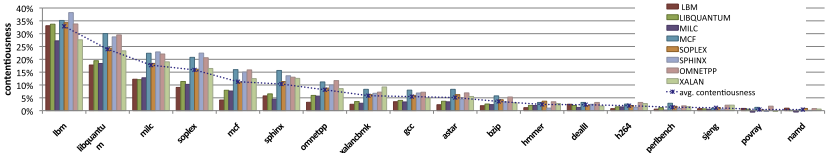
## 2.3 Experiment Design, Results, and Insights

To evaluate these issues, we have performed a series of experiments using 18 benchmarks of SPEC CPU2006 benchmarks suite. These benchmarks represent a diverse range of application workloads and memory behaviors, including different working set sizes, cache misses, and offcore traffic. All experiments were conducted on Intel Core i7 920 (Nehalem) Quad Core with 2.67 GHZ processors, 8 MB last level cache shared by four cores and 4 GB memory. For each experiment, we selected two of the 18 benchmarks, co-located them on neighboring two cores, and measured each benchmark's contentiousness and sensitivity in each experiment using Eqs. (3) and (5). We then calculated each benchmark's average contentiousness and sensitivity using Eqs. (4) and (6). We conducted exhaustive co-running of all possible co-running pairs, which is a total of  $162(\frac{18 \times 18}{2})$  co-running experiments executed to completion on `ref` inputs. Each experiment was conducted three times to calculate the average. Note that SPEC runs are fairly stable and there is little variance between runs.

### 2.3.1 Contentiousness

Figure 1 presents our benchmarks' *contentiousness*. This contentiousness is calculated using Eq. (5), which indicates the performance degradation each of the 18 benchmarks causes to its co-runner. The 18 benchmarks are shown on the  $x$ -axis. For each of the 18 benchmarks, we show its measured contentiousness when it is co-running with each of the 8 most contentious co-runners respectively. Each bar represents a co-runner. Only 8 co-runners are shown in the figure because of the space limit. The dotted line shows the average contentiousness of each benchmark, computed by averaging each benchmark's 18 contentiousness values across 18 co-runners using Eq. (6). The 18 benchmarks on the  $x$ -axis are then sorted by their average contentiousness. The line graph for average contentiousness shows a general descending trend.

Figure 1 demonstrates that contentiousness is a consistent characteristic of an application. The relative order of benchmarks' contentiousness stays fairly consistent regardless of which co-runner is present. For example, when



**Fig. 1.** Contentiousness. Each bar shows the performance degradation of a co-runner caused by the application across  $x$ -axis.

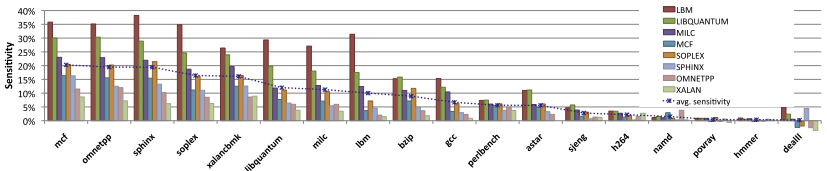
comparing each benchmark’s contentiousness when it is co-running with `lbm`, shown by the first bar for each 18 benchmark, we notice that the contentiousness of 18 benchmarks are almost all in descending order along the  $y$ -axis mirroring the dotted line. This also applies to all other co-runners as well. The graph also shows that `lbm` is the most contentious benchmark among the 18 benchmarks.

**2.3.2 Sensitivity**

Similar to Fig. 1, Fig. 2 shows the sensitivity to contention of each of the 18 benchmarks when co-located with the most contentious applications. This sensitivity is calculated using Eq. (3), indicating how much degradation the 8 co-runners cause to each of the 18 benchmarks. These 18 benchmarks are sorted according to their average sensitivity, calculated using Eq. (4). Similar to Fig. 1, this figure shows that sensitivity is also consistent for each application. Although the descending trend is not as consistent as Fig. 1, the general trend is strong.

**2.3.3 Contentiousness vs. Sensitivity**

In Fig. 3, we juxtapose contentiousness and sensitivity. In this graph, for each application across the  $x$ -axis, the first bar shows the average contentiousness



**Fig. 2.** Sensitivity. Each bar shows the performance degradation of the application across  $x$ -axis caused by each of the 8 different co-runners.

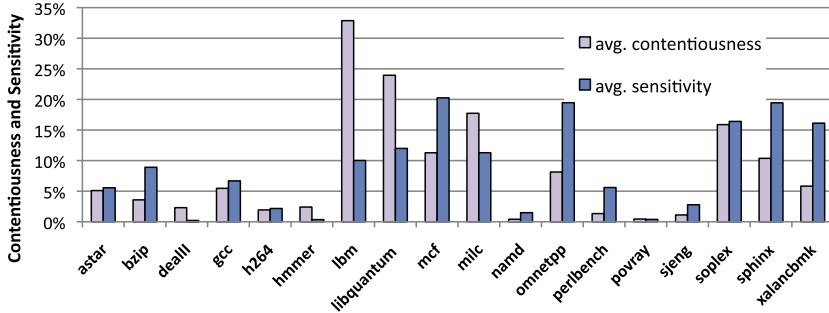


Fig. 3. Average contentiousness vs. sensitivity.

of this application with the 18 co-runners presented in Figs. 1 and 2. The second bar shows each benchmark’s average sensitivity to the same set of co-runners. Figure 3 clearly demonstrates a large disparity between application contentiousness and sensitivity. As shown in the figure, applications such as *lbm* and *libquantum* are highly contentious and only mildly sensitive, while other applications such as *omnetpp* and *xalan* are highly sensitive, and slightly contentious. Also notice that, in Figs. 1 and 2, the sorted ordering of the 18 benchmarks ( $x$ -axis) are almost completely different. In fact, the correlation coefficient between contentiousness and sensitivity using linear regression is 0.48, which further shows they are not strongly correlated.

To summarize, through our experimentation we find,

1. Contentiousness and sensitivity are an application’s consistent characteristics. Figure 1 shows that applications with higher contentiousness tend to be consistently more contentious regardless of co-runners. This general trend also applies to sensitivity, as shown in Fig. 2.
2. Contentiousness and sensitivity of general purpose applications are not strongly correlated as shown in Fig. 3. While we do not observe applications that are only sensitive or only contentious, four outcomes occur in practice; applications can be (1) contentious and sensitive; (2) not contentious and insensitive; (3) contentious but not highly sensitive; (4) not highly contentious but sensitive.

We present analysis as why contentiousness and sensitivity are different in Section 4.1. Section 4.2.2 also presents more experimental data on a different set of benchmarks to demonstrate the difference between contentiousness and sensitivity.



### 3. LLC MISSES AS AN INDICATOR?

The ability to predict application contention characteristics is important for contention-aware runtime systems. In this chapter, we focus on indicators using performance monitoring units (PMUs). Last level cache (LLC) miss rate is one of the most commonly used indicators of an application's contentiousness and is used to classify applications to achieve sensible co-scheduling decisions [14, 37] and detect contention online [21]. In this section we evaluate the effectiveness of using last level cache misses to indicate an application's level of contentiousness and sensitivity to contention.

Both LLC *miss rate*, the number of misses for a given amount of time, and *miss ratio*, the number of misses for a given number of instructions, have been used by prior work to perform contention-aware scheduling. To evaluate whether LLC miss rate or ratio is a good indicator for an application's contentiousness, we measure LLC miss rate and miss ratio for the 18 SPEC2006 benchmarks used in Section 2, and compare each benchmark's rate and ratio against the average degradation it causes to its co-runners. We also compare each benchmark's miss rate and ratio to the average degradation it suffers due to contention to evaluate if LLC miss is a good indicator for sensitivity. Experiment set up is as described in Section 2.3. Both the LLC miss rate and ratio are collected when each benchmark is running alone using pfmon [7].

Figures 4 and 5 compare the average contentiousness and sensitivity of the benchmarks with their LLC misses per million instructions. Figures 6 and

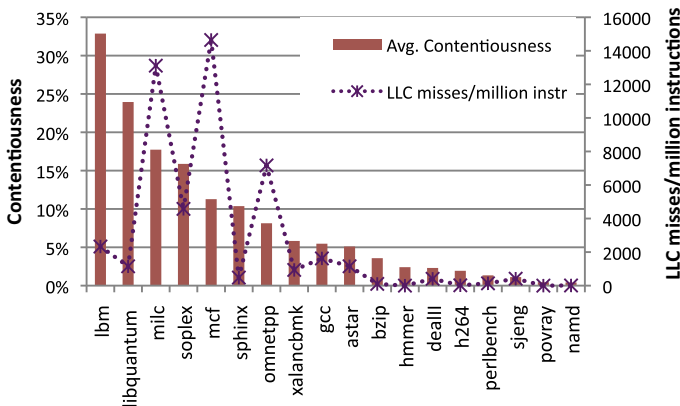


Fig. 4. LLC miss ratio vs. average contentiousness.

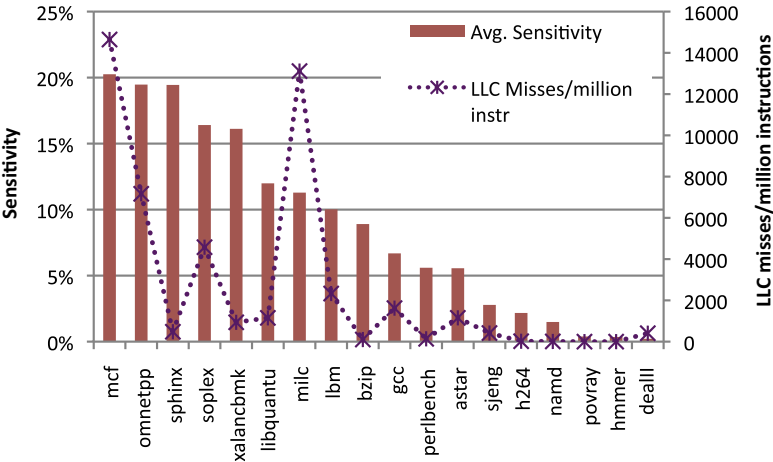


Fig. 5. LLC miss ratio vs. average sensitivity.

7 compare the average contentiousness and sensitivity with LLC misses per millisecond. In these four figures, each bar shows the average contentiousness or sensitivity of each application as measured from our experimentation in Section 2. The dotted line shows the each benchmark’s LLC miss rate or ratio. We use line graphs to better demonstrate the difference between the trend of LLC misses and each application’s contentiousness or sensitivity. The left  $y$ -axis shows the contentiousness and sensitivity, respectively, and the right  $y$ -axis shows the LLC misses rate and ratio.

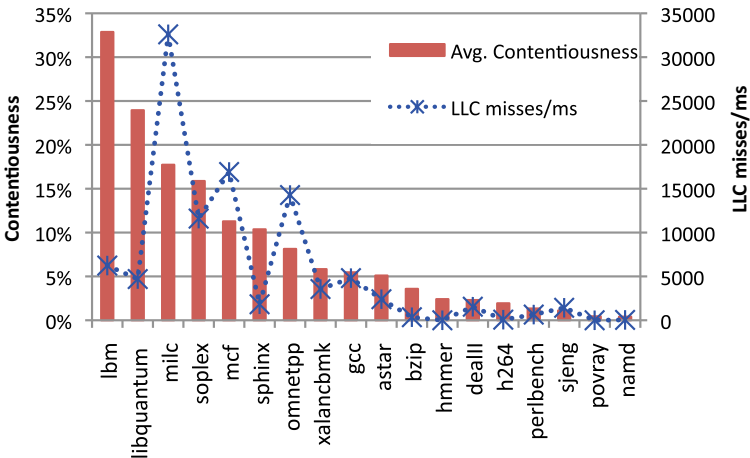


Fig. 6. LLC miss rate vs. average contentiousness.

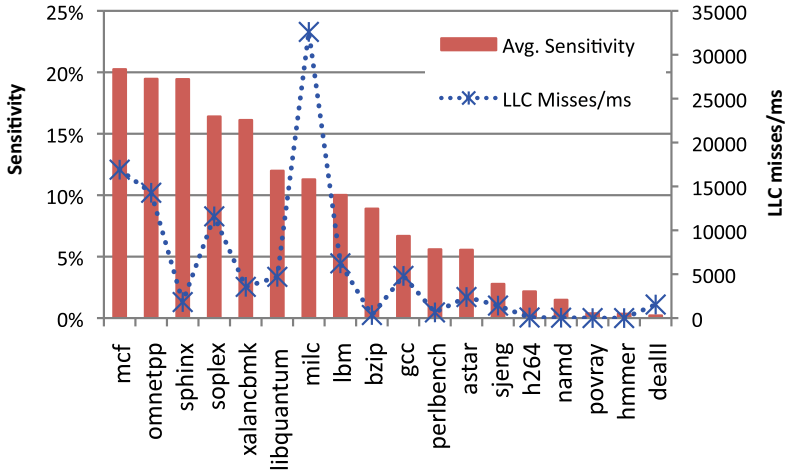


Fig. 7. LLC miss rate vs. average sensitivity.

These figures demonstrate two key observations. The first observation is that **LLC miss rate and ratio are good indicators to distinguish CPU-bound applications and memory-bound applications**. Applications that are shown to the right of each figure, such as `hmm`, `sjeng`, and `povray` are CPU bound applications. They tend to have little contentiousness or sensitivity to contention, and this is accurately predicted by their extremely low cache miss rate/ratio. This insight indicates that a contention-aware runtime system that uses LLC misses to predict performance degradation or detect contention may be quite effective for workloads that contain a balanced mix of CPU bound applications and memory bound applications, as co-scheduling CPU bound and memory bound applications does indeed minimize contention effectively. The scheduler would simply have to pair low LLC miss benchmarks (e.g., `povray`) with high LLC misses benchmarks (e.g., `milc`). Also contention may not occur when CPU bound applications are executing and thus using LLC miss rate can fairly effectively detect when contention is not occurring. This observation may explain the good results in prior work.

The second key observation is that **LLC cache miss rate and ratio are not good at predicting the degree of contentiousness and sensitivity for memory bound applications**. This is demonstrated by the mismatch between the dotted line and bars for benchmarks to the left of each figure. These benchmarks exhibit various levels of contentiousness and sensitivity ranging from 5% to 35% for contentiousness and 20% for

sensitivity). However using LLC miss rate and ratio gives little indication of the magnitude of the contentiousness or sensitivity. For example, in Figs. 4 and 6, *lbm* and *libquantum* are the two most contentious benchmarks, yet their LLC misses rate and ratio are quite low. In Figs. 5 and 7, *sphinx* is shown to be one of the most sensitive benchmarks, yet its LLC miss rate and ratio are almost negligible. From these observations, we conclude that LLC miss rate or ratio is not a good indicator for predicting the magnitude of contentiousness and sensitivity of a memory bound application, and therefore is not suitable for scheduling workloads that are memory bound biased (containing more memory bound applications than cpu bound) or detecting the severity of contention among such workloads.

Section 4.2.1 presents more details on why LLC miss rate is not a good indicator for contentiousness or sensitivity.



## 4. PREDICTING CONTENTION CHARACTERISTICS

In this section we construct models to indicate contention characteristics for all types of workloads including memory bound workloads. One of the key insights of this work is that *because an application's contentiousness and its sensitivity to contention are two distinct characteristics, we need separate predictors for each*. Also, based on the results presented in Section 3, we conclude that contention also occurs in other shared components in the memory subsystem in addition to last level caches. Therefore, understanding the contention characteristics of an application requires a holistic view of the memory subsystem and a comprehensive predictor must capture how an application uses and relies on the shared resources beyond last level cache such as memory bandwidth, prefetchers, memory controllers, etc.

In this section, we first construct general models to estimate an application's contention characteristics that take sharing of multiple memory components into consideration. We then select PMUs that can reflect application's activity in regard to these shared memory components. Finally, we determine the detailed prediction models using regression analysis between an application's selected PMUs profile and its contention characteristics.

### 4.1 Modeling Contention Characteristics

#### Why are applications' contentiousness and sensitivity are different?

The fundamental difference between contentiousness and sensitivity is that *contentiousness reflects how much pressure an application puts on the shared resources; meanwhile, application sensitivity to contention reflects an application's*

**reliance** on the shared memory resources. Last level caches and prefetchers are both essentially performance optimization mechanisms whose effectiveness is depending on application's data reuse patterns, therefore for these two resources, there is a difference between an application's pressure and reliance on them. **Pressure** is directly linked to how much the shared resource (LLC or prefetcher) an application is using; while **reliance** is how much an application's progress is benefiting from using the shared resource.

For example, an application's working set may occupy a great amount of LLC space but the application's performance may or may not rely on the LLC, depending on whether it reuses its data residing in the LLC. Another example is that an application can issue a large amount of prefetching requests but may not benefit or only benefit slightly from these requests. In this case, the application is heavily using but not depending on the prefetcher. For other components such as memory bandwidth, pressure, and reliance can be more correlated.

To model an application's contention characteristics, we use a linear model to combine the effect of shared resources, including LLC, memory bandwidth and prefetchers. We also consider contentiousness and sensitivity to contention separately.

**Contentiousness.** An application's contentiousness is determined by the amount of pressure it puts on the shared memory subsystem. Thus it can be directly predicted using the application's usage of the shared resources.

$$C = a_1 \times \text{LLC\_usage} + b_1 \times \text{BW\_usage} + c_1 \times \text{Pref\_usage}, \quad (7)$$

where  $C$  stands for contentiousness, BW is bandwidth, and Pref is prefetchers. It is fairly easy to quantify and measure the bandwidth usage (e.g., bus transactions per second). However, it is difficult to quantify cache usage because it is multifaceted. For example, both the cache access frequency and the cache footprint reflect dimensions of the cache usage.

Each application may have a different combination of cache, bandwidth, and prefetch usage. For example, a cache-intensive application whose working set is similar to the size of the LLC has a heavy LLC usage and probably little bandwidth usage. Streaming applications may have little to medium cache usage but heavy bandwidth usage. How contentious these applications are relative to each other depends on the relative importance between the cache contention and the bandwidth contention. Note that the goal of the prediction model is to rank the relative contentiousness of a group of applications to make scheduling decisions, instead of predicting the *exact* average contentiousness or the *exact* performance degradation. Therefore

identifying the relative importance of contention in shared caches, bandwidth, and prefetchers, reflected as coefficients  $a_1$ ,  $b_1$ , and  $c_1$ , is one of the main objectives of the modeling and regression. The next section will present the regression analysis for determining coefficients of the model. It is worth noting that  $a_1$ ,  $b_1$ ,  $c_1$  are architecture-specific.

**Sensitivity.** A good prediction model for sensitivity should capture how much the application is relying on the shared memory system. However, this is much more challenging than predicting contentiousness using PMUs.

$$S = a_2 \times \text{LLC\_usage} + b_2 \times \text{BW\_usage} + c_2 \times \text{Pref\_usage}, \quad (8)$$

As shown in Eq. (8), to capture the difference between contentiousness and sensitivity, we use difference coefficients (e.g.,  $a_1$  vs.  $a_2$ ). In addition to being architecture-specific, coefficients  $a_2$ ,  $b_2$ , and  $c_2$  are also application specific. This is because, as we discussed earlier, even with the same amount of resource usage, how much an application relies on the shared resources is different. And it is heavily depending on how applications reuse data.

## 4.2 Approximation Using PMUs

In this section, we identify performance counters (PMUs) to estimate the usage of memory resources including LLC and memory bandwidth. We then profile a set of synthetic benchmarks to collect the selected performance counters as well as the contention characteristics of these benchmarks on a real architecture. Using performance counter profiles to estimate resource usages in Eqs. (7) and (8), we can use regression analysis to determine coefficients of the models. The platform we use in this section is a quad-core Intel Core i7 described in Section 2.3.

### 4.2.1 PMUs for Memory Resource Usage

**Contentiousness.** On our Intel Core i7 platform, we identify the number of cache lines the last level cache (LLC) brings in per millisecond (L3LinesIn/ms), as shown in Fig. 8, to measure the memory bandwidth usage. This is because that LLC lines in rate can better capture the actual aggregate pressure an application is putting on the bandwidth than LLC miss rate or ratio because it includes prefetchers' effect on the bandwidth. We identify (L2LinesIn–L3LinesIn)/ms to estimate the shared cache (L3) usage. (L2LinesIn–L3LinesIn) rate shows how much data is used in an interval that is coming from only L3 and not the DRAM. However, unlike using L3LinesIn/ms to estimate the bandwidth usage, (L2LinesIn–L3LinesIn) rate is an approximation of the L3 cache usage. As we discussed, both the cache

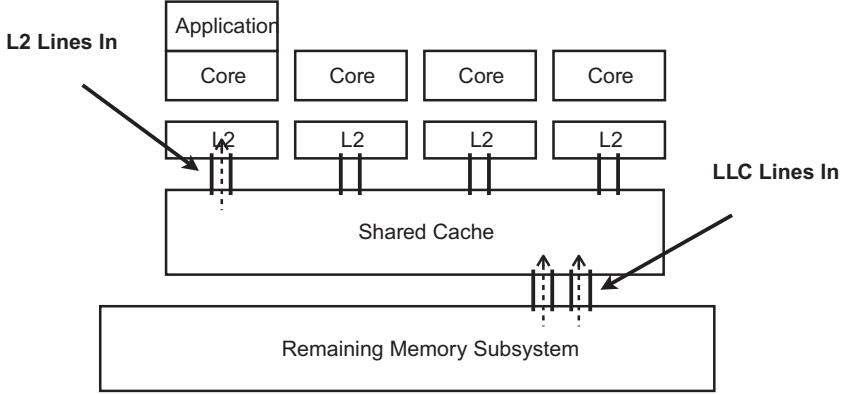


Fig. 8. PMUs used for predicting contention characteristics.

footprint and the access frequency are dimensions of the cache usage. Bigger footprint and higher access frequency indicate more pressure on the cache. (L2LinesIn–L3LinesIn) rate only reflects the frequency but may not fully reflect the application’s footprint in the L3 cache because PMUs do not reflect the amount of data reuse. However, we will show that this is a sufficient approximation when indicating contentiousness. Prefetcher usage is manifested in both cache and bandwidth usage. The main impact of prefetchers is the increased bandwidth and the cache space the prefetched data occupy. Because both L3LinesIn and L2LinesIn include the prefetchers’ traffic, we do not need an extra PMU to measure the prefetcher usage. Although we use an Intel Core i7 platform here, the reasoning of selecting PMUs should be general for other multicore architectures. Using the above PMUs, Eq. (7) becomes:

$$C = a_1 \times (\text{L2LinesIn\_rate} - \text{L3LinesIn\_rate}) + b_1 \times \text{L3LinesIn\_rate} \quad (9)$$

### Why LLC miss rate is not a good indicator for contentiousness?

Our experiments in Section 3 show that solo LLC miss rate and ratio do not accurately indicate an application’s level of contentiousness. There are two main reasons that our model (Eq. (9)) can be more accurate. Firstly, LLC miss rate does not fully reflect the contention for the memory bandwidth or prefetcher. LLC miss rate or ratio, as an architectural performance monitoring event on most platforms, does not capture the prefetching bandwidth, which often consumes a large portion of the memory bandwidth on modern architectures. Secondly, LLC miss rate and ratio also cannot accurately capture cache contentiousness of an application. An application can have a

working set that fits in the L3 cache. The application can frequently access its working set without incurring many cache misses. However, since it is heavily using the shared cache, it can be very cache contentious when co-located and causing cache misses to its co-runners. LLC miss rate cannot accurately predict cache-intensive applications' contentiousness but (L2LinesIn - L3LinesIn) rate can.

**Sensitivity.** We use the same PMUs to estimate the reliance an application has on the shared resources and to predict the application's sensitivity to contention.

$$S = a_2 \times (\text{L2LinesIn\_rate} - \text{L3LinesIn\_rate}) + b_2 \times \text{L3LinesIn\_rate}. \quad (10)$$

As we mentioned in Section 4.1,  $a_2$  and  $b_2$  are application-specific coefficients that are related to how an application reuses its data. However, due to the limitation of current available PMUs on most hardware, we cannot accurately measure data reuse. Therefore, the goal of the regression analysis for the sensitivity model is to investigate whether these application specific factors are not negligible when predicting an application's sensitivity.

**Why LLC miss rate is not a good indicator for sensitivity?** As we discussed in the previous section, LLC miss rate does not always reflect the reliance an application has on LLC or the rest of the shared memory system. Firstly, an application can be highly relying on the LLC, occupying large portion of the LLC and frequently accessing it without incurring LLC misses. This type of applications actually may be highly sensitive. However, the LLC miss rate does not reflect that. Secondly, multiple shared memory components also need to be considered for sensitivity to contention. For example, sensitivity to bandwidth contention is not considered previously. Streaming applications are considered to be contentious but not necessarily sensitive because they are already having cache misses when it is running alone. So it would seem that co-running with other applications would not make the situation worse. However, they are highly sensitive to memory bandwidth contention although not to cache contention. These applications also may not have high LLC miss rates.

#### 4.2.2 Regression to Determine Coefficients

In this section, we use multiple regression to determine the coefficients in Eqs. (9) and (10). The goal of the regression analysis is to firstly test that whether there is a strong correlation between an application's resource usage and its contention characteristics; and secondly to determine the relative importance of contention in various resources.



**Table 1** Synthetic benchmarks.

Benchmark	Footprint	Description
Bst	4 mb, 8 mb, 50 mb	Random accessing a binary search tree
Naive	4 mb, 8 mb, 50 mb	Random accessing an array
Er-naive	4 mb, 8 mb, 50 mb	Fast random accessing an array
Blockie	Small, medium, large	A number of large 3D arrays. A portion of one array is continuously copied to another
Sledge	Small, medium, large	Two large arrays, copies data back and forth between arrays with this sledgehammer pattern

**Synthetic Benchmarks.** To conduct regression analysis, we collect PMU profiles and contention characteristics of a suite of synthetic benchmarks. Table 1 presents our synthetic benchmarks. `Bst`, `naive`, `blockie`, and `sledge` are from the contention benchmark suite developed by Mars and Soffa [16]. The benchmarks are memory intensive applications with various memory access patterns. They are run using three different inputs with different working set sizes to stress different memory resources. The only difference between `naive` and `er-naive` is that `er-naive` uses a much faster random number generator. The goal is to test how contention characteristics would change when an application’s cache access frequency increases but everything else remains the same. Figure 9 presents each benchmark’s average contention characteristics calculated using Eqs. (4) and (6). As the figure shows, the benchmark suite presents a fairly wide range of contentiousness and sensitivity. Also this figure again demonstrates that an application’s contentiousness and sensitivity are not strongly correlated.

**Regression.** We conduct multiple linear regression on Eq. (9) using each benchmark’s `L2LinesIn` rate, `L3LinesIn` rate, and average `C` (contentiousness), shown in Fig. 9. The regression result for contentiousness is:

$$C = 1.663 \times (\text{L2LinesIn}/ns - \text{L3LinesIn}/ns) + 8.890 \times \text{L3LinesIn}/ns + 0.044 \quad (11)$$

The  $p$  value for  $(\text{L2LinesIn}/ns - \text{L3LinesIn}/ns)$  is 0.018; for  $\text{L3LinesIn}/ns$ ,  $5.11\text{e-}07$ ; for the entire regression,  $2.015\text{e-}06$ ; all smaller than 0.5, indicating statistically significant effects. The  $R$ -squared is 0.8876, indicating a strong fit. The coefficients show the relative importance between the bandwidth usage and the LLC usage, indicating that memory bandwidth contention has a more dominating effect. Figure 10 presents benchmarks’

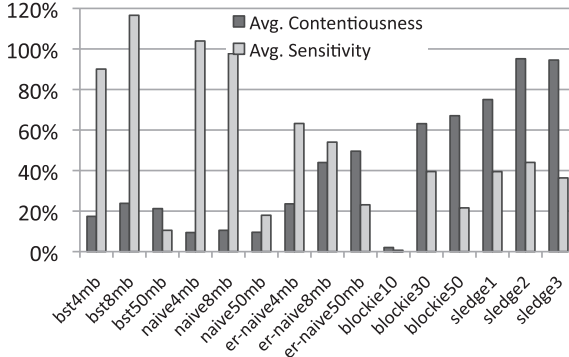


Fig. 9. Average contentiousness and sensitivity of synthetic benchmarks.

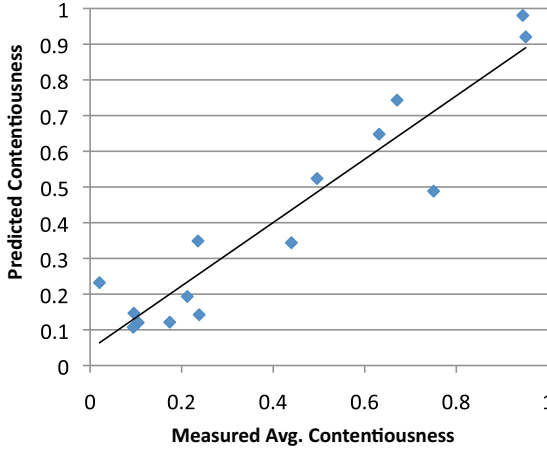


Fig. 10. Regression Result for Contentiousness using L2LinesIn and L3LinesIn. The figure shows a strong linear fit.

predicted contentiousness values using the regression model (Eq. (11)) comparing against the measured actual average contentiousness.

Figures 11 and 12 demonstrate the relative importance of contention in the memory bandwidth and contention in the LLC. Figure 11 shows that for most benchmarks, L3LinesIn rate can be very indicative of an application’s contentiousness. Applications with high L3LinesIn\_rate are in general causing more performance degradation to its co-runners. This is true except for a few benchmarks including er-naive4mb, er-naive8mb, bst4mb, and bst8mb. Those benchmarks have minimum L3LinesIn\_rate but they have medium levels of contentiousness. This is because they are

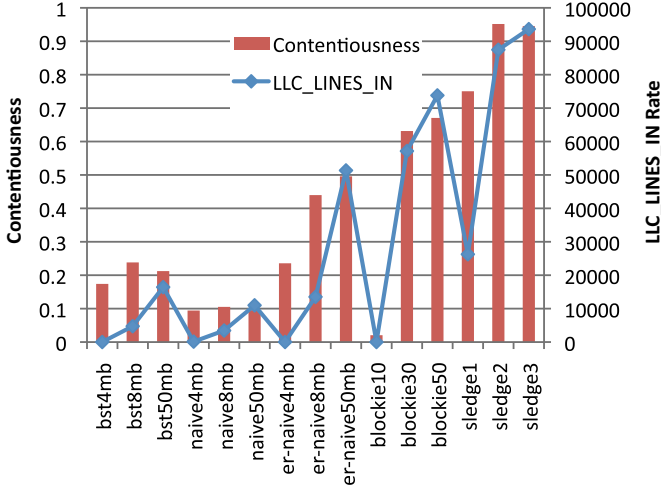


Fig. 11. Benchmarks' average contentiousness vs. their L3LinesIn/ms.

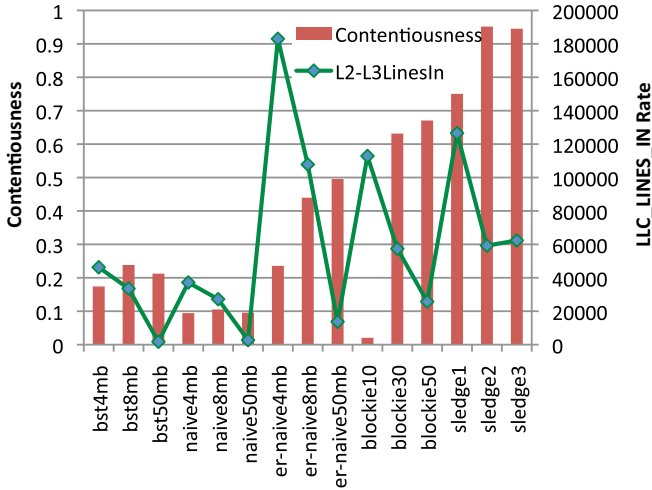


Fig. 12. Benchmarks' average contentiousness vs. their (L2LinesIn-L3LinesIn)/ms.

contentious for the shared cache instead of contentious for the bandwidth. Figure 12 shows that these benchmarks have a medium to high (L2LinesIn-L3LinesIn) rate, indicating that (L2LinesIn-L3LinesIn) rate can capture their potential cache contentiousness. It is not as accurate as to predicting bandwidth contention because as we mentioned, cache usage is more difficult to capture using PMUs. Note that their contentiousness level is mild

comparing to benchmarks such as `blockie` and `sledge`. This is consistent with the regression results that bandwidth usage is more important (has a much higher coefficient) than cache usage.

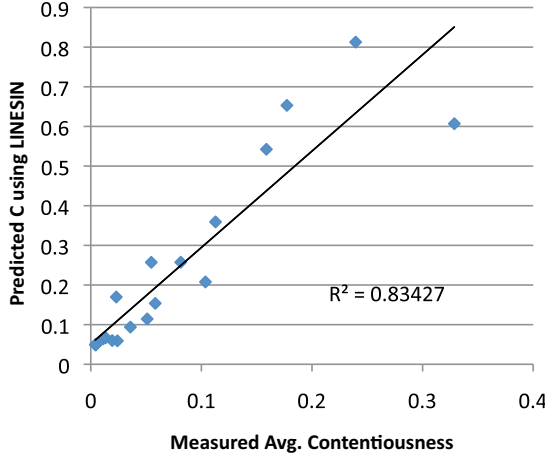
However, regression for Eq. (10) cannot establish a linear model for sensitivity, indicating that application-specific factors such as locality play a non-negligible role in deciding applications' sensitivity and PMU alone may not be a good candidate for an accurate prediction model. It is worth noting that predicting sensitivity is challenging using other approaches too. Reuse distance profile can capture the application's locality characteristics. However, most works [12, 37] using reuse distance profile only consider contention in the last level cache, and it may be difficult to simulate and combine the contention effect in various other resources such as sophisticated prefetchers. One promising approach is through direct empirical measurement instead of indirect PMU indicators. `Cipe`, proposed by Mars et al. [20], empirically measures application's sensitivity in a controlled synthesised environment.

**Summary.** In summary, an application's contentiousness is determined by the pressure the application places on the shared memory subsystem. On the Intel Core i7, a combination of `L2LinesIn` and `L3LinesIn` rate is a better indicator of contention characteristics instead of LLC misses. One key insight is, because the fundamental difference between an application's contentiousness and sensitivity to contention (e.g., contentiousness is directly related to resource usage but sensitivity is related to the dependence on the resource), it is easier to predict an application's contentiousness using PMUs. However, PMU alone may not be sufficient for an accurate sensitivity prediction. In addition, because of the complexity of the memory system design on modern multicore architectures, a good predictor for contentiousness needs to fully reflect the aggregate usage of a number of resources including shared caches, memory bandwidth, prefetchers, etc.



## 5. EVALUATION

In this section, we evaluate our prediction model for application's contentiousness (Eq. (11)) using SPEC CPU2006 benchmarks. All experiments are conducted on quad-core Intel Core i7 described in Section 2.3. Each benchmark's contentiousness is measured as described in Section 2.3.1,



**Fig. 13.** Predicted contentiousness using our model is highly correlated with the real measured contentiousness for SPEC benchmarks.

shown in Fig. 3. We also measure each benchmark’s solo L2LinesIn\_rate and L3LinesIn\_rate. Using the PMU profiles, we calculate the predicted contentiousness using Eq. (11).

Figure 13 presents our prediction results compared to the real measured contentiousness for SPEC CPU2006 benchmarks. The linear correlation coefficient  $R$  is 0.91, indicating our prediction is highly correlated with the real measured contentiousness. Note that we are not predicting the actual value of contentiousness because most contention-aware runtime systems only need to rank applications according to their contentiousness levels. Based on the ranking, the scheduler then co-locates highly contentious applications with applications that are not so contentious. The strong correlation (0.91) demonstrates that our prediction model can successfully rank the contentiousness levels of applications and thus can greatly improve scheduling decisions. For comparison, Fig. 14 shows the results when using LLC miss rate to predict applications’ contentiousness. Figure 15 shows the prediction results using LLC reference rate. Zhuravlev et al. [37] propose using LLC reference rate to predict an application’s intensity (contentiousness). The correlation coefficients ( $R$ ) are 0.47 and 0.28, respectively, showing that neither LLC miss rate nor LLC reference rate can accurately indicate application contentiousness.

Our evaluation shows that our prediction model can indicate applications’ contentiousness much more accurately than the state-of-the-art LLC

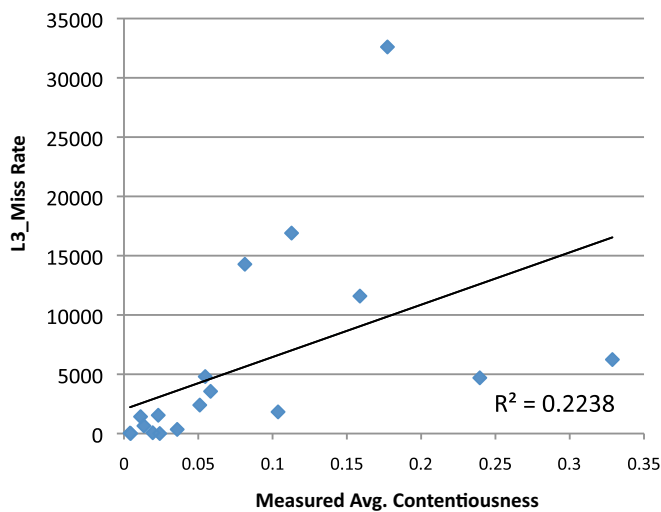


Fig. 14. L3 Miss rate is not strongly correlated with the real contentiousness.

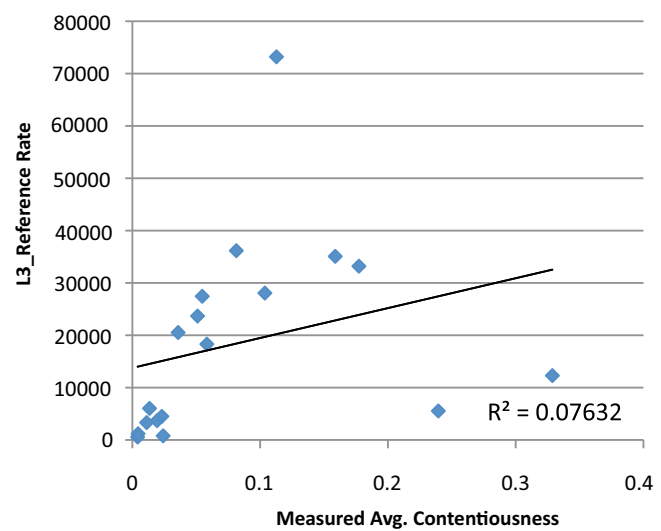


Fig. 15. L3 reference rate is not strongly correlated with the real contentiousness.

miss rate indicator. And our contentiousness model can improve contention-aware runtime solutions that base on PMUs to indicate applications' contentiousness levels.



## 6. RELATED WORK

There has been a wealth of research on the challenge of shared resource contention on multicore processors. Contention-aware runtime systems have been proposed to mitigate the effect of contention [2, 12, 14, 21, 22, 34, 37]. Jiang et al. develop a locality model to predict co-running applications' degradation and use the model for co-scheduling to reduce performance degradation and unfairness [12]. Zhuravlev et al. demonstrate that cache contention is not the dominant cause for performance degradation of co-running applications on CMPs; contention that happens in many components of the memory subsystem all contributes to the performance degradation. They also conclude that last level cache miss ratio is one of the best predictor for co-running applications' performance degradation [37]. Jiang et al. and Tian et al. study the theoretical complexity of co-scheduling and provide approximate algorithms [11, 32]. Also, there has been a number of contention-aware scheduling schemes proposed that guarantee fairness and *Quality-of-Service* for multiprogrammed and multithreaded applications [1, 8, 14]. Fedorova et al. use cache model prediction to enhance the OS scheduler to provide performance isolation by allocating CPU resources according to contention interference [8]. Hardware techniques and related algorithms to enable cache management such as cache partitioning and memory scheduler have been proposed [4, 13, 23, 25, 28]. Iyer et al. proposed a QoS-enabled memory architecture for CMP platforms to allocate memory resources such as cache and memory bandwidth [10]. Other hardware solutions have been developed to guarantee fairness and QoS [9, 15, 24, 26]. Related to novel cache designs and architectural support, analytical models to predict the impact of cache sharing are also proposed by Chandra et al. [3]. In addition to new hardware cache management, other approaches manage the shared cache through the OS [5, 8, 27, 36].



## 7. SUMMARY

In this chapter, we performed a thorough study of contention characteristics to develop an improved predictor for contention-aware runtime systems. We studied the two aspects of an application's contention characteristics: an application's *contentiousness*, e.g., the amount of degradation it tends to cause to its co-runners due to its demand on shared resources, and an application's *sensitivity*, e.g., the amount of degradation the application is likely to

suffer due to co-running with contentious applications. Our study found that although these two characteristics are consistent for each application, they are not strongly correlated for general purpose applications. We also found that although last level cache miss rate is a commonly perceived good indicator for application contention characteristic, it could often be misleading. Based on the findings and insights, we then present prediction models that comprehensively consider contention in various memory resources. Our regression analysis establishes an accurate model to predict application contentiousness. Further evaluation using SPEC CPU2006 benchmarks shows that our predictor significantly outperforms the state-of-the-art PMU indicators.

## REFERENCES

- [1] M. Banikazemi, D. Poff, B. Abali, PAM: a novel performance/power aware meta-scheduler for multi-core systems, in: SC '08: Proceedings of the 2008 ACM/IEEE Conference on Supercomputing, 2008.
- [2] M. Bhaduria, S. McKee, An approach to resource-aware co-scheduling for cmps, in: ICS '10: Proceedings of the 24th ACM International Conference on Supercomputing, June 2010.
- [3] D. Chandra, F. Guo, S. Kim, Y. Solihin, Predicting inter-thread cache contention on a chip multi-processor architecture, in: HPCA '05: Proceedings of the 11th International Symposium on High-Performance Computer Architecture, 2005.
- [4] J. Chang, G.S. Sohi, Cooperative cache partitioning for chip multiprocessors, in: ICS '07: Proceedings of the 21st Annual International Conference on Supercomputing, 2007.
- [5] S. Cho, L. Jin, Managing distributed shared L2 caches through OS-level page allocation, in: MICRO 39: Proceedings of the 39th Annual IEEE/ACM International Symposium on Microarchitecture, 2006.
- [6] C. Delimitrou, C. Kozyrakis, Paragon: QoS-aware scheduling for heterogeneous data-centers, in: Proceedings of the 18th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), March 2013.
- [7] S. Eranian, What can performance counters do for memory subsystem analysis? in: Proceedings of the 2008 ACM SIGPLAN Workshop on Memory Systems Performance and Correctness: Held in Conjunction with the Thirteenth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'08), 2008, pp. 26–30.
- [8] A. Fedorova, M. Seltzer, M.D. Smith, Improving performance isolation on chip multiprocessors via an operating system scheduler, in: PACT '07: Proceedings of the 16th International Conference on Parallel Architecture and Compilation Techniques, 2007.
- [9] A. Herdrich, R. Illikkal, R. Iyer, D. Newell, V. Chadha, J. Moses, Rate-based QoS techniques for cache/memory in CMP platforms, in: ICS '09: Proceedings of the 23rd International Conference on Supercomputing, 2009.
- [10] R. Iyer, L. Zhao, F. Guo, R. Illikkal, S. Makineni, D. Newell, Y. Solihin, L. Hsu, S. Reinhardt, QoS policies and architecture for cache/memory in CMP platforms, in: ACM SIGMETRICS, Performance Evaluation Review, vol. 35, 2007.
- [11] Y. Jiang, X. Shen, J. Chen, R. Tripathi, Analysis and approximation of optimal co-scheduling on chip multiprocessors, in: PACT '08: Proceedings of the 17th International Conference on Parallel Architectures and Compilation Techniques, 2008.



- [12] Y. Jiang, K. Tian, X. Shen, Combining locality analysis with online proactive job co-scheduling in chip multiprocessors, in: *High Performance Embedded Architectures and Compilers*, 2010.
- [13] S. Kim, D. Chandra, Y. Solihin, Fair cache sharing and partitioning in a chip multiprocessor architecture, in: *PACT '04: Proceedings of the 13th International Conference on Parallel Architectures and Compilation Techniques*, 2004.
- [14] R. Knauerhase, P. Brett, B. Hohlt, T. Li, S. Hahn, [Using OS observations to improve performance in multicore systems](#), *IEEE Micro* 28 (3) (2008) .
- [15] J. Lin, Q. Lu, X. Ding, Z. Zhang, X. Zhang, P. Sadayappan, Gaining insights into multicore cache partitioning: bridging the gap between simulation and real systems, in: *The IEEE 14th International Symposium on High Performance Computer Architecture*, 2008, pp. 367–378.
- [16] J. Mars, M.L. Soffa, Synthesizing contention, in: *Workshop on Binary Instrumentation and Applications*, 2009.
- [17] J. Mars, L. Tang, R. Hundt, in: *ISCA '13: Proceedings of the 40th Annual International Symposium on Computer Architecture*, IEEE/ACM, 2013.
- [18] Jason Mars, Lingjia Tang, Kevin Skadron, Mary Lou Soffa, Robert Hundt, Increasing utilization in modern warehouse-scale computers using bubble-up, *IEEE Micro* 32 (3) (May 2012) 88–99, < <http://dx.doi.org/10.1109/MM.2012.22>> (issn: 0272-1732).
- [19] J. Mars, L. Tang, R. Hundt, K. Skadron, M. Soffa, [Bubble-up: increasing utilization in modern warehouse scale computers via sensible co-locations](#), in: *MICRO '11: Proceedings of the 44th Annual IEEE/ACM International Symposium on Microarchitecture*, ACM, New York, NY, USA, 2011.
- [20] J. Mars, L. Tang, M.L. Soffa, [Directly characterizing cross core interference through contention synthesis](#), in: *Proceedings of the Sixth International Conference on High Performance and Embedded Architectures and Compilers, HiPEAC '11*, ACM, New York, NY, USA, 2011, pp. 167–176.
- [21] J. Mars, N. Vachharajani, R. Hundt, M. Soffa, Contention aware execution: online contention detection and response, in: *CGO '10: Proceedings of the Eighth Annual IEEE/ACM International Symposium on Code Generation and Optimization*, April, 2010.
- [22] A. Merkel, J. Stoess, F. Bellosa, Resource-conscious scheduling for energy efficiency on multicore processors, in: *EuroSys '10: Proceedings of the Fifth European Conference on Computer Systems*, April 2010.
- [23] K.J. Nesbit, N. Aggarwal, J. Laudon, J.E. Smith, Fair queuing memory systems, in: *MICRO 39: Proceedings of the 39th Annual IEEE/ACM International Symposium on Microarchitecture*, 2006.
- [24] K.J. Nesbit, J. Laudon, J.E. Smith, [Virtual private caches](#), in: *ISCA '07: Proceedings of the 34th Annual International Symposium on Computer Architecture*, vol. 35(2), 2007.
- [25] M.K. Qureshi, Y.N. Patt, Utility-based cache partitioning: a low-overhead, high-performance, runtime mechanism to partition shared caches, in: *MICRO 39: Proceedings of the 39th Annual IEEE/ACM International Symposium on Microarchitecture*, 2006.
- [26] N. Rafique, W.-T. Lim, M. Thottethodi, Architectural support for operating system-driven CMP cache management, in: *PACT '06: Proceedings of the 15th International Conference on Parallel Architectures and Compilation Techniques*, 2006.
- [27] L. Soares, D. Tam, M. Stumm, Reducing the harmful effects of last-level cache polluters with an OS-level, software-only pollute buffer, in: *MICRO '08: Proceedings of the 2008 41st IEEE/ACM International Symposium on Microarchitecture*, 2008.
- [28] G.E. Suh, S. Devadas, L. Rudolph, A new memory monitoring scheme for memory-aware scheduling and partitioning, in: *HPCA '02: Proceedings of the Eighth International Symposium on High-Performance Computer Architecture*, 2002.

- [29] L. Tang, J. Mars, M.L. Soffa, Compiling for niceness: mitigating contention for QoS in warehouse scale computers, in: CGO '12: Proceedings of the 2012 International Symposium on Code Generation and Optimization, ACM, New York, NY, USA, 2012.
- [30] L. Tang, J. Mars, N. Vachharajani, R. Hundt, M.L. Soffa, The impact of memory subsystem resource sharing on datacenter applications, in: ISCA '11: Proceeding of the 38th Annual International Symposium on Computer Architecture, ISCA '11, ACM, New York, NY, USA, 2011, pp. 283–294.
- [31] L. Tang, J. Mars, W. Wang, T. Dey, M.L. Soffa, Reqos: reactive static/dynamic compilation for QoS in warehouse scale computers, in: ASPLOS '13: Proceedings of the 18th International Conference on Architectural Support for Programming Languages and Operating Systems, ACM, 2013.
- [32] K. Tian, Y. Jiang, X. Shen, A study on optimally co-scheduling jobs of different lengths on chip multiprocessors, in: CF '09: Proceedings of the Sixth ACM Conference on Computing Frontiers, 2009.
- [33] Y. Xie, G.H. Loh, Dynamic classification of program memory behaviors in CMPs, The Second Workshop on Chip Multiprocessor Memory Systems and Interconnects, 2008.
- [34] D. Xu, C. Wu, P. Yew, On mitigating memory bandwidth contention through bandwidth-aware scheduling, in: Proceedings of the 19th International Conference on Parallel Architectures and Compilation Techniques, December 2010.
- [35] H. Yang, A. Breslow, J. Mars, L. Tang, Bubble-flux: precise on-line QoS management for increased utilization in warehouse scale computers, in: ISCA '13: Proceedings of the 40th Annual International Symposium on Computer Architecture, IEEE/ACM, 2013.
- [36] X. Zhang, S. Dwarkadas, K. Shen, Towards practical page coloring-based multicore cache management, in: EuroSys '09: Proceedings of the Fourth ACM European Conference on Computer Systems, 2009.
- [37] S. Zhuravlev, S. Blagodurov, A. Fedorova, Addressing shared resource contention in multicore processors via scheduling, in: ASPLOS '10: Proceedings of the 15th Edition of ASPLOS on Architectural Support for Programming Languages and Operating Systems, vol. 38, 2010.

## ABOUT THE AUTHORS

**Jason Mars** is an assistant professor of electrical engineering and computer science at the University of Michigan. His research interests include adaptive systems in software and hardware, warehouse-scale computer architecture, and software/hardware code-signed architectures. Mars has a Ph.D. in computer science from the University of Virginia. He is a member of IEEE and the ACM.

**Lingjia Tang** is an assistant professor of electrical engineering and computer science at the University of Michigan. Her research interests include compilers, runtime systems and computer architecture, especially for emerging systems such as large-scale datacenters and energy constrained mobile devices. Tang has a Ph.D. in computer science from the University of Virginia. She is a member of IEEE and the ACM.