

# MESS: Memory Performance Debugging on Embedded Multi-core Systems

Sudipta Chattopadhyay

Linköping University, Sweden

sudipta.chattopadhyay@liu.se

<https://bitbucket.org/sudiptac/mess>

**Abstract.** Multi-core processors have penetrated the modern computing platforms in several dimensions. Desktop machines, handheld devices and advanced embedded systems are now equipped with high-performance and energy-efficient multi-core systems. Multi-core systems aim to achieve high-performance via running computations in parallel. However, such systems also employ shared resources, such as shared caches and shared buses. The presence of parallel computations increases the congestion in such shared resources, leading to poor execution time. For embedded and real-time software, such performance loss is particularly undesirable. This is due to the reason that most embedded systems need to satisfy some extra-functional constraints, such as time.

In modern computing platforms, memory subsystems are several magnitudes slower than the processor. In order to bridge such performance gap between processor and memory, *caches* are employed by designers. Therefore, software performance may degrade drastically due to the congestion in shared caches, on multi-core systems. The congestion in the shared-cache critically depends on the order of memory-access operations. In this paper, we propose MESS, a performance debugging framework for embedded, multi-core systems. MESS systematically discovers the order of memory-access operations that expose performance bugs due to shared caches. We build a compositional approach that initially monitors the performance of each core *in isolation* and generates a performance summary for each core. Subsequently these summaries are used to build a constraint system. The solution of the constraint system reveals the interleaved memory-access-pattern that leads to a performance bug. Our baseline framework does not generate any *false positive*. Besides, its failure to find a solution highlights the absence of performance bugs due to shared caches, for a given input. Our baseline framework can also be employed to derive the memory access-order that leads to the worst-case shared-cache performance, for a given input. Finally, we propose an approximate solution that dramatically reduces debugging time, at the cost of a reasonable amount of false positives. We have implemented our entire framework using `SimpleScalar` simulator and `Z3` constraint solver. Experiments with several embedded software and a real-life robot controller suggest that we can discover performance bugs in a reasonable time.

## 1 Introduction

It is notoriously difficult to understand and discover performance bugs in software. Whereas performance bugs may appear in any application, these bugs are critical for

certain class of software, such as embedded and real-time software. Embedded and real-time applications are, in general, constrained via several temporal requirements. For hard real-time applications, violation of such temporal constraints may lead to catastrophic effects, often costing human lives. Apart from hard real-time applications, the existence of performance bugs may substantially impact the quality of soft real-time applications (*e.g.* media players) as well as web applications.

As the computing world is moving towards the multi-core era, it has become a critical problem to develop correct and efficient software on multi-core platforms. On the one hand, an application can utilize the potential of multi-core platforms by performing computations in parallel. On the other hand, parallel execution may dramatically increase the amount of non-determinism compared to a sequential execution. This happens due to the presence of an exponential number of interleaved execution patterns. Therefore, to validate the correctness and efficiency of applications on multi-core platforms, it is crucial that the validation methodologies consider different interleaving patterns. In this paper, broadly, we concentrate on the efficiency of applications which run on multi-core platforms.

In multi-threaded execution, software functionality might be disrupted due to the non-deterministic order in accessing shared data [13]. Similarly, the performance of multi-core systems may highly vary due to the non-deterministic order in accessing *shared resources*, such as shared caches. Caches are managed at runtime and they store copies of memory blocks from the main memory. In current generation computing platforms, caches are several magnitudes faster than accessing the main memory. As a result, cache memory is a crucial component to bridge the performance gap between the processor and main memory, and to improve the overall performance of applications. However, since caches are managed at runtime, the order of memory-access patterns play a crucial role in deciding the content of caches. For instance, consider a shared cache which can hold only one memory block. If memory accesses  $m_1$  and  $m_2$  are interleaved in parallel, the ordering  $(m_1 \cdot m_2)^*$  will *always* lead to cache misses. This is because  $m_1$  and  $m_2$  will always replace each other from the cache. On the contrary, for the ordering  $(m_1^* \cdot m_2^*)$ , *only the first* accesses of  $m_1$  and  $m_2$  will suffer cache misses. In summary, depending on the order of memory accesses, there might be a high variation on cache performance, which dramatically impacts the overall performance of software.

In this paper, we propose a novel approach to discover interleaving patterns that violate a given temporal constraint. For a given program input, our framework *automatically* discovers the order of memory accesses that highlights a performance bug. These bugs happen due to the sharing of caches in multi-core systems and they may lead to serious performance issues at runtime. A typical usage of our framework is the reproduction of performance bugs on multi-core systems and subsequently, improve the overall performance via classic cache management techniques, such as cache locking [22]. We leverage on the recent advances in constraint solving and *satisfiability modulo theory* (SMT) to systematically explore memory-access patterns. We propose a baseline framework which does not generate any *false* alarm. Moreover, if our baseline framework terminates without a solution, then we can guarantee the *absence* of performance bugs (*i.e.* violation of the temporal constraint), for the given input. We also propose an approximation that systematically partitions the set of constraints and solve each par-

tion in parallel. Such a strategy dramatically improves the solver performance. Our approximation guarantees *soundness*, meaning that the absence of a solution highlights the absence of performance bugs. However, the price of approximation might be paid via pessimism (*i.e.* false alarms). Our evaluation reveals that the magnitude of such pessimism is reasonable.

The generation of a performance-stressing interleaving pattern involves many technical challenges. Unlike the functionality of an application, its performance is not directly annotated in the code. Moreover, it is infeasible to execute an application for all possible interleaving patterns, due to an exponential number of possibilities. To resolve such challenges, we propose a compositional approach to discover performance bugs. Our framework broadly contains two stages. In the first stage, we monitor the performance of each core *in isolation*. Therefore, in this stage, we ignore interferences in the shared cache from other cores. The output of the first stage is a performance-summary for each core, where the timing to access the shared cache is replaced by a symbolic variable. In the second stage, we formulate constraints that relate the order of memory accesses with the delay to access the shared-cache. In particular, we formulate constraints that *symbolically encode necessary and sufficient conditions for a memory block to be evicted from the shared-cache*. As a result, using these constraints, we could determine whether a given memory block is available in the shared-cache, when it is being accessed. In other words, we can use such constraints to bound the delay to access the shared-cache and thereby, constraining the value of symbolic variables, which were introduced in the first stage of our framework. Finally, the temporal constraint is also provided as a quantifier-free formula. All the constraints, together with the temporal constraint, is given to an SMT solver. If the solver finds a solution, the resulting solution highlights an interleaving pattern that violates the temporal constraint. Since SMT technology is continuously evolving, we believe that such a compositional approach will be appealing to discover performance bugs in multi-core systems.

To tackle the complexity of our systems, we also propose an approximate solution that significantly improves the performance of our proposed framework. For shared caches, we observed that the set of all constraints can be partitioned systematically to solve in parallel. The general intuition is to consider partitions of memory accesses which can contend in the shared-cache and solve the constraints generated for each partition independently. By increasing the size of each partition, the designer can reduce the number of *false positives* at the cost of debugging time. Therefore, our framework gives designer the flexibility to fine tune the precision, with respect to debugging time.

*Contribution* In summary, we propose a performance debugging framework that exposes performance issues due to shared caches. We leverage on single-core performance profiling and symbolic-constraint solving, in order to discover the interleaving pattern that violates a given temporal constraint. Our baseline framework does not generate any *false positive* and it can also be used to prove the absence of performance bugs for a given input. Moreover, for time-critical code fragments, our baseline framework can be employed to derive the worst-case interleaving pattern (in terms of shared-cache performance), for a given input. To tackle the complexity of our constraint-based framework, we have also proposed an approximation that dramatically increases the solver performance. To show the generality of our approach, we have instantiated our framework

for two different caches (i) caches with least-recently-used (LRU) replacement policy and (ii) caches with first-in-first-out (FIFO) policy. We have implemented our entire framework on top of `simplescalar` [6] – an open-source, cycle-accurate, processor simulator and `Z3` [5] – an open source, SMT solver. Our experiments with several embedded software reveals the effectiveness of our approach. Last but not the least, we have evaluated our framework with a real-life robot controller (available in [3]). On average, our baseline framework was able to check a variety of temporal constraints for the controller within *3 minutes* and our approximation took only *20 seconds* on average to check the same set of constraints. This makes the idea of constraint-based formulation in performance debugging quite appealing for research in future.

## 2 Overview

In this Section, we shall first give a brief background on caches. Subsequently, we shall use a simple example to illustrate the challenges involved in accurately computing the shared-cache access delay on multi-core systems.

*Background on caches* Caches are employed between the CPU and the main memory (DRAM) to bridge the performance gap between the CPU and the DRAM. A cache can be described as a three tuple  $\langle \mathcal{A}, \mathcal{S}, \mathcal{L} \rangle$ , where  $\mathcal{A}$  is the associativity of the cache,  $\mathcal{S}$  is the number of cache sets and  $\mathcal{L}$  is the line size (in bytes). Each cache set can hold  $\mathcal{A}$  cache lines, leading to a total cache size of  $(\mathcal{A} \cdot \mathcal{S} \cdot \mathcal{L})$  bytes. When  $\mathcal{A} = 1$ , the respective caches are called to be *directly mapped*. Data is fetched into caches at the granularity of line size ( $\mathcal{L}$ ). Therefore, for an arbitrary memory address  $x$ ,  $\mathcal{L}$  contiguous bytes are fetched into the cache starting from address  $\lfloor \frac{x}{\mathcal{L}} \rfloor$  and we say that  $x$  belongs to the *memory block*  $\lfloor \frac{x}{\mathcal{L}} \rfloor$ . The number of cache sets ( $\mathcal{S}$ ) decides the location where a particular memory block would be placed in the cache. For instance, a memory block, starting at address  $M$ , is always mapped to the cache set  $M \bmod \mathcal{S}$ . Since each cache set can hold only  $\mathcal{A}$  cache lines, a cache line needs to be replaced when the number of memory blocks mapping to a cache set exceeds  $\mathcal{A}$ . In order to accomplish this, a replacement policy is employed when  $\mathcal{A} \geq 2$ . In this paper, we instantiate our framework for two widely used replacement policies – LRU and FIFO. In LRU policy, the memory block, that was not *accessed* for the longest period of time, is replaced from the cache to make room for other memory blocks. In FIFO policy, the memory block, which is *residing* in the cache for the longest period of time, is replaced to make room for other blocks. In general, the performance of a cache may greatly depend on the underlying replacement policy.

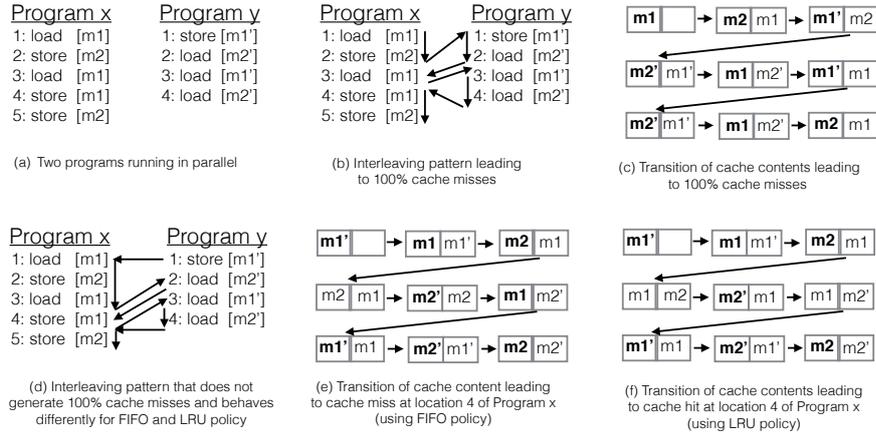
*Terminologies* We use the following terminologies on caches throughout the paper.

1. *memory block*: For an arbitrary memory reference to address  $x$ , we say that it belongs to memory block  $\lfloor \frac{x}{\mathcal{L}} \rfloor$  ( $\mathcal{L}$  is the line size of cache, in bytes), in order to distinguish different cache lines.
2. *cache hit/miss*: For an arbitrary memory reference, we say that it is a cache hit (miss) if the referenced memory block is found (not found) in the cache.

3. *cache conflict*: Two memory blocks  $M_1$  and  $M_2$  conflict in the cache if they map to the same cache set. In other words,  $M_1$  is conflicting to  $M_2$  (and vice versa). These conflicting memory blocks might be accessed within the same core (intra-core) or across different cores (inter-core).
4. *cache-set state*: Ordered  $\mathcal{A}$ -tuple capturing the content of a cache set. For instance,  $\langle m_1, m_2 \rangle$  captures such a tuple for caches with associativity 2. The relative position of a memory block in the tuple decides the number of unique cache conflicts required to evict the same from the cache. For instance, in  $\langle m_1, m_2 \rangle$ ,  $m_1$  requires two unique cache conflicts to be evicted from the cache, whereas  $m_2$  requires only one. The generation of cache conflicts critically depends on the replacement policy and the order of memory accesses.

### Motivation and challenges

Figure 1 captures an example where two programs are executing in parallel on different processor cores and sharing a cache. In general, multi-core systems employ several levels of private caches and a last-level shared cache (e.g. ARM MPCore). For the sake of simplicity, let us assume that all the instructions in both Program x and Program y access the same shared-cache set. In Figure 1(a), the memory block accessed by each instruction is shown within the brackets. In the following discussion, we shall capture the location  $i$  of Program x via  $x^i$  and the same of Program y via  $y^i$ .



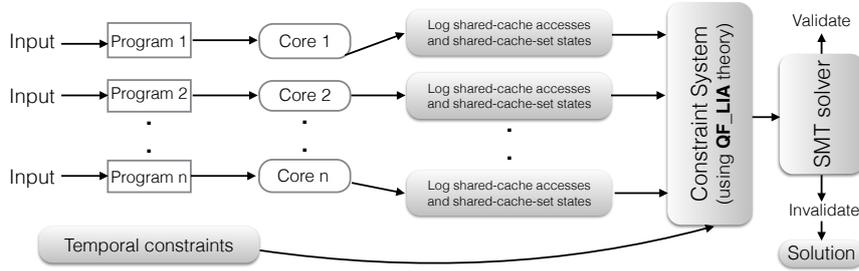
**Fig. 1.** An example showing the impact of interleaving pattern on shared-cache performance. The direction of an arrow captures the *happens-before* relation. Cache misses are highlighted in bold.

Let us assume a cache with associativity ( $\mathcal{A}$ ) two and employing FIFO replacement policy. The delay to access the shared-cache might highly vary due to an exponential number of possible interleaving patterns. For instance, let us first assume that we want to check whether *all instructions in both programs can face cache misses*. Figure 1(b) captures an interleaving pattern which leads to 100% cache misses in both programs. The progression of the cache content for this interleaving pattern is captured via Figure 1(c). It is worthwhile to note that many interleaving patterns will fail to generate

100% cache misses in both programs. Figure 1(d) captures one such interleaving pattern. As a result, if the set of memory accesses (*cf.* Figure 1(a)) appears within a loop, the memory-access delay might change dramatically depending on the interleaving pattern. This is due to the huge performance gap between processor and memory. The respective cache contents for the interleaving pattern in Figure 1(d) are shown via Figure 1(e). In general, it is infeasible to perform an exhaustive search over the set of all possible interleaving patterns, due to an exponential number of possibilities. As a result, a systematic method is required to check performance-related constraints, in the context of multi-core systems.

Let us now assume that we want to check whether location  $x^4$  can face a *cache miss*. Such a behaviour can also take place only for a few interleaving patterns. Figure 1(d) captures an interleaving pattern which lead to a cache miss at location  $x^4$  (*cf.* Figure 1(e) for the transition of cache contents). Unfortunately, if we replay the same interleaving pattern for LRU replacement policy, it will not lead to a cache miss at location  $x^4$ . This behaviour is captured via Figure 1(f), which demonstrates the modification of cache contents in the presence of LRU policy. This shows the influence of the *cache replacement policy* to check or invalidate temporal constraints.

To summarize, due to the presence of shared caches in multi-core systems, it is challenging to check the validity of temporal constraints or reproduce any violation of temporal constraints in a production run. This phenomenon occurs due to the non-determinism in the order of interleaved memory-accesses, which, in turn leads to non-determinism in cache contention and variability in memory-access delay. In the following, we shall give an outline of our performance debugging framework.



**Fig. 2.** Performance debugging framework for multi-core systems

*Overall framework* Figure 2 outlines our overall strategy in order to check the validity or violation of temporal constraints on multi-core systems. For a given input to each program running in parallel, our framework is used to check the temporal constraints. We monitor the execution on each core *in isolation*, ignoring any interference from other cores. As a result, this monitoring phase can be carried out in parallel for each core. While monitoring, we record accesses to the shared cache by each core. This can be performed either via instrumentation on real hardware or using a simulator which models the underlying memory hierarchy. At the end of the monitoring phase, we obtain a sequence of shared-cache accesses  $\langle i^1, i^2, \dots, i^{\mathcal{V}_i-1}, i^{\mathcal{V}_i} \rangle$  for each core  $i$ , where  $\mathcal{V}_i$  is the total number of shared-cache accesses by core  $i$ . We also collect the shared-cache-set states at these access points. Using the information obtained from the monitoring phase,

we build a constraint system. Intuitively, this constraint system relates the order of memory accesses with the delay to access the shared cache. The size of our constraint system is *polynomial*, with respect to the number of accesses to the shared cache. Finally, the temporal constraint can be provided to the constraint system via quantifier-free predicates. The entire constraint system, along with the temporal constraints, is provided to an SMT solver. If the constraint system is *satisfiable*, then the solution returned by the SMT solver captures an interleaving pattern that *violates* certain temporal constraints. This solution can further be used for debugging performance on multi-core systems.

*System model* We assume a sequentially-consistent, multi-core system where *each core may have several levels of private caches and only the last-level cache is shared across cores*. Therefore, a shared-cache miss will lead to an access to the slow DRAM. Such a design of memory-hierarchy is typical in embedded multi-core processors [1]. In this paper, we do not address the problem of cache coherency and any cache misses resulting from the same. Such cache misses might appear due to the invalidation of cache lines that hold outdated data. Besides, additional cache misses might appear due to *false sharing* [2]. In summary, we first assume that programs, running on different cores, have disjoint memory spaces. We argue that, even in the absence of cache coherency, debugging shared-cache performance is sufficiently complex. In Section 4, we shall discuss the required modifications in our framework in the presence of data sharing. In the following, we first present our baseline framework. Subsequently, we shall discuss an approximation scheme to reduce the debugging time.

### 3 Methodologies

In this section, we shall introduce the formal foundation of our framework. Recall that the outcome of our framework is to compute a memory-access ordering, leading to a specific performance problem (specified by a quantifier-free predicate). This ordering is captured among all accesses to the shared cache.

Let us assume that we have a total of  $\mathcal{N}$  cores, each of which might exhibit a different sequence of shared-cache accesses. We use the notation  $i^j$  to capture the  $j$ -th shared-cache access by  $i$ -th core and  $\mathcal{V}_i$  to capture the total number of shared-cache accesses by core  $i$ . Besides, we shall use the following notations to formulate our performance debugging framework.

- $\sigma_i^j$  : The memory block accessed by the shared-cache access  $i^j$ .
- $\pi(m)$  : Cache set where memory block  $m$  is mapped.
- $\zeta_i^j$  : Shared-cache-set state for cache set  $\pi(\sigma_i^j)$ , immediately before the access  $i^j$ .
- $\mathcal{C}_i^j$  : The set of memory blocks, other than  $\sigma_i^j$ , mapping to the same cache set as  $\sigma_i^j$  in the shared cache. Therefore, for any  $m' \in \mathcal{C}_i^j$ , we have  $m' \neq \sigma_i^j$  and  $\pi(m') = \pi(\sigma_i^j)$ .
- $\mathcal{O}_i^j$  : The position of the shared-cache access  $i^j$  in the ordering among all accesses to the shared cache.
- $\delta_i^j$  : The delay suffered by the shared-cache access  $i^j$ .

For instance, in Figure 1(b),  $\sigma_x^1 = m1$ ,  $\sigma_y^1 = m1'$ ,  $\zeta_y^1 = \langle m2, m1 \rangle$  and the interleaving pattern is captured as follows:  $\mathcal{O}_x^1 < \mathcal{O}_x^2 < \mathcal{O}_y^1 < \mathcal{O}_y^2 < \mathcal{O}_x^3 < \mathcal{O}_y^3 < \mathcal{O}_y^4 < \mathcal{O}_x^4 < \mathcal{O}_x^5$ . The outcome of our framework is such an interleaving pattern.

**Profiling each core in isolation** As outlined in the preceding section, our framework initially records the performance of each core *in isolation*. The primary purpose of this recording phase is to accurately identify accesses to the shared cache, for each core. Therefore, while profiling each core in isolation,  $\zeta_i^j$  contains memory blocks accessed *only* within core  $i$  and ignores all memory blocks accessed within core  $\bar{i} \neq i$ .

Let us assume  $age_i^j$  denotes the relative position of  $\sigma_i^j$  within  $\zeta_i^j$ , while profiling each core in isolation. If  $\sigma_i^j \notin \zeta_i^j$  (i.e.  $i^j$  suffers a shared-cache miss), we assign  $age_i^j$  to  $\mathcal{A} + 1$ , where  $\mathcal{A}$  is the associativity of the shared-cache. Subsequently, for each core  $i$ , we encode a performance-summary  $\alpha_i$  as a sequence of pairs. Each such pair captures a shared-cache access  $i^j$ , along with  $age_i^j$  as follows:

$$\alpha_i \equiv \langle (i^1, age_i^1), (i^2, age_i^2), \dots, (i^{\mathcal{V}_i-1}, age_i^{\mathcal{V}_i-1}), (i^{\mathcal{V}_i}, age_i^{\mathcal{V}_i}) \rangle \quad (1)$$

For any shared-cache access  $i^j$ , it is a shared-cache miss if and only if  $\sigma_i^j \notin \zeta_i^j$ , leading to the value of  $age_i^j$  set to  $\mathcal{A} + 1$ . Such a cache miss can happen because of the following reasons:

1.  $\sigma_i^j$  was accessed for the first time
2.  $\sigma_i^j$  was evicted from the shared-cache by some other memory block

Recall that programs running on different cores have disjoint memory spaces. As a result, while profiling each core in isolation, we can accurately identify shared-cache misses when  $\sigma_i^j$  was accessed for the first time. This is because,  $\sigma_i^j$  was not accessed by any other core except core  $i$ . In subsequent sections, we shall describe our constraint system, which formulates necessary and sufficient conditions for evicting memory blocks from the shared-cache, leading to shared-cache misses.

**Program order constraints** These constraints are generated to capture the program order on each core. Note that  $\langle i^1, i^2, \dots, i^{\mathcal{V}_i-1}, i^{\mathcal{V}_i} \rangle$  captures the sequence of shared-cache accesses by core  $i$ . Therefore, the following constraints are generated to capture the program order.

$$\Theta_{order} \equiv \bigwedge_{i \in [1, \mathcal{N}]} \bigwedge_{j \in [2, \mathcal{V}_i]} \left( \mathcal{O}_i^j > \mathcal{O}_i^{j-1} \right) \quad (2)$$

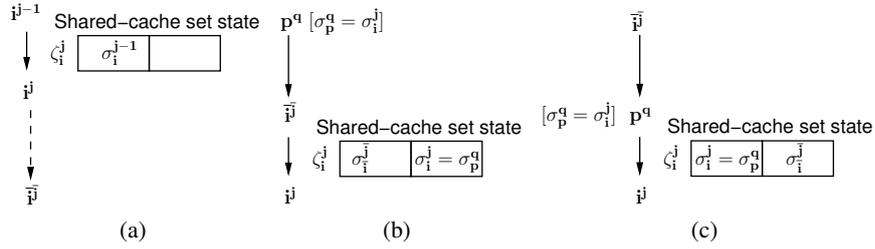
Program-order constraints are generated irrespective of the cache replacement policy. In the following, we now instantiate the constraint formulation for LRU and FIFO policies.

### 3.1 Constraint system for LRU caches

A shared-cache access  $i^j$  is a cache hit if and only if  $\zeta_i^j$  contains  $\sigma_i^j$ . Otherwise,  $i^j$  suffers a shared-cache miss. Therefore, to accurately determine the shared-cache performance, it is crucial to track all feasible states of  $\zeta_i^j$ . We accomplish this by relating the order of memory accesses with the changes in cache-set states. In order to understand the relationship between the memory-access order and cache-set states, we first define the notion of cache-conflict generation between two shared-cache accesses.

**Definition 1** (*Cache Conflict Generation*) Consider a shared-cache access  $\bar{i}^j$ , which requests memory block  $\bar{m}$  (i.e.  $\sigma_{\bar{i}}^j = \bar{m}$ ). A shared-cache access  $\bar{i}^j$  generates (cache) conflict to  $i^j$ , only if accessing  $\sigma_{\bar{i}}^j$  can affect the relative position of  $\sigma_i^j$  within  $\zeta_i^j$ . For instance, in Figure 1(d), accesses to  $m1'$  and  $m2'$  do not generate cache conflict to  $x^3$ , but an access to  $m2$  does (at  $x^2$ ).

We introduce a variable  $\Psi_i^j(\bar{m})$  to capture whether any access to memory block  $\bar{m}$  generates conflict to the shared-cache access  $i^j$ . As stated in Definition 1, the memory block  $\bar{m}$  might be accessed more than once and therefore, the formulation of  $\Psi_i^j(\bar{m})$  must consider all possible places where  $\bar{m}$  was accessed. Consider one such place  $\bar{i}^j$ , where  $\bar{m}$  was accessed. Therefore,  $\sigma_{\bar{i}}^j = \bar{m}$ . Figure 3 illustrates different scenarios in LRU policy, with respect to the generation of cache conflicts.



**Fig. 3.** The direction of arrow captures the total order between accesses to the shared cache. The left-most position in  $\zeta_i^j$  captures the *most recently used* memory block. (a)  $\bar{i}^j$  cannot affect shared-cache set state  $\zeta_i^j$  and therefore, it cannot generate cache conflict to  $i^j$ , if  $\bar{i}^j$  happens after  $i^j$ . (b)  $\bar{i}^j$  can affect  $\zeta_i^j$  only if  $\bar{i}^j$  happens before  $i^j$ . (c) shared-cache access  $p^q$  accesses the same memory block as that of  $i^j$  (i.e.  $\sigma_p^q = \sigma_i^j$ ) and therefore, access  $\bar{i}^j$  cannot affect the relative position of  $\sigma_i^j$  within  $\zeta_i^j$ .

In particular, Figures 3(a)-(b) capture the *happens-before* relationship between accesses  $\bar{i}^j$  and  $i^j$ . It is impossible for  $\bar{i}^j$  to affect the cache-set state  $\zeta_i^j$ , if  $i^j$  happens before  $\bar{i}^j$ . Moreover, if the memory block  $\sigma_i^j$  is accessed after  $\bar{i}^j$  and before  $i^j$ , then such an access will hide the cache conflict between  $\bar{i}^j$  and  $i^j$ . Figure 3(c) captures one such situation, where shared-cache access  $p^q$  accesses the memory block  $\sigma_i^j$  and prevents  $\bar{i}^j$  to affect the relative position of  $\sigma_i^j$  within cache-set state  $\zeta_i^j$ .

In the following, we describe the formulation of constraints for an arbitrary shared-cache access  $i^j$ . The primary purpose of these constraints is to compute the delay  $\delta_i^j$ . Considering the intuition provided in Figure 3, we can state that a shared-cache access  $\bar{i}^j$  generates conflict to the shared-cache access  $i^j$ , only if the following conditions hold:

- $\psi_{cft}^{lru}(\bar{i}^j, i^j)$  : Shared-cache access  $\bar{i}^j$  happens before the shared-cache access  $i^j$ . Therefore,  $\mathcal{O}_{\bar{i}}^j < \mathcal{O}_i^j$ . This is illustrated via Figures 3(a)-(b).
- $\psi_{ref}^{lru}(\bar{i}^j, i^j)$  : There does not exist any shared-cache access  $p^q$ , such that  $p^q$  accesses memory block  $\sigma_i^j$  from the shared-cache,  $p^q$  happens before  $i^j$  and  $\bar{i}^j$  happens before  $p^q$ . Therefore, for any shared-cache access  $p^q$ , where  $\sigma_p^q = \sigma_i^j$ , conditions  $\mathcal{O}_p^q < \mathcal{O}_i^j$  and  $\mathcal{O}_{\bar{i}}^j < \mathcal{O}_p^q$  cannot be satisfiable together. Otherwise, note that  $p^q$  will hide the cache conflict between  $\bar{i}^j$  and  $i^j$ , as illustrated via Figure 3(c).

$\psi_{cft}^{lru}(\bar{i}^j, i^j)$  and  $\psi_{ref}^{lru}(\bar{i}^j, i^j)$  can be formalized via the following constraints:

$$\psi_{cft}^{lru}(\bar{i}^j, i^j) \equiv \mathcal{O}_{\bar{i}^j} < \mathcal{O}_i^j \quad (3)$$

$$\psi_{ref}^{lru}(\bar{i}^j, i^j) \equiv \bigwedge_{p,q: \sigma_p^q = \sigma_i^j} \neg (\mathcal{O}_{\bar{i}^j} < \mathcal{O}_p^q \wedge \mathcal{O}_p^q < \mathcal{O}_i^j) \quad (4)$$

We combine Constraint (3) and Constraint (4) to formulate the generation of shared-cache conflict. Recall that  $\mathcal{C}_i^j$  captures the set of memory blocks that map to the same shared-cache set as  $\sigma_i^j$ . Therefore, Constraints (3)-(4) need to be generated for each memory block in  $\mathcal{C}_i^j$ . Formally, for each shared-cache access  $i^j$ , we generate the following constraints to capture cache conflicts generated across cores.

$$\Theta_1^{lru}(i, j) \equiv \bigwedge_{\bar{i} \neq i: \sigma_{\bar{i}}^j \in \mathcal{C}_i^j} \left( (\psi_{cft}^{lru}(\bar{i}^j, i^j) \wedge \psi_{ref}^{lru}(\bar{i}^j, i^j)) \Rightarrow (\Psi_i^j(\sigma_{\bar{i}}^j) = 1) \right) \quad (5)$$

The absence of inter-core cache conflict is captured via the negation of Constraint (5). In particular, for any memory block  $\bar{m} \in \mathcal{C}_i^j$ , we need to consider the set of locations  $\bar{i}^j$  where  $\bar{m}$  is accessed (*i.e.*  $\sigma_{\bar{i}^j}^j = \bar{m}$ ). If none of these locations satisfy either Constraint (3) or Constraint (4), we can conclude that accesses to memory block  $\bar{m}$  do not generate any cache conflict at location  $i^j$ . This behaviour can be captured via the following constraints:

$$\Theta_0^{lru}(i, j) \equiv \bigwedge_{\bar{m} \in \mathcal{C}_i^j} \left( \bigwedge_{\bar{i} \neq i: \sigma_{\bar{i}}^j = \bar{m}} (\neg \psi_{cft}^{lru}(\bar{i}^j, i^j) \vee \neg \psi_{ref}^{lru}(\bar{i}^j, i^j)) \Rightarrow (\Psi_i^j(\bar{m}) = 0) \right) \quad (6)$$

Finally, we need to link Constraints (5)-(6) to the absolute latency suffered by shared-cache access  $i^j$  (*i.e.*  $\delta_i^j$ ). Let us assume *HIT* and *MISS* capture the shared-cache hit latency and miss penalty, respectively. To compute the latency, we need to check whether the set of cache conflicts generated at  $i^j$  could evict the memory block  $\sigma_i^j$ . Therefore, we generate the following constraints to formulate the delay suffered at location  $i^j$ .

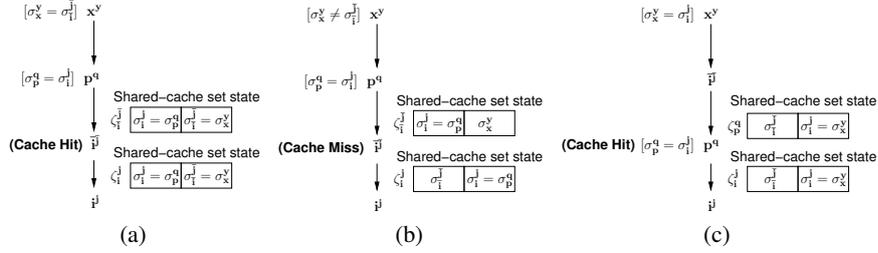
$$\Theta_{miss}^{lru}(i, j) \equiv \left( \sum_{\bar{i} \neq i: \sigma_{\bar{i}}^j \in \mathcal{C}_i^j} \Psi_i^j(\sigma_{\bar{i}}^j) \geq \mathcal{A} - age_i^j + 1 \right) \Rightarrow (\delta_i^j = MISS) \quad (7)$$

$$\Theta_{hit}^{lru}(i, j) \equiv \left( \sum_{\bar{i} \neq i: \sigma_{\bar{i}}^j \in \mathcal{C}_i^j} \Psi_i^j(\sigma_{\bar{i}}^j) \leq \mathcal{A} - age_i^j \right) \Rightarrow (\delta_i^j = HIT) \quad (8)$$

$age_i^j$  denotes the relative position of  $\sigma_i^j$  within  $\zeta_i^j$  and  $age_i^j = \mathcal{A} + 1$ , if  $\sigma_i^j \notin \zeta_i^j$ . The value  $age_i^j$  was collected while profiling each core *in isolation* (*cf.* Equation (1)). Therefore,  $age_i^j$  already captures cache conflicts generated within core  $i$  and the quantity  $(\mathcal{A} - age_i^j + 1)$  captures the minimum number of unique, inter-core cache conflicts (as formulated via Constraint (5)) to evict  $\sigma_i^j$  from the shared cache.

### 3.2 Constraint system for FIFO caches

Unlike LRU policy, cache-set state remains unchanged for all *cache hits* in FIFO policy (cf. Figure 1(e)). As a result, the necessary conditions to generate cache conflicts (cf. Constraints (3)-(4)) need to be modified for FIFO policy.



**Fig. 4.** The direction of arrow captures the total order between accesses to the shared cache. The left-most position in  $\zeta_i^j$  captures the *most recent memory block inserted* into  $\zeta_i^j$ . (a)  $\bar{i}^j$  cannot affect shared-cache set state  $\zeta_i^j$  as  $\bar{i}^j$  is a cache hit. Therefore,  $\bar{i}^j$  cannot generate cache conflict to  $i^j$ , (b)  $\bar{i}^j$  can affect  $\zeta_i^j$  only if  $\bar{i}^j$  happens before  $i^j$  and it is a cache miss, (c) shared-cache access  $p^q$  accesses the same memory block as that of  $i^j$  (i.e.  $\sigma_p^q = \sigma_i^j$ ), however,  $p^q$  is a cache hit. Therefore,  $p^q$  cannot hide the cache conflict generated between  $\bar{i}^j$  and  $i^j$ .

To illustrate the difference between LRU and FIFO policy, let us consider the scenarios in Figure 4. For instance, in Figure 4(a), shared-cache access  $\bar{i}^j$  happens before the access  $i^j$ . However,  $\bar{i}^j$  cannot affect the relative position of  $\sigma_i^j$  within  $\zeta_i^j$  and therefore,  $\bar{i}^j$  cannot generate cache conflict to  $i^j$  (cf. Definition 1). It is worthwhile to note that,  $\bar{i}^j$  would have generated conflict to  $i^j$ , in the presence of LRU policy. Figure 4(b) captures a scenario, where  $\bar{i}^j$  was a cache miss, leading to the generation of cache conflict to  $i^j$ . Recall that, for LRU policy, if the memory block  $\sigma_i^j$  was accessed between  $\bar{i}^j$  and  $i^j$ , then  $\bar{i}^j$  could not generate cache conflict to  $i^j$  (cf. Constraint (4)). However in FIFO policy, as shown in Figure 4(c), even though access  $p^q$  references  $\sigma_i^j$  and it occurs between  $\bar{i}^j$  and  $i^j$ ,  $p^q$  cannot hide the cache conflict between  $\bar{i}^j$  and  $i^j$ . This is because  $p^q$  was a cache hit and therefore, it does not affect the relative position of  $\sigma_i^j$  within  $\zeta_i^j$ .

In summary, a shared-cache access must be a cache miss if it affects the cache-set state  $\zeta_i^j$ . In order to realize this intuition, we formulate the following constraints, which capture the necessary conditions for  $\bar{i}^j$  generating cache conflict to  $i^j$ .

$$\psi_{cft}^{fif}(\bar{i}^j, i^j) \equiv (\mathcal{O}_i^{\bar{j}} < \mathcal{O}_i^j) \wedge (\delta_i^{\bar{j}} = \text{MISS}) \quad (9)$$

$$\psi_{ref}^{fif}(\bar{i}^j, i^j) \equiv \bigwedge_{p,q: \sigma_p^q = \sigma_i^j} \neg \left( (\mathcal{O}_i^{\bar{j}} < \mathcal{O}_p^q) \wedge (\mathcal{O}_p^q < \mathcal{O}_i^j) \wedge (\delta_p^q = \text{MISS}) \right) \quad (10)$$

Constraint (9) ensures that  $\bar{i}^j$  incurs a cache miss, in order to generate cache conflict to  $i^j$  (cf. Figures 4(a)-(b)). Similarly, Constraint (10) ensures that access  $p^q$  needs to be a cache miss to hide the cache conflict between  $\bar{i}^j$  and  $i^j$  (cf. Figure 4(c)).

The outcome of Constraints (9)-(10) may depend on the interleaving pattern, even within a single core (i.e.  $\bar{i} = i$ ). This is because, values of  $\delta_i^{\bar{j}}$  and  $\delta_p^q$  may depend on the

interleaving pattern. As a result, the generation of cache conflicts, even within a core, may be affected with FIFO policy. Hence, unlike LRU policy, we need to formulate cache conflict both within a core and across cores. This is accomplished by modifying Constraints (5)-(6), so that the resulting constraints also consider cache conflicts within cores. In particular, we remove the condition  $\bar{i} \neq i$  from Constraints (5)-(6) as follows.

$$\Theta_1^{fiffo}(i, j) \equiv \bigwedge_{\bar{i}, \bar{j}: \sigma_{\bar{i}}^{\bar{j}} \in \mathcal{C}_i^j} \left( \left( \psi_{cft}^{fiffo}(\bar{i}^{\bar{j}}, i^j) \wedge \psi_{ref}^{fiffo}(\bar{i}^{\bar{j}}, i^j) \right) \Rightarrow \left( \Psi_i^j(\sigma_{\bar{i}}^{\bar{j}}) = 1 \right) \right) \quad (11)$$

$$\Theta_0^{fiffo}(i, j) \equiv \bigwedge_{\bar{m} \in \mathcal{C}_i^j} \left( \bigwedge_{\bar{i}, \bar{j}: \sigma_{\bar{i}}^{\bar{j}} = \bar{m}} \left( \neg \psi_{cft}^{fiffo}(\bar{i}^{\bar{j}}, i^j) \vee \neg \psi_{ref}^{fiffo}(\bar{i}^{\bar{j}}, i^j) \right) \Rightarrow \left( \Psi_i^j(\bar{m}) = 0 \right) \right) \quad (12)$$

Finally, we link Constraints (11)-(12) to compute the memory-access latency. Intuitively, we check whether the total amount of cache conflict can evict the memory block accessed by  $i^j$ . This can be formalized via the following constraints.

$$\Theta_{miss}^{fiffo}(i, j) \equiv \left( \left( \sum_{\bar{m} \in \mathcal{C}_i^j} \Psi_i^j(\bar{m}) \geq \mathcal{A} \right) \vee \left( age_i^j = \mathcal{A} + 1 \right) \right) \Rightarrow \left( \delta_i^j = MISS \right) \quad (13)$$

$$\Theta_{hit}^{fiffo}(i, j) \equiv \left( \left( \sum_{\bar{m} \in \mathcal{C}_i^j} \Psi_i^j(\bar{m}) < \mathcal{A} \right) \wedge \left( age_i^j \neq \mathcal{A} + 1 \right) \right) \Rightarrow \left( \delta_i^j = HIT \right) \quad (14)$$

$\mathcal{A}$  is the associativity of the cache. Recall that  $age_i^j = \mathcal{A} + 1$ , if  $\sigma_i^j \notin \zeta_i^j$  and  $age_i^j$  was measured while investigating each core in isolation (cf. Equation 1). Therefore, the condition  $age_i^j = \mathcal{A} + 1$  guarantees to include the first-ever cache miss of  $\sigma_i^j$ . Once  $\sigma_i^j$  enters the cache, it takes at least  $\mathcal{A}$  unique cache-conflicts to evict it from the cache.  $\sum_{\bar{m} \in \mathcal{C}_i^j} \Psi_i^j(\bar{m})$  accounts all unique cache-conflicts faced by  $\sigma_i^j$ , since it enters the cache and till  $i^j$ . Therefore, Constraint (13) precisely captures all possibilities of a cache miss at  $i^j$ . The violation of Constraint (13) will result in a cache hit at  $i^j$ , as shown in Constraint (14).

**Providing temporal constraints** For embedded software, temporal constraints can be provided in the form of an assertion. Therefore, our framework will search for an ordering on symbolic variables  $\mathcal{O}_i^j$  that violates such assertions. In particular, we consider assertions that check the execution time against a threshold  $\tau$ . In our framework, the non-determinism in timing behaviour appears due to the accesses to shared caches. Therefore, in our evaluation, we search for a solution that satisfy the following constraint:  $\left( \sum_{i,j} \delta_i^j \geq \tau \right)$ . Recall that  $\delta_i^j$  symbolically captures the delay suffered by shared-cache access  $i^j$ . It is worthwhile to note that we can also check the timing behaviour of a code fragment, instead of checking the same for the entire system. In such cases, we consider only a subset of  $\delta_i^j$  variables relating to the code fragment.

**Putting it all together** Our formulated constraints, along with the temporal constraint, is provided to an off-the-shelf SMT solver. As a result, any ongoing and future improvements in the solver technology will directly boost the efficiency of our approach. The SMT solver searches for a satisfying solution of the following constraints:

$$\Phi \equiv \Theta_{order} \wedge \bigwedge_{i,j} (\Theta_1^x(i,j) \wedge \Theta_0^x(i,j) \wedge \Theta_{miss}^x(i,j) \wedge \Theta_{hit}^x(i,j)) \wedge \left( \sum_{i,j} \delta_i^j \geq \tau \right) \quad (15)$$

where  $x \in \{lru, fifo\}$ , depending on the cache replacement policy. The solution of the solver captures concrete values of symbolic variables  $\mathcal{O}_i^j$  that satisfy  $\Phi$ . Such concrete values can be used to derive the total order among all accesses to the shared-cache.

**Complexity of constraints** The complexity of our constraints  $\Phi$  (cf. Constraint (15)) is dominated by the number of constraints to formulate cache conflicts. For instance, in LRU policy, Constraints (5)-(6) dominate the total number of constraints. Let us assume that the total number of shared-cache accesses across all cores is  $\mathcal{K}$ . Therefore, the size of Constraints (2) has a complexity of  $O(\mathcal{K})$ . Similarly, the total size of Constraints (7)-(8), for LRU policy (respectively, the total size of Constraints (13)-(14) for FIFO policy) has a size of  $O(\mathcal{K})$ . Finally, during the formulation of cache conflict, each shared-cache access can be compared with all conflicting shared-cache accesses. Therefore,  $\Theta_1^{lru}(i,j)$ ,  $\Theta_0^{lru}(i,j)$ ,  $\Theta_1^{fifo}(i,j)$  and  $\Theta_0^{fifo}(i,j)$  have a worst-case size-complexity  $O(\mathcal{K}^2)$ . Since there exists a total of  $\mathcal{K}$  shared-cache accesses, the total size of Constraints (5)-(6) has a complexity of  $O(\mathcal{K}^3)$ . Putting everything together, our constraint system has a worst-case size-complexity  $O(\mathcal{K}^3)$ . However, our evaluation reveals that the size of our constraint system is substantially lower than the worst-case complexity.

### 3.3 Approximate solution

In the preceding, we observed that the number of constraints is polynomial with respect to the number of accesses to the shared cache. However, for long traces, the solver may impose a bottleneck to handle a large number of constraints. To address this, we propose an approximation, which potentially improves the scalability by several magnitudes. The general intuition of our approximation is based on the design principle of caches. In particular, we leverage the fact that two different cache sets never interfere with each other, in terms of cache conflict. Therefore, we model the constraints for each cache set separately and solve them in parallel. In the following, we shall formalize the concept.

**Finding a slice of constraints** The key idea for the approximation is to find a slice of constraints that could be solved *independently*. Recall that the symbolic variable  $\delta_i^j$  captures the delay suffered by shared-cache access  $i^j$ . It is worthwhile to note that the memory block accessed at  $i^j$  (i.e.  $\sigma_i^j$ ) can be evicted from the shared-cache only by memory blocks conflicting to  $\sigma_i^j$ . A memory block  $\bar{m}$  conflicts to  $\sigma_i^j$  in the cache if and only if  $\bar{m}$  and  $\sigma_i^j$  map to the same cache set. Therefore, we first group shared-cache accesses with respect to different cache sets and generate the respective constraints. For instance, consider that we are generating constraints with respect to cache set  $s$ . We shall use  $\pi(m)$  to capture the cache set in which memory block  $m$  is mapped.

We slice out the program-order constraints by considering only the memory blocks which map to cache set  $s$ . Therefore, the set of program-order constraints, with respect to cache set  $s$ , can be defined as follows.

$$\Gamma_{order}(s) \equiv \bigwedge_{i \in [1, \mathcal{N}]} \left( \bigwedge_{j, k \in [1, \mathcal{V}_i]: j < k \wedge (\pi(\sigma_i^j) = \pi(\sigma_i^k) = s) \wedge (\forall m \in [j+1, k]: \pi(\sigma_i^j) \neq \pi(\sigma_i^m))} \mathcal{O}_i^k > \mathcal{O}_i^j \right) \quad (16)$$

Let us now consider LRU cache replacement policy. The set of constraints, with respect to cache set  $s$ , considers constraints that only influence the memory blocks mapped to cache set  $s$ . Therefore, for cache set  $s$ , we extract the constraints formulated in Equations (5)-(8) as follows.

$$\Gamma_1^{lru}(s) \equiv \bigwedge_{i, j: \pi(\sigma_i^j) = s} \Theta_1^{lru}(i, j); \quad \Gamma_0^{lru}(s) \equiv \bigwedge_{i, j: \pi(\sigma_i^j) = s} \Theta_0^{lru}(i, j) \quad (17)$$

$$\Gamma_{miss}^{lru}(s) \equiv \bigwedge_{i, j: \pi(\sigma_i^j) = s} \Theta_{miss}^{lru}(i, j); \quad \Gamma_{hit}^{lru}(s) \equiv \bigwedge_{i, j: \pi(\sigma_i^j) = s} \Theta_{hit}^{lru}(i, j) \quad (18)$$

Finally, we gather all constraints with respect to cache set  $s$ . Our goal is to maximize the delay faced by accessing memory blocks mapped to  $s$ . This is performed via the following constraints and objective function.

$$\Gamma(s) \equiv \Gamma_{order}(s) \wedge \Gamma_1^{lru}(s) \wedge \Gamma_0^{lru}(s) \wedge \Gamma_{miss}^{lru}(s) \wedge \Gamma_{hit}^{lru}(s) \quad (19)$$

$$\Delta(s) = \text{maximize} \quad \sum_{i, j: \pi(\sigma_i^j) = s} \delta_i^j \quad (20)$$

Note that  $\Gamma(s)$  includes *all* constraints that could influence  $\Delta(s)$ . We can use recent development in SMT solving [21] to maximize the objective function captured via Equation (20). It is also worthwhile to mention that the preceding process can be carried out in an exactly same fashion for FIFO policy. As a result, our approximation strategy is generic, with respect to the replacement policy employed in a cache.

For each cache set  $s$ , we formulate  $\Gamma(s)$  and obtain the value of  $\Delta(s)$  using [21]. If  $s_1, s_2, \dots, s_q$  are all different sets in the shared cache,  $\sum_{r \in [1, q]} \Delta(s_r)$  over-approximates the total delay in accessing the shared cache. More precisely, we state the crucial property of our approximation scheme as follows (see Appendix for the proof).

**Property 2** *Let us assume  $\{s_1, s_2, \dots, s_q\}$  are different sets in the shared cache. For a given temporal constraint  $\sum_{i, j} \delta_i^j < \tau$ , if our baseline constraint system  $\Phi$  (cf. Constraint (15)) is satisfiable, then  $\sum_{r \in [1, q]} \Delta(s_r) \geq \tau$ . In other words, our approximation scheme will never miss the violation of any temporal constraint.*

However, it is worthwhile to mention that our approximation scheme may generate *false positives*. In particular,  $\sum_{r \in [1, q]} \Delta(s_r)$  might over-approximate the maximum value of  $\sum_{i, j} \delta_i^j$ . This is due to the reason that interleaving patterns, which lead to the maximum delay for individual cache sets, may not be feasible together. In our evaluation, we empirically evaluate the amount of pessimism in our approximation scheme.

## 4 Extension

*Applications with shared variables* Recall that we monitor the performance of each core *in isolation* and replace the delay to access the shared cache via a symbolic variable. In particular, our framework handles interferences in the shared resources, but, not in the shared variables. As a result, we do not catch the scenario when the program control-flow changes due to updates to shared variables. However, many embedded applications are designed by a number of independent components and the communication occurs in terms of reading sensor inputs or writing to output ports. In our evaluation, we show a real-life robot controller which operates via two independent tasks – balance and navigation. Moreover, shared memory-space across cores often bypass caches, to avoid power consumption due to coherence traffic [18]. If accessing the shared memory-space bypasses cache, our framework can be easily extended for general applications with shared variables. In order to accomplish this, we need to generate additional constraints, which encode the program control-flow observed during a failure run (*i.e.* an execution scenario violating certain temporal constraints). This can be achieved in an exactly same fashion as shown in [19].

It is slightly more complex when accessing the shared memory-space goes through caches. In particular, we need to add constraints that capture cache misses due to data coherency and false sharing [2]. This can be accomplished by correlating writes and reads to the same memory block. Besides, we need to distinguish the *first-ever shared-cache miss* for a memory block. Without data sharing, such cache misses can be detected during the inspection of each core in isolation. However, with data sharing, we need to detect first-ever shared-cache misses using the following constraints, for any replacement policy:

$$\bigwedge_{i,j} \left( \bigwedge_{p,q: \sigma_p^q = \sigma_i^j} (\mathcal{O}_p^q > \mathcal{O}_i^j) \Rightarrow (\delta_i^j = \text{MISS}) \right) \quad (21)$$

Constraint (21) encodes the scenario of  $i^j$  being the first shared-cache access to request memory block  $\sigma_i^j$ . This, in turn, leads to a shared-cache miss. We are currently extending MESS to handle data sharing and cache coherency.

*Performance debugging for a class of inputs* With minor changes, our framework can be extended for performance debugging on a class of inputs. The key to such extension is to collect *path conditions* [16], while monitoring the performance of each core *in isolation*. For each core, such a *path condition* captures the set of all inputs which lead to the respective execution scenario. However, depending on the value of input  $x$ , the statement  $a[x]$  might access different memory blocks, for the same path condition. Therefore, we need to generate constraints for each such memory block, satisfying the respective path constraint. Let us assume that array  $a$  might access memory block  $m_1$  if  $0 \leq x \leq 2$  and it accesses memory block  $m_2$  if  $2 < x \leq 5$ . Subsequently, to formulate cache conflicts generated by memory blocks (*i.e.* Constraints (5)-(6) for LRU policy and Constraints (11)-(12) for FIFO policy)  $m_1$  and  $m_2$ , we additionally constrain via conditions  $(0 \leq x \leq 2)$  and  $(2 < x \leq 5)$ , respectively. For instance, we modify  $\Theta_1^{lru}(i, j)$

to  $\Theta_1^{lru}(i, j) \wedge (0 \leq x \leq 2)$  for memory block  $m_1$  and to  $\Theta_1^{lru}(i, j) \wedge (2 < x \leq 5)$  for memory block  $m_2$ . In future, we aim to build such extension to instantiate performance debugging on a set of inputs, which are captured symbolically by path conditions.

## 5 Evaluation

We have implemented MESS using `simplescalar` [6] and `Z3` constraint solver [5]. In our evaluation, we configure a multi-core system with dual-core processor, where each core has a private level-one cache and all the cores share a level-two cache. This is a typical design in many embedded systems, such as devices using Exynos 5250 [4], which, in turn, contains a dual-core, ARM Cortex-A15 [1] chip. We configure 1 KB level-one caches with associativity 2 and 2 KB level-two cache with associativity 4. All caches have a line size of 32 bytes. Cache sizes are chosen in a fashion such that we obtain enough accesses to the shared cache and therefore, generate a reasonable number of constraints in our framework (see the Appendix for experiments with different cache configurations). To evaluate our framework, we have chosen medium to large size programs from [17], which are generally used to validate timing analyzers. We have also used a robot controller from [3], which contains two tasks — `balance` (to help the robot to keep it in upright position) and `navigation` (to drive the robot through rough terrain). These two tasks are assigned to different cores in our configured dual-core system.

**Experimental setup** For our evaluation with programs from [17], we run `jfdctint` on one core and choose different programs to run on the other core. We use such a setup in order to check the influence of the same inter-core cache conflicts on different programs. For the robot controller, we run `balance` and `navigation` on two different cores. The first two columns in Table 1 list the set of programs and the respective size of source code. We monitor the execution on each core by instrumenting memory accesses in `Simplescalar`. At the end of the execution, we generate a summary of memory performance for each core, which, in turn are used to generate constraints. The generated constraints are solved via `Z3`. All evaluations have been performed on an Intel I7 machine, having 8 GB of RAM and running ubuntu 14.04 operating systems.

**Basic results** Table 1 outlines the basic evaluation of our framework. We set the shared-cache miss-penalty (hit-latency) to be 100 (1) cycles. Recall that we aim to check the validity of temporal constraints  $\sum_{i,j} \delta_i^j < \tau$ . We generate a number of temporal constraints by varying  $\tau$  from 200 to 3100 cycles, at a step of 100 cycles and for each such temporal constraint, we invoke our framework. Note that  $\tau$  captures all possibilities between two to thirty one shared-cache misses. Besides, in  $\sum_{i,j} \delta_i^j$ , we only consider shared-cache accesses  $i^j$ , whose latency were unknown during the investigation of each core in isolation (*cf.* Column 4 in Table 1). Therefore, any shared-cache access  $i^j$ , which incurs the first-ever cache miss of the respective memory block  $\sigma_i^j$ , is not included in  $\sum_{i,j} \delta_i^j$ . In Table 1, we report the maximum and geometric mean over the time to check all temporal constraints. For several cases, this maximum time was recorded for a *valid* temporal constraint, meaning that the solver failed to find a violation. We can observe

| Program     | Total lines of C code | Shared-cache repl. policy | #shared-cache access | Size of constraints | #violations | Time to generate constraints (secs) | Solver time (secs) Max. / Geo. Mean |
|-------------|-----------------------|---------------------------|----------------------|---------------------|-------------|-------------------------------------|-------------------------------------|
| cnt         | 642                   | LRU                       | 432                  | 2111                | 22          | 1.17                                | 25.01 / 1.84                        |
| +jfdctint   |                       | FIFO                      | 432                  | 6586                | 22          | 9.52                                | 161.83 / 15.73                      |
| expint      | 532                   | LRU                       | 433                  | 2166                | 23          | 1.22                                | 10.84 / 2.16                        |
| +jfdctint   |                       | FIFO                      | 433                  | 6643                | 23          | 9.62                                | 576.56 / 20.02                      |
| qurt        | 541                   | LRU                       | 448                  | 2817                | 30          | 1.88                                | 24.81 / 3.16                        |
| +jfdctint   |                       | FIFO                      | 448                  | 7272                | 30          | 9.38                                | 31.77 / 11.59                       |
| matmult     | 538                   | LRU                       | 436                  | 2283                | 28          | 1.31                                | 244.39 / 1.91                       |
| +jfdctint   |                       | FIFO                      | 436                  | 6758                | 28          | 9.69                                | 15495.83 / 12.82                    |
| fdct        | 614                   | LRU                       | 479                  | 3943                | 30          | 2.99                                | 17.49 / 5.01                        |
| +jfdctint   |                       | FIFO                      | 479                  | 8418                | 30          | 11.85                               | 44.31 / 21.44                       |
| nsichneu    | 4628                  | LRU                       | 1679                 | 40087               | 30          | 49.2                                | 17120.46 / 7904.08                  |
| +jfdctint   |                       | FIFO                      | 1679                 | 44562               | 30          | 15.35                               | 27534.20 / 15174.8                  |
| balance     | 2098                  | LRU                       | 772                  | 3881                | 30          | 0.23                                | 155.17 / 63.94                      |
| +navigation |                       | FIFO                      | 773                  | 6770                | 30          | 0.56                                | 389.68 / 184.32                     |

**Table 1.** Evaluation of our baseline framework: “lines of C code” considers the sum of source code of two programs running on two different cores, “#violations” captures the number of violations within the set of 30 temporal constraints  $\{\sum_{i,j} \delta_i^j < 200, \dots, \sum_{i,j} \delta_i^j < 3100\}$ .

that, for many scenarios, the solver returns a solution in reasonable time. However, with large number of constraints, the solver takes long time to find a solution. For instance, with program `nsichneu`, such a scenario happens due to its large size and a substantial number of accesses to the shared-cache. In general, finding a solution for FIFO policy takes longer time compared to LRU policy, due to a larger constraint-size.

| Program            | Replacement policy of the shared-cache | Max. #constraints |        | Solver time (in seconds) |        | Max. delay ( $\max \sum_{i,j} \delta_i^j$ ) (in CPU cycles) |        |
|--------------------|--|-------------------|--------|--------------------------|--------|---|--------|
|                    |  | baseline          | approx | baseline                 | approx | baseline  | approx |
| cnt+jfdctint       | LRU                                    | 2111              | 154    | 23.58                    | 4.39   | 2394  | 3285   |
|                    | FIFO                                   | 6586              | 513    | 116.49                   | 14.35  | $2300 < X < 2400$   | 3285   |
| expint+jfdctint    | LRU                                    | 2166              | 207    | 10.84                    | 4.77   | 2494  | 3385   |
|                    | FIFO                                   | 6643              | 526    | 409.39                   | 14.58  | $2400 < X < 2500$   | 3385   |
| qurt+jfdctint      | LRU                                    | 2817              | 305    | 565.91                   | 9.2    | 3884  | 6161   |
|                    | FIFO                                   | 7272              | 631    | <i>TO</i>                | 29.03  | $\geq 3900$   | 6061   |
| matmult+jfdctint   | LRU                                    | 2283              | 154    | 244.39                   | 5.23   | 2988  | 4473   |
|                    | FIFO                                   | 6758              | 513    | 15495.83                 | 15.98  | $2900 < X < 3000$   | 4473   |
| fdct+jfdctint      | LRU                                    | 3943              | 304    | <i>TO</i>                | 22.31  | $\geq 6200$   | 10116  |
|                    | FIFO                                   | 8418              | 599    | <i>TO</i>                | 66.4   | $\geq 6200$   | 10116  |
| nsichneu+jfdctint  | LRU                                    | 40087             | 2862   | <i>TO</i>                | 764.56 | $\geq 10000$  | 31500  |
|                    | FIFO                                   | 44562             | 3137   | <i>TO</i>                | 926.45 | $\geq 10000$  | 31500  |
| balance+navigation | LRU                                    | 3881              | 442    | 93.32                    | 12.81  | $12800 < X < 12900$   | 13200  |
|                    | FIFO                                   | 6770              | 818    | 182.68                   | 25.08  | 12200   | 12200  |

**Table 2.** Efficiency and precision of our approximation. *TO* denotes timeout ( $> 5$  hours). “Max. #constraints” capture the maximum number of constraints solved by Z3 over all invocations.

**Evaluation of the approximate solution** Table 2 compares our approximation and the baseline framework. As clearly observed, our approximation dramatically reduces the debugging time, compared to the baseline framework. This is due to the partitioning of constraints with respect to different cache sets. Such constraint partitioning drastically reduces the number of constraints to be solved together, leading to a substantial reduction of pressure to Z3. As our approximation may generate false positives, we

also compare the precision of our approximation compared to the baseline framework. In order to do this, we compare the maximum delay computed by our approximation with the maximum delay computed by the baseline framework. This maximum delay captures the sum of all delays to access the shared-cache. For our baseline framework, obtaining such maximum delay may incur large overhead (we used `symba` [21] to compute the maximum delay). In such cases, we use the time taken by the solver to validate a temporal constraint  $\sum_{i,j} \delta_i^j < \tau$ . This means that the maximum delay cannot exceed  $\tau - 1$ . For instance, in Table 2,  $2300 \leq X < 2400$  indicates that the solver found a solution for  $\sum_{i,j} \delta_i^j \geq 2300$ , but not for  $\sum_{i,j} \delta_i^j \geq 2400$ . The respective debugging-time captures the time taken by the solver for  $\sum_{i,j} \delta_i^j \geq 2400$ . Finally, we use a timeout of five hours for the solver. For instance, the timeout event happens for the program `fdct`. From Table 2, we also observe that the precision of our approximation scheme is reasonable, in the context of validating embedded software. Finally, we note that with the current state-of-the-art solutions (e.g. using [21]), discovering the exact worst-case ordering among memory accesses (in terms of performance), is not very efficient.

**Notes on scalability** We have implemented a *proof-of-concept* of MESS. We have also shown an approximation, which dramatically improves the solver performance, with a reasonable loss of precision. We believe that several optimizations are still possible. In particular, as shown in [19], other optimizations for parallel constraint-solving is feasible. We are exploring such techniques to further improve the efficiency of MESS.

## 6 Related work

Testing and debugging of multi-threaded applications has been an active topic of research for the last few years [24,12,19,29,26,23,30,7]. For multi-threaded applications, the key challenge is to detect thread scheduling patterns that may lead to software functionality bugs. On the contrary, we systematically detect the order of memory accesses that may lead to performance bugs on multi-core platforms. Our approach concentrates on resource sharing in parallel architectures, rather than data sharing in parallel applications. However, to consider shared data in our framework, an approach similar to [19] can be integrated easily into our constraint system.

Modeling shared-cache performance has been an active topic of research in the past decade [25,27,31,15,14,28]. These works aim to provide an average confidence on the performance of shared caches. In contrast, we aim to bring the power of formal methods to provide strong guarantees on the presence or absence of performance bugs due to shared-caches. We believe that such guarantees are crucial for time-critical code fragments and our work is a preliminary step to establish such guarantees.

A different line of work [9] aims to synthesize correct and optimal concurrent programs from their specification and a performance model. The goal of the proposed approach is to synthesize partial programs, which are not only correct, but also optimal, with respect to the given performance model. Our work in this paper is orthogonal to efforts in program synthesis, such as the approach taken in [9]. Instead of generating correct and optimal programs from their specification, we aim to discover performance bugs in the original implementation of software. The presence of performance bugs in

implementation is, in general, inevitable. Therefore, our approach is inclined towards performance debugging. However, we believe that our approach can be complimentary to synthesis approaches. For instance, a performance debugger can highlight the code fragments that could be refactored optimally.

Recent works on performance testing [8,20] have either targeted sequential applications or the targeted bugs can be discovered by exploring exactly one interleaving, namely *canonical schedule* [20]. However, the order of memory accesses may dramatically influence the overall performance. Therefore, our primary goal is to highlight the order of shared-cache accesses that expose a performance bug (*i.e.* violates a given temporal constraint). Our previous work [11] aims to generate performance-stressing execution in embedded GPUs and it is not suitable for reproducing or debugging performance bugs. Besides, in this paper, we provide strong guarantees on the absence of performance bugs, when a given temporal constraint is not invalidated by the solver.

Works on worst case execution time (WCET) analysis have recently made progress in the context of multi-core platforms [10]. As the name suggests, WCET captures the maximum execution time of an application over all possible inputs and interleaving patterns. In this paper, our goal is orthogonal to approaches taken in WCET analysis. In particular, for a given program input, we aim to discover the interleaving pattern that causes the violation of temporal constraints. Therefore, our work has a significant testing and debugging flavour compared to the approaches proposed via WCET analysis.

In summary, previous works on automated debugging have mostly concentrated on functionality bugs or performance bugs on single-core systems. In this paper, we propose a systematic debugging approach that highlights performance bugs on multi-core systems, with a specific focus on shared caches.

## 7 Conclusion

In this paper, we have proposed MESS, a constraint-based framework to debug memory performance in multi-core systems. MESS systematically finds the interleaving pattern that causes the violation of temporal constraints. An appealing feature of our framework is its ability to provide guarantees on the absence of performance bugs, such as the validity of temporal constraints, for a given input. We have also integrated an approximation scheme, which, with a reasonable loss of precision, improves the debugging time by several magnitudes. In general, this opens up several opportunities to improve the debugging time enforced by MESS. Our evaluation with several embedded software and also with a real-life robot controller shows the effectiveness of our approach. Finally, since the performance of constraint solvers is continuously improving, we believe that MESS proposes a promising approach for performance debugging on multi-core systems. In future, we aim to build on our approach to consider shared data and other crucial shared resources in multi-core systems, such as shared buses. We also aim to use MESS to automatically synthesize fixes of performance bugs. One possible approach would be to synthesize barriers. The primary purpose of such barriers will be to satisfy a given temporal constraint, via restricting certain interleaving patterns.

**Acknowledgement** We thank the anonymous reviewers for their insightful comments and feedback. This work is partially supported by the Swedish National Graduate School on Computer Science (CUGS). This support is gratefully acknowledged.

## References

1. ARM Cortex-A5 Processor. <http://www.arm.com/products/processors/cortex-a/cortex-a5.php>.
2. Avoiding and identifying false sharing among threads. <https://software.intel.com/en-us/articles/avoiding-and-identifying-false-sharing-among-threads>.
3. Ballybot balancing robots. <http://robotics.ee.uwa.edu.au/eyebot/doc/robots/ballybot.html>.
4. Samsung Exynos Processor. [http://www.samsung.com/global/business/semiconductor/file/product/Exynos\\_5\\_Dual\\_User\\_Manual\\_Public\\_REV100-0.pdf](http://www.samsung.com/global/business/semiconductor/file/product/Exynos_5_Dual_User_Manual_Public_REV100-0.pdf).
5. Z3 Constraint Solver. <http://z3.codeplex.com/>.
6. Todd Austin, Eric Larson, and Dan Ernst. SimpleScalar: An infrastructure for computer system modeling. *Computer*, 35(2), 2002.
7. Thomas Ball, Sebastian Burckhardt, Katherine E. Coons, Madanlal Musuvathi, and Shaz Qadeer. Preemption sealing for efficient concurrency testing. In *TACAS*, 2010.
8. Abhijeet Banerjee, Sudipta Chattopadhyay, and Abhik Roychoudhury. Static analysis driven cache performance testing. In *RTSS*, 2013.
9. Pavol Cerný, Krishnendu Chatterjee, Thomas A. Henzinger, Arjun Radhakrishna, and Rohit Singh. Quantitative synthesis for concurrent programs. In *CAV*, 2011.
10. Sudipta Chattopadhyay, Lee Kee Chong, Abhik Roychoudhury, Timon Kelter, Peter Marwedel, and Heiko Falk. A unified WCET analysis framework for multi-core platforms. *TECS*, 13(4s), 2014.
11. Sudipta Chattopadhyay, Petru Eles, and Zebo Peng. Automated software testing of memory performance in embedded gpus. In *EMSOFT*, 2014.
12. Dongdong Deng, Wei Zhang, and Shan Lu. Efficient concurrency-bug detection across inputs. In *OOPSLA*, 2013.
13. Joseph Devietti, Brandon Lucia, Luis Ceze, and Mark Oskin. DMP: deterministic shared memory multiprocessing. In *ASPLOS*, 2009.
14. Chen Ding, Xiaoya Xiang, Bin Bao, Hao Luo, Ying-Wei Luo, and Xiao-lin Wang. Performance metrics and models for shared cache. *J. Comput. Sci. Technol.*, 29(4), 2014.
15. David Eklov, David Black-Schaffer, and Erik Hagersten. Fast modeling of shared caches in multicore systems. In *HiPEAC*, 2011.
16. Patrice Godefroid, Nils Klarlund, and Koushik Sen. DART: directed automated random testing. In *PLDI*, 2005.
17. Jan Gustafsson, Adam Betts, Andreas Ermedahl, and Björn Lisper. The mälardalen WCET benchmarks: Past, present and future. In *WCET*, 2010.
18. Bryce Holton, Ke Bai, Aviral Shrivastava, and Harini Ramaprasad. Construction of GCCFG for inter-procedural optimizations in software managed manycore (SMM) architectures. In *CASES*, 2014.
19. Jeff Huang, Charles Zhang, and Julian Dolby. CLAP: recording local executions to reproduce concurrency failures. In *PLDI*, 2013.
20. Guodong Li, Peng Li, Geoffrey Sawaya, Ganesh Gopalakrishnan, Indradeep Ghosh, and Sreeranga P. Rajan. GKLEE: Concolic verification and test generation for GPUs. In *PPoPP*, 2012. [http://www.cs.utah.edu/formal\\_verification/GKLEE/](http://www.cs.utah.edu/formal_verification/GKLEE/).

21. Yi Li, Aws Albarghouthi, Zachary Kincaid, Arie Gurfinkel, and Marsha Chechik. Symbolic optimization with SMT solvers. In *POPL*, 2014.
22. Yun Liang and Tulika Mitra. Instruction cache locking using temporal reuse profile. In *DAC*, 2010.
23. Madanlal Musuvathi and Shaz Qadeer. Iterative context bounding for systematic testing of multithreaded programs. In *PLDI*, 2007.
24. Santosh Nagarakatte, Sebastian Burckhardt, Milo M. K. Martin, and Madanlal Musuvathi. Multicore acceleration of priority-based schedulers for concurrency bug detection. In *PLDI*, 2012.
25. Xiaoyue Pan and Bengt Jonsson. Modeling cache coherence misses on multicores. In *ISPASS*, 2014.
26. Malavika Samak and Murali Krishna Ramanathan. Trace driven dynamic deadlock detection and reproduction. In *PPoPP*, 2014.
27. Andreas Sandberg, David Black-Schaffer, and Erik Hagersten. Efficient techniques for predicting cache sharing and throughput. In *PACT*, 2012.
28. Andreas Sandberg, Andreas Sembrant, Erik Hagersten, and David Black-Schaffer. Modeling performance variation due to cache sharing. In *HPCA*, 2013.
29. Koushik Sen. Race directed random testing of concurrent programs. In *PLDI*, 2008.
30. Dasarath Weeratunge, Xiangyu Zhang, and Suresh Jagannathan. Analyzing multicore dumps to facilitate concurrency bug reproduction. In *ASPLOS*, 2010.
31. Chi Xu, Xi Chen, Robert P. Dick, and Zhuoqing Morley Mao. Cache contention and application performance prediction for multi-core systems. In *ISPASS*, 2010.

## Appendix

The content of this section is not part of the main paper and it can be skipped. We provide this additional content to include illustrated examples, proofs and additional experiments. Readers, who want more detail, may read it at their discretion.

### Examples of our constraint system

In the following, we show the constraints generated for the example in Figure 1(a).

#### Example: constraint system for LRU caches

Consider a shared-cache with associativity 2 and employing LRU replacement policy. For the sake of simplicity, we shall assume that all accesses in Figure 1(a) go through the shared cache. At the first phase, executions of `Program x` and `Program y` are recorded *in isolation*. In particular, for `Program x`, we record the following information:

- $\zeta_x^1 = \langle \phi, \phi \rangle$ ,  $\sigma_x^1 = m1$ , and  $age_x^1 = 3$  (cache miss),
- $\zeta_x^2 = \langle m1, \phi \rangle$ ,  $\sigma_x^2 = m2$  and  $age_x^2 = 3$  (cache miss),
- $\zeta_x^3 = \langle m2, m1 \rangle$ ,  $\sigma_x^3 = m1$  and  $age_x^3 = 2$ ,
- $\zeta_x^4 = \langle m1, m2 \rangle$ ,  $\sigma_x^4 = m1$  and  $age_x^4 = 1$ , and
- $\zeta_x^5 = \langle m1, m2 \rangle$ ,  $\sigma_x^5 = m2$  and  $age_x^5 = 2$ .

Let us now consider the formulation of constraints for  $x^4$ . The memory block accessed at  $x^4$  (i.e.  $\sigma_x^4$ ) is  $m1$ . The set of memory blocks conflicting with  $\sigma_x^4$ , or  $\mathcal{C}_x^4$  is  $\{m2, m1', m2'\}$ . Within  $\mathcal{C}_x^4$ ,  $\{m1', m2'\}$  may generate inter-core cache conflicts. Consider the memory block  $m1'$ .  $m1'$  can be accessed at  $y^1$  and  $y^3$ . Besides,  $x^1$  and  $x^3$  access memory block  $m1$  prior to the shared-cache access  $x^4$  and they might hide the cache conflict generated by memory block  $m1'$ . Therefore,  $\psi_{cft}^{lru}(y^1, x^4)$  and  $\psi_{ref}^{lru}(y^1, x^4)$  are formulated as follows.

$$\psi_{cft}^{lru}(y^1, x^4) \equiv \mathcal{O}_y^1 < \mathcal{O}_x^4 \quad (22)$$

$$\psi_{ref}^{lru}(y^1, x^4) \equiv \neg(\mathcal{O}_y^1 < \mathcal{O}_x^1 \wedge \mathcal{O}_x^1 < \mathcal{O}_x^4) \wedge \neg(\mathcal{O}_y^1 < \mathcal{O}_x^3 \wedge \mathcal{O}_x^3 < \mathcal{O}_x^4) \quad (23)$$

Since both  $y^1$  and  $y^3$  access memory block  $m1'$ ,  $\psi_{cft}^{lru}(y^3, x^4)$  and  $\psi_{ref}^{lru}(y^3, x^4)$  are formulated similarly as follows.

$$\psi_{cft}^{lru}(y^3, x^4) \equiv \mathcal{O}_y^3 < \mathcal{O}_x^4 \quad (24)$$

$$\psi_{ref}^{lru}(y^3, x^4) \equiv \neg(\mathcal{O}_y^3 < \mathcal{O}_x^1 \wedge \mathcal{O}_x^1 < \mathcal{O}_x^4) \wedge \neg(\mathcal{O}_y^3 < \mathcal{O}_x^3 \wedge \mathcal{O}_x^3 < \mathcal{O}_x^4) \quad (25)$$

Recall that the generation of cache conflict by  $m1'$  to  $x^4$  is captured via  $\Psi_x^4(m1')$ . Hence,  $\Psi_x^4(m1')$  is set via the following constraints.

$$\left( \psi_{cft}^{lru}(y^1, x^4) \wedge \psi_{ref}^{lru}(y^1, x^4) \right) \Rightarrow (\Psi_x^4(m1') = 1) \quad (26)$$

$$\left(\psi_{cft}^{lru}(y^3, x^4) \wedge \psi_{ref}^{lru}(y^3, x^4)\right) \Rightarrow (\Psi_x^4(m1') = 1) \quad (27)$$

The absence of cache conflict from  $m1'$  needs to consider both  $y^1$  and  $y^3$ , and it is, therefore, formalized as follows.

$$\left(\neg\psi_{cft}^{lru}(y^1, x^4) \vee \neg\psi_{ref}^{lru}(y^1, x^4)\right) \wedge \left(\neg\psi_{cft}^{lru}(y^3, x^4) \vee \neg\psi_{ref}^{lru}(y^3, x^4)\right) \Rightarrow (\Psi_x^4(m1') = 0) \quad (28)$$

In a similar fashion, we can generate constraints to formulate the value of  $\Psi_x^4(m2')$ . Finally, the delay faced by  $x^4$  (i.e.  $\delta_x^4$ ) can be determined using the following constraints.

$$(\Psi_x^4(m1') + \Psi_x^4(m2') \geq 2) \Rightarrow (\delta_x^4 = MISS) \quad (29)$$

$$(\Psi_x^4(m1') + \Psi_x^4(m2') \leq 1) \Rightarrow (\delta_x^4 = HIT) \quad (30)$$

Consider the interleaving pattern in Figure 1(e). Such an interleaving pattern is captured via the ordering  $\mathcal{O}_y^1 < \mathcal{O}_x^1 < \mathcal{O}_x^2 < \mathcal{O}_x^3 < \mathcal{O}_y^2 < \mathcal{O}_x^4 < \mathcal{O}_y^3 < \mathcal{O}_y^4 < \mathcal{O}_x^5$ . Since  $\mathcal{O}_y^1 < \mathcal{O}_x^3 < \mathcal{O}_x^4$ ,  $\psi_{ref}^{lru}(y^1, x^4)$  will be evaluated to *false* (cf. Constraint (23)). In a similar fashion, since  $\mathcal{O}_x^4 < \mathcal{O}_y^3$ ,  $\psi_{cft}^{lru}(y^3, x^4)$  will also be evaluated to *false* (cf. Constraint (24)). As a consequence, we shall observe that  $\Psi_x^4(m1')$  will be evaluated to *zero*. Hence,  $m1'$  does not generate any cache conflict at  $x^4$ . This will eventually result in a *shared-cache hit* at  $x^4$ , as was observed in our example (cf. Figure 1(f)).

### Example: constraint system for FIFO caches

Consider our example in Figure 1 along with a shared cache with associativity 2 and employing FIFO replacement policy. In the first stage, we record the information from each core *in isolation*. At the end of first stage, we obtain the following information for Program  $x$  and Program  $y$ :

- $\zeta_x^1 = \langle \phi, \phi \rangle$ ,  $\sigma_x^1 = m1$  and  $age_x^1 = 3$  (cache miss),
- $\zeta_x^2 = \langle m1, \phi \rangle$ ,  $\sigma_x^2 = m2$  and  $age_x^2 = 3$  (cache miss),
- $\zeta_x^3 = \langle m2, m1 \rangle$ ,  $\sigma_x^3 = m1$  and  $age_x^3 = 2$ ,
- $\zeta_x^4 = \langle m2, m1 \rangle$ ,  $\sigma_x^4 = m1$  and  $age_x^4 = 2$ ,
- $\zeta_x^5 = \langle m2, m1 \rangle$ ,  $\sigma_x^5 = m2$  and  $age_x^5 = 1$ ,
- $\zeta_y^1 = \langle \phi, \phi \rangle$ ,  $\sigma_y^1 = m1'$  and  $age_y^1 = 3$  (cache miss),
- $\zeta_y^2 = \langle m1', \phi \rangle$ ,  $\sigma_y^2 = m2'$  and  $age_y^2 = 3$  (cache miss),
- $\zeta_y^3 = \langle m2', m1' \rangle$ ,  $\sigma_y^3 = m1'$  and  $age_y^3 = 2$ , and
- $\zeta_y^4 = \langle m2', m1' \rangle$ ,  $\sigma_y^4 = m2'$  and  $age_y^4 = 1$ .

Note that, for FIFO policy, there exists a minor difference for shared-cache access  $x^4$ . In the following, we shall illustrate how the constraint generation for  $x^4$  differs from the same in LRU policy. The set of memory blocks conflicting to  $\sigma_x^4$  is  $\{m2, m1', m2'\}$ . Therefore,  $\mathcal{C}_x^4 = \{m2, m1', m2'\}$ . Within  $\mathcal{C}_x^4$ , consider the memory block  $m1'$ .  $m1'$  is accessed at  $y^1$  and  $y^3$ . Therefore, we formulate  $\psi_{cft}^{fif}(y^1, x^4)$  and  $\psi_{ref}^{fif}(y^1, x^4)$  as follows.

$$\psi_{cft}^{fif}(y^1, x^4) \equiv (\mathcal{O}_y^1 < \mathcal{O}_x^4) \wedge (\delta_y^1 = MISS) \quad (31)$$

$$\begin{aligned} \psi_{ref}^{fiffo}(y^1, x^4) &\equiv \neg(\mathcal{O}_y^1 < \mathcal{O}_x^1 \wedge \mathcal{O}_x^1 < \mathcal{O}_x^4 \wedge \delta_x^1 = MISS) \\ &\bigwedge \neg(\mathcal{O}_y^1 < \mathcal{O}_x^3 \wedge \mathcal{O}_x^3 < \mathcal{O}_x^4 \wedge \delta_x^3 = MISS) \end{aligned} \quad (32)$$

Note that we additionally guard Constraints (31)-(32) via conditions ( $\delta_y^1 = MISS$ ), ( $\delta_x^1 = MISS$ ) and ( $\delta_x^3 = MISS$ ), respectively. This is because, for FIFO policy, cache conflict is generated only in the presence of cache misses (*cf.* Section 3.2).

Since both  $y^1$  and  $y^3$  access memory block  $m1'$ ,  $\psi_{cft}^{fiffo}(y^3, x^4)$  and  $\psi_{ref}^{fiffo}(y^3, x^4)$  are formulated similarly as follows.

$$\psi_{cft}^{fiffo}(y^3, x^4) \equiv (\mathcal{O}_y^3 < \mathcal{O}_x^4) \wedge (\delta_y^3 = MISS) \quad (33)$$

$$\begin{aligned} \psi_{ref}^{fiffo}(y^3, x^4) &\equiv \neg(\mathcal{O}_y^3 < \mathcal{O}_x^1 \wedge \mathcal{O}_x^1 < \mathcal{O}_x^4 \wedge \delta_x^1 = MISS) \\ &\bigwedge \neg(\mathcal{O}_y^3 < \mathcal{O}_x^3 \wedge \mathcal{O}_x^3 < \mathcal{O}_x^4 \wedge \delta_x^3 = MISS) \end{aligned} \quad (34)$$

The generation of cache conflict by  $m1'$  to  $x^4$  is captured via  $\Psi_x^4(m1')$ . Hence,  $\Psi_x^4(m1')$  is set via the following constraints.

$$(\psi_{cft}^{fiffo}(y^1, x^4) \wedge \psi_{ref}^{fiffo}(y^1, x^4)) \Rightarrow (\Psi_x^4(m1') = 1) \quad (35)$$

$$(\psi_{cft}^{fiffo}(y^3, x^4) \wedge \psi_{ref}^{fiffo}(y^3, x^4)) \Rightarrow (\Psi_x^4(m1') = 1) \quad (36)$$

From Section 3.2, recall that for FIFO policy, we need to formulate constraints that take into account the cache-conflict generation both within cores and across cores. Therefore, we also generate constraints to set the value of  $\Psi_x^4(m2')$  and  $\Psi_x^4(m2)$  as follows.

$$(\psi_{cft}^{fiffo}(y^2, x^4) \wedge \psi_{ref}^{fiffo}(y^2, x^4)) \Rightarrow (\Psi_x^4(m2') = 1) \quad (37)$$

$$(\psi_{cft}^{fiffo}(y^4, x^4) \wedge \psi_{ref}^{fiffo}(y^4, x^4)) \Rightarrow (\Psi_x^4(m2') = 1) \quad (38)$$

$$(\psi_{cft}^{fiffo}(x^2, x^4) \wedge \psi_{ref}^{fiffo}(x^2, x^4)) \Rightarrow (\Psi_x^4(m2) = 1) \quad (39)$$

It is worthwhile to note that only shared-cache access  $x^2$  (accessing  $m2$ ) can generate cache conflict to  $x^4$  within a single core. This is because  $x^2$  appears before  $x^4$  in program order, but  $x^5$  does not appear before  $x^4$  in program order. In particular,  $\psi_{cft}^{fiffo}(x^2, x^4)$  and  $\psi_{ref}^{fiffo}(x^2, x^4)$  are determined in a similar fashion as follows.

$$\psi_{cft}^{fiffo}(x^2, x^4) \equiv (\mathcal{O}_x^2 < \mathcal{O}_x^4) \wedge (\delta_x^2 = MISS) \quad (40)$$

$$\begin{aligned} \psi_{ref}^{fiffo}(x^2, x^4) &\equiv \neg(\mathcal{O}_x^2 < \mathcal{O}_x^1 \wedge \mathcal{O}_x^1 < \mathcal{O}_x^4 \wedge \delta_x^1 = MISS) \\ &\bigwedge \neg(\mathcal{O}_x^2 < \mathcal{O}_x^3 \wedge \mathcal{O}_x^3 < \mathcal{O}_x^4 \wedge \delta_x^3 = MISS) \end{aligned} \quad (41)$$

The formulation to capture the absence of cache conflict is similar to LRU policy, as shown in the following.

$$\begin{aligned} &(\neg\psi_{cft}^{fiffo}(y^1, x^4) \vee \neg\psi_{ref}^{fiffo}(y^1, x^4)) \wedge (\neg\psi_{cft}^{fiffo}(y^3, x^4) \vee \neg\psi_{ref}^{fiffo}(y^3, x^4)) \\ &\Rightarrow (\Psi_x^4(m1') = 0) \end{aligned} \quad (42)$$

$$\begin{aligned} & \left( \neg\psi_{cft}^{ffo}(y^2, x^4) \vee \neg\psi_{ref}^{ffo}(y^2, x^4) \right) \wedge \left( \neg\psi_{cft}^{ffo}(y^4, x^4) \vee \neg\psi_{ref}^{ffo}(y^4, x^4) \right) \\ & \Rightarrow (\Psi_x^4(m2') = 0) \end{aligned} \quad (43)$$

$$\left( \neg\psi_{cft}^{ffo}(x^2, x^4) \vee \neg\psi_{ref}^{ffo}(x^2, x^4) \right) \Rightarrow (\Psi_x^4(m2) = 0) \quad (44)$$

Finally, the delay faced by  $x^4$  (i.e.  $\delta_x^4$ ) can be determined using the following constraints.

$$(\Psi_x^4(m1') + \Psi_x^4(m2') + \Psi_x^4(m2)) \geq 2 \Rightarrow (\delta_x^4 = MISS) \quad (45)$$

$$(\Psi_x^4(m1') + \Psi_x^4(m2') + \Psi_x^4(m2)) \leq 1 \Rightarrow (\delta_x^4 = HIT) \quad (46)$$

To understand the difference between LRU and FIFO policy, consider the interleaving pattern  $\mathcal{O}_y^1 < \mathcal{O}_x^1 < \mathcal{O}_x^2 < \mathcal{O}_x^3 < \mathcal{O}_y^2 < \mathcal{O}_x^4 < \mathcal{O}_y^3 < \mathcal{O}_y^4 < \mathcal{O}_x^5$  in Figure 1(e). Even though  $\mathcal{O}_x^2 < \mathcal{O}_x^3 < \mathcal{O}_x^4$ ,  $\psi_{ref}^{ffo}(x^2, x^4)$  will be evaluated to *true* (cf. Constraint (41)). This is due to the guard condition  $\delta_x^3 = MISS$ , which will be *false* for the respective interleaving pattern. As a result, unlike LRU policy, the cache conflict generated by  $m2$  will not be hidden by  $x^3$ . This will eventually result in a cache miss at  $x^4$ , as was observed in our example (cf. Figure 1(e)).

**Property 1** *Let us assume  $\{s_1, s_2, \dots, s_q\}$  are different sets in the shared cache. For a given temporal constraint  $\sum_{i,j} \delta_i^j < \tau$ , if our baseline constraint system  $\Phi$  (cf. Constraint (15)) is satisfiable, then  $\sum_{r \in [1,q]} \Delta(s_r) \geq \tau$ . In other words, our approximation scheme will never miss the violation of any temporal constraint.*

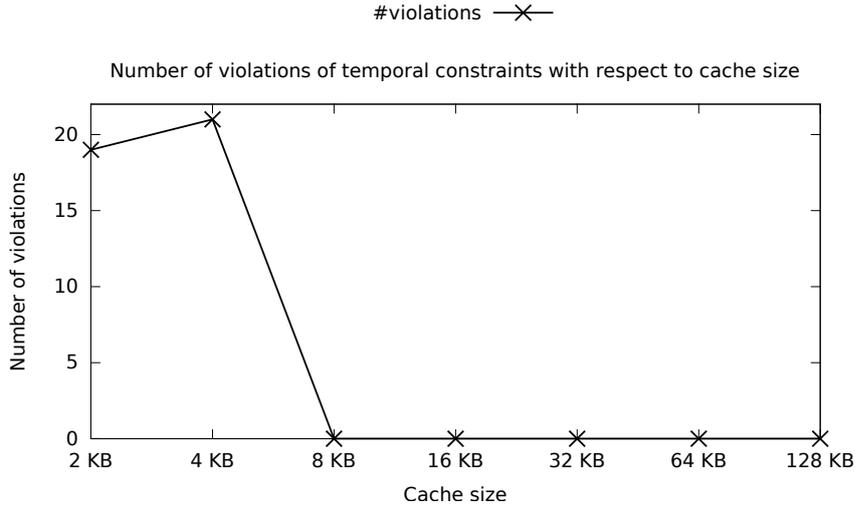
*Proof.* For an arbitrary cache set  $s_k$ , we shall show that  $\Gamma(s_k)$  (cf. Constraint (19)) includes all constraints that could influence the value  $\sum_{i,j: \pi(\sigma_i^j)=s_k} \delta_i^j$ . Therefore,  $\Delta(s_k)$  captures a sound over-approximation of the delay  $\sum_{i,j: \pi(\sigma_i^j)=s_k} \delta_i^j$ . Since this result holds for an arbitrary cache set  $s_k$ ,  $\sum_{k \in [1,q]} \Delta(s_k)$  is a sound over-approximation of the total delay to access the shared cache, i.e. of  $\sum_{i,j} \delta_i^j$ .

Let us now consider an arbitrary cache set  $s_k$ . Therefore, from Equation (20), we can observe that  $\Delta(s_k)$  includes the symbolic delay  $\delta_i^j$ , if and only if  $\sigma_i^j$  is mapped to cache set  $s_k$  (i.e.  $\pi(\sigma_i^j)=s_k$ ). Consider one such symbolic delay  $\delta_i^j$ . We first need to show that  $\Gamma(s_k)$  includes all constraints that might influence the value of  $\delta_i^j$ . The value of  $\delta_i^j$  is influenced by Constraints (5)-(8) for LRU policy and by Constraints (11)-(14) for FIFO policy. Note that Constraints (5)-(8) and Constraints (11)-(14) consider shared-cache access  $i^j$  and any shared-cache access  $\bar{i}^j$ , only if  $\sigma_i^j$  is conflicting to  $\sigma_{\bar{i}}^j$  (i.e.  $\sigma_{\bar{i}}^j \in \mathcal{C}_i^j$ ). Therefore,  $\Gamma(s_k)$  must include all constraints (from our baseline framework) that consider shared-cache access  $i^j$  and any shared-cache access  $\bar{i}^j$ , where  $\sigma_{\bar{i}}^j \in \mathcal{C}_i^j$ . From Constraints (16)-(18), we can observe that  $\Gamma(s_k)$  considers any shared-cache access  $i^j$ , where  $\sigma_i^j$  is mapped to cache set  $s_k$ . From the design principle of caches, we know that  $\sigma_{\bar{i}}^j \in \mathcal{C}_i^j$  if and only if  $\sigma_{\bar{i}}^j$  and  $\sigma_i^j$  are mapped to the same cache set or  $\pi(\sigma_{\bar{i}}^j) = \pi(\sigma_i^j)$ . Since  $\sigma_i^j$  is mapped to cache set  $s_k$ ,  $\sigma_{\bar{i}}^j$  is also mapped to cache set  $s_k$ . Hence, Constraints (17)-(18) include all constraints that might influence the value of  $\delta_i^j$ . Since we started with an arbitrary  $\delta_i^j$ , this result holds for any  $\delta_i^j$  included in  $\Delta(s_k)$ . Therefore, we conclude that  $\Gamma(s_k)$  includes all constraints that can influence the value  $\sum_{i,j: \pi(\sigma_i^j)=s_k} \delta_i^j$

(cf. Equation (20)). Since  $s_k$  is an arbitrary cache set,  $\Delta(s_k)$  captures a sound over-approximation of  $\sum_{i,j: \pi(\sigma_i^j)=s_k} \delta_i^j$ , for each cache set  $s_k$ . Therefore, our approximation scheme never misses the violation of any temporal constraint.

## Additional experiments

We evaluate our framework with multiple cache configurations. In order to do this, we choose the robot controller and instantiate our framework for multiple cache configurations. Figures 5-7 capture the evaluation for LRU cache replacement policy. From Table 2, we observe that  $12800 \leq \sum_{i,j} \delta_i^j < 12900$ , for LRU policy. We aim to check the influence of cache size on the threshold of  $\sum_{i,j} \delta_i^j$ . Therefore, we generate temporal constraints of the form  $\sum_{i,j} \delta_i^j < \tau$ , by varying  $\tau$  from 12000 to 13000 CPU cycles, at a step of 100 cycles. In a similar fashion, for FIFO policy, we generate the set of temporal constraints  $\{\sum_{i,j} \delta_i^j < 10000, \sum_{i,j} \delta_i^j < 10100, \dots, \sum_{i,j} \delta_i^j < 11900, \sum_{i,j} \delta_i^j < 12000\}$  (since  $\sum_{i,j} \delta_i^j = 12200$  for FIFO policy, see Table 2). Figures 8-10 capture our evaluation for FIFO policy.



**Fig. 5.** Number of violations of temporal constraints with respect to shared-cache size (LRU)

Note that, in  $\sum_{i,j} \delta_i^j$ , we only consider shared-cache accesses  $i^j$ , whose latency were unknown during the investigation of each core in isolation. Therefore, any shared-cache access  $i^j$ , which incurs a shared-cache miss during the investigation of each core in isolation, is not included in  $\sum_{i,j} \delta_i^j$ . For instance, in Figure 5 and in Figure 8, we observe that the number of violations increases when the cache size is changed from 2

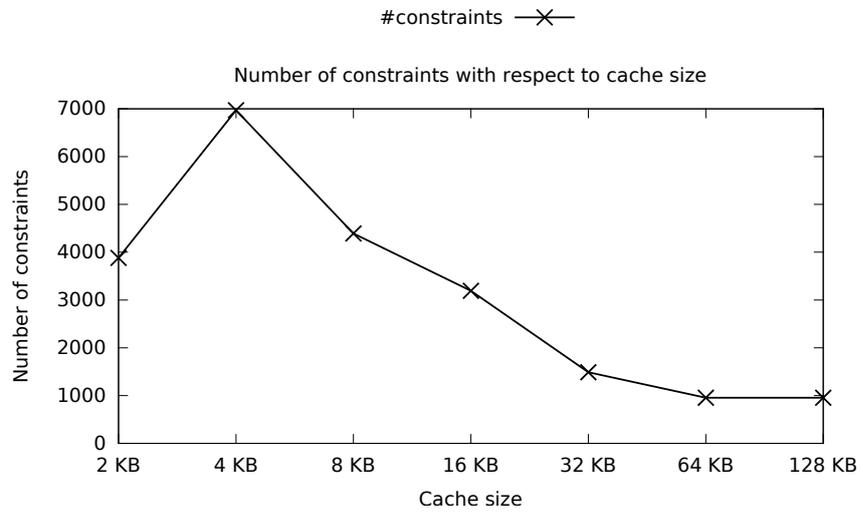


Fig. 6. Number of constraints with respect to shared-cache size (LRU)

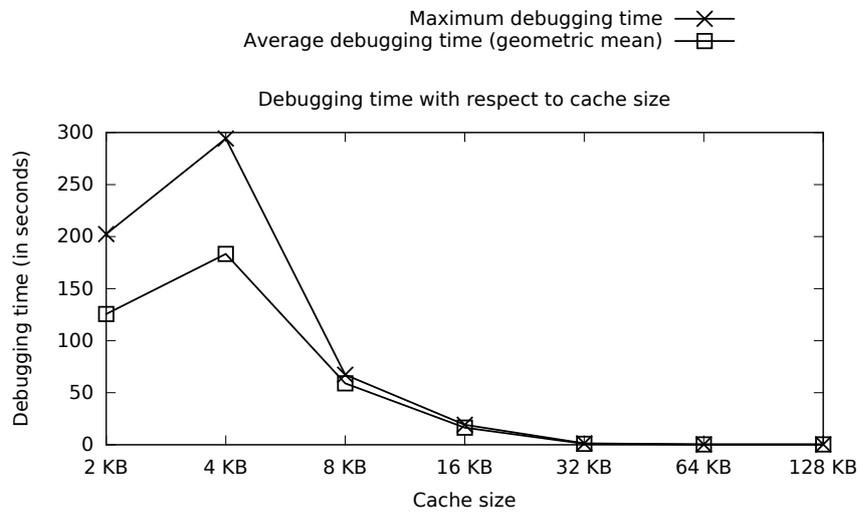


Fig. 7. Time taken by our framework with respect to shared-cache size (LRU)

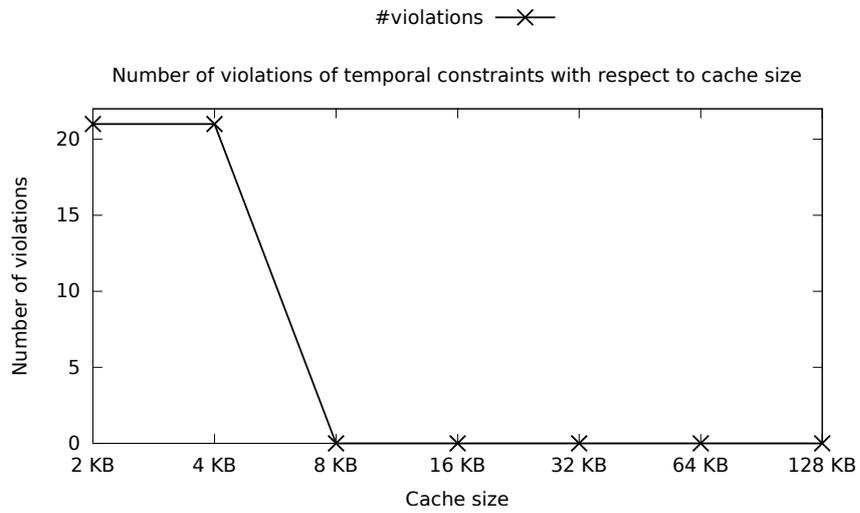


Fig. 8. Number of violations of temporal constraints with respect to shared-cache size (FIFO)

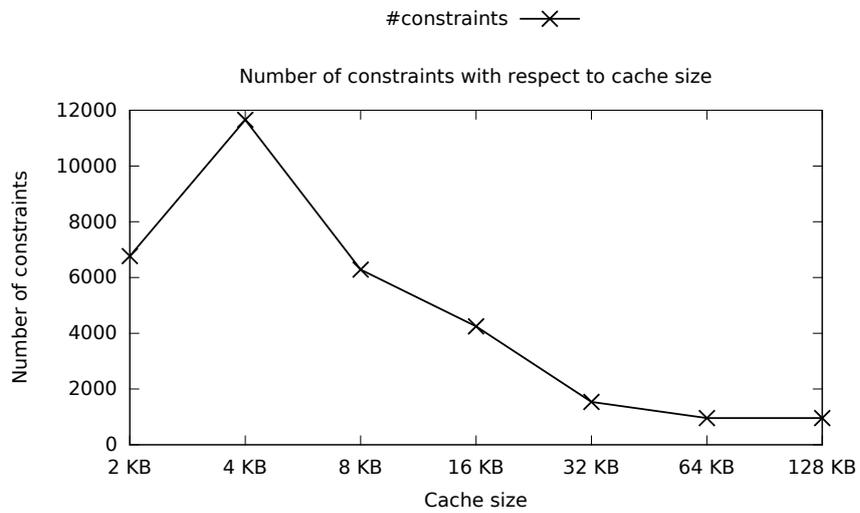
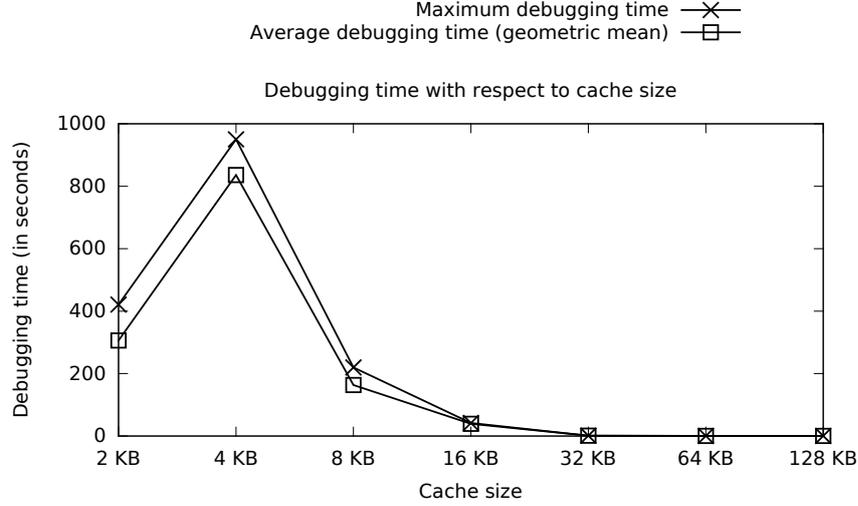


Fig. 9. Number of constraints with respect to shared-cache size (FIFO)



**Fig. 10.** Time taken by our framework with respect to shared-cache size (FIFO)

KB to 4 KB. This is because, an increased cache size may reduce the number of shared-cache misses within cores, leading to more accesses  $i^j$  to be included in  $\sum_{i,j} \delta_i^j$ . This may potentially lead to more violations of temporal constraints in our framework. It is, however, important to realize that increased cache size (with the same associativity and line-size) will always lead to reduced latency. As captured in Figure 5 and in Figure 8, the number of violations drops to zero when the cache size is changed to 8 KB or more. *Thus, we can conclude that shared-cache-misses, due to inter-core cache conflicts, do not reach the timing threshold captured by any temporal constraint.* This result also helps to choose an appropriate cache size for the respective application, given an input and timing threshold. For instance, given a timing threshold of 10000 CPU cycles, we can choose an 8 KB, FIFO cache for the robot controller.

Number of constraints may increase with cache size. This is because, the number of shared-cache accesses, which were not cache misses during the investigation of each core in isolation, may increase. As a result, we need to generate additional constraints for such shared-cache accesses. This behaviour was observed when the cache size was changed from 2 KB to 4 KB (*cf.* Figure 6 and Figure 9). However, the number of constraints decreases when the cache size is increased beyond 4 KB. This is because, an increased cache also reduces conflicting memory blocks to a particular shared-cache access. Recall that we need to generate constraints for each conflicting memory block (*cf.* Constraints (5)-(6) and Constraints (11)-(12)). As a result, increasing the cache size will also reduce the number of constraints to be considered in our framework.

Since the time taken by our framework directly depends on the number of generated constraints, we observe that the solver time increases when cache size is increased to 4 KB. Subsequently, the time taken by the solver reduces by several factors, and eventually becomes negligible with very large caches (*cf.* Figure 7 and Figure 10).

Finally, the number of constraints generated in FIFO policy is relatively larger than in LRU policy. As a result, the time taken by our framework for FIFO policy, is relatively longer compared to LRU policy, as observed from our evaluation.