

# 1 Impact of DM-LRU on WCET: a Static Analysis 2 Approach

3 Renato Mancuso

4 Boston University, USA

5 rmancuso@bu.edu

6 Heechul Yun

7 University of Kansas, USA

8 heechul.yun@ku.edu

9 Isabelle Puaut

10 University of Rennes 1/IRISA, France

11 isabelle.puaut@irisa.fr

## 12 — Abstract —

---

13 Cache memories in modern embedded processors are known to improve average memory access  
14 performance. Unfortunately, they are also known to represent a major source of unpredictability for  
15 hard real-time workload. One of the main limitations of typical caches is that content selection and  
16 replacement is entirely performed in hardware. As such, it is hard to control the cache behavior in  
17 software to favor caching of blocks that are known to have an impact on an application's worst-case  
18 execution time (WCET).

19 In this paper, we consider a cache replacement policy, namely DM-LRU, that allows system  
20 designers to prioritize caching of memory blocks that are known to have an important impact  
21 on an application's WCET. Considering a single-core, single-level cache hierarchy, we describe an  
22 abstract interpretation-based timing analysis for DM-LRU. We implement the proposed analysis in  
23 a self-contained toolkit and study its qualitative properties on a set of representative benchmarks.  
24 Apart from being useful to compute the WCET when DM-LRU or similar policies are used, the  
25 proposed analysis can allow designers to perform WCET impact-aware selection of content to be  
26 retained in cache.

27 **2012 ACM Subject Classification** Computer systems organization → Real-time systems; Theory of  
28 computation → Caching and paging algorithms

29 **Keywords and phrases** real-time, static cache analysis, abstract interpretation, LRU, deterministic  
30 memory, static cache locking, dynamic cache locking, cache profiling, WCET analysis

31 **Digital Object Identifier** 10.4230/LIPIcs.ECRTS.2019.20

## 32 **1** Introduction

33 Most modern embedded processors include cache(s) to improve average performance by  
34 reducing average memory access cost. However, a well-known downside of using caches is  
35 that it makes timing analysis difficult because software has little, if any, control over whether  
36 a certain memory block is in the cache or not, as it is determined by the hardware—the cache  
37 replacement policy and the state of the cache. This is problematic because analyzing precise  
38 and tight worst-case timing is necessary for real-time systems. While there are timing analysis  
39 techniques for well-known cache replacement policies [42], they cannot take advantage of  
40 programmer's insights (e.g., important data used in time-critical loops), potentially resulting  
41 in pessimistic timing.

42 On the other hand, a scratchpad memory is similar to a cache as it offers high-speed  
43 temporary storage for a processor, but the key difference is that it is entirely managed by  
44 software. For real-time systems, the fact that software, not hardware, has full control over  
45 its management is highly beneficial because accurate timing analysis is possible. However,  
46 the downside of scratchpad is that it is generally more difficult to use than cache due to  
47 its high programming complexity [3]. Alternatively, some cache designs support selective



48 cache locking, which enables programmers to lock certain cache-lines in the cache at a  
 49 fine-granularity (typically a cache line) [2, 7, 13]. A locked cache-line stays in the cache until  
 50 it is explicitly unlocked by the programmer, which guarantees predictable timing. However,  
 51 because the cache size is limited, the programmer must carefully select which cache-lines to be  
 52 locked [5, 40]. Dynamic cache-locking techniques [39, 48] can help alleviate the size limitation  
 53 problem of static cache-locking, but at the cost of increased complexity (for selecting locked  
 54 cache lines) and overhead (to change cache contents dynamically).

55 In this paper, we consider a new cache architecture, which can leverage programmers' high-  
 56 level insights on access frequency of memory blocks, and propose an abstract interpretation-  
 57 based static analysis method to reason on the worst-case execution time (WCET) of applica-  
 58 tions. Our approach is based on a new memory abstraction, called Deterministic Memory  
 59 (DM). Deterministic Memory enables classification of a program's address space into two  
 60 distinct memory types—DM and non-DM [10], where the DM type indicates predictability is  
 61 more important while the non-DM type indicates average performance is more important.  
 62 The DM abstraction allows effective and extensible software/hardware co-designs, some of  
 63 which are demonstrated in the context of providing efficient hardware isolation in multi-  
 64 core [10]. In this work, we instead focus on a single-core with a private cache, and study how  
 65 static guarantees on cache hits/misses can be derived for a DM-aware LRU cache replacement  
 66 policy, which we call DM-LRU.

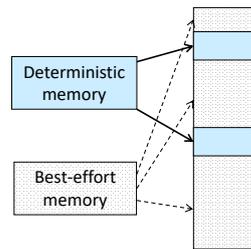
67 We first describe the DM-LRU cache replacement algorithm, which is a single-core  
 68 adaptation of the DM-aware cache initially proposed in [10]. Next, we generalize an abstract  
 69 interpretation-based analysis for LRU caches to reason on the worst-case behavior of DM-  
 70 LRU. We integrated DM-LRU support in Heptane [23], an academic static WCET analysis  
 71 tool, in order to evaluate the effectiveness of DM-LRU in lowering tasks' WCET. Our results  
 72 show that with DM-LRU WCET improvements up to 23.7% can be achieved, compared to  
 73 vanilla LRU. The WCET improvements are comparable to static and dynamic cache locking  
 74 techniques while significantly lowering programming complexity. Our contributions are as  
 75 follows:

- 76 ■ We extend LRU abstract interpretation-based analysis to perform static WCET timing  
 77 analysis for DM-LRU.
- 78 ■ We implement DM-LRU support in the Heptane static WCET analysis tool.
- 79 ■ We provide experimental evaluation results showing the WCET benefits and complexity  
 80 reduction of the DM-LRU based approach.
- 81 ■ We propose a WCET-driven heuristic approach to select content to be preferentially  
 82 cached using DM-LRU.

83 The remainder of the paper is organized as follows. Section 2 introduces necessary  
 84 background on caches and the deterministic memory abstraction. Next, the DM-LRU policy  
 85 is described in Section 3 and the proposed static timing analysis is described in Section 4.  
 86 A comprehensive example on how to apply the proposed analysis is presented in Section 5.  
 87 Comparison and differences with cache locking techniques are briefly highlighted in Section 6,  
 88 while the WCET of a set of representative benchmarks is evaluated in Section 7. Section 8  
 89 discuss related work and we conclude in Section 9.

## 90 **2** Background

91 In this section, we provide necessary background on memory abstractions, cache replacement  
 92 algorithms, and cache timing analysis.



■ **Figure 1** High-level application's memory view, where DM and BE memory coexist.

## 93 2.1 Deterministic Memory Abstraction

94 Traditionally, operating systems and hardware have provided a simple uniform memory  
 95 abstraction to applications. While the simple abstraction is convenient for programmability,  
 96 its downside is that programmer's insights on memory characteristics (e.g., time-criticality of  
 97 certain data structures) cannot be explicitly expressed to enable better resource management.

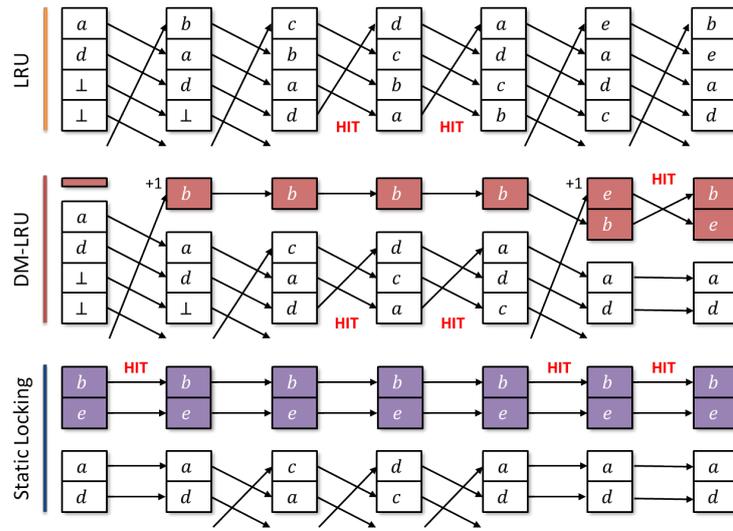
98 Recently, a new memory abstraction, called Deterministic Memory abstraction, was  
 99 proposed to explore the possibilities of more expressive memory abstractions [10]. In essence,  
 100 the abstraction allows a programmer to associate (tag) a single bit of information to each  
 101 memory block in the system, which classifies the memory block as either “deterministic  
 102 memory” (DM) or “best-effort memory” (BE). Figure 1 shows an example address space  
 103 of a task using both deterministic and best-effort memory. In [10], the memory tagging is  
 104 implemented at the page granularity, although more fine-granularity tagging is also possible  
 105 (e.g., [45]).

106 Once a task's memory blocks are tagged, the information can then be used by the operating  
 107 system and the hardware to apply different resource management policies depending on the  
 108 memory tag information. In [10], the DM abstraction is used to achieve hardware isolation  
 109 among the cores in multicore, focusing on effective isolation of shared cache and DRAM.

## 110 2.2 DM-LRU Cache Replacement Policy

111 In this paper, we consider a deterministic memory-aware *private* cache design and show how  
 112 such a design enables tighter static WCET cache timing analysis. We assume the cache  
 113 controller has a mean to distinguish whether a certain cache-line corresponds to deterministic  
 114 memory or best-effort one. This can be implemented as an additional bit in the auxiliary tag  
 115 store of each cache-line, as in [10], or as a set of separately located architectural hardware  
 116 range registers as in [27]. The cache implements an extended least recently used (LRU) cache  
 117 replacement algorithm, which defines two *eviction classes* using the DM/BE abstractions  
 118 and applies LRU-based replacement to DM lines and to BE lines separately. Allocation of a  
 119 DM line can cause eviction of a BE line, but the opposite is not allowed. Note that prior  
 120 work that implements a similar cache replacement policy exists [27]. In this paper, we call  
 121 the extended LRU as Deterministic Memory-aware Least Recently Used, or DM-LRU for  
 122 short. A more formal definition of DM-LRU is given in Section 3.

123 Figure 2 illustrates the difference between traditional LRU, DM-LRU, and static locking.  
 124 For simplicity, the example considers a single set of a 4-way set-associative cache. In the  
 125 first step, only  $a$  and  $d$  are cached, and 0 lines are allocated for DM blocks under DM-LRU.  
 126 Moreover, blocks  $b$  and  $e$  are set as DM blocks under DM-LRU, and pre-allocated in cache  
 127 in case of static locking. The figure tracks the evolution of the cache state for the same  
 128 access sequence  $b, c, d, a, e, b$ . A miss for a DM block triggers an increase of the number of  
 129 ways allocated for the DM class. This is depicted in step 2 (miss on  $b$ ) and 6 (miss on  $e$ ).



■ **Figure 2** Comparison between traditional LRU, DM-LRU, and a statically locked LRU cache over the same access pattern  $b, c, d, a, e, b$ , where  $b$  and  $e$  are DM memory blocks, or statically locked.

130 Traditional LRU simply ignores the DM/BE tag of the considered memory blocks. First,  
 131 note that DM-LRU results in fewer misses compared to LRU, as the DM marked memory  
 132 block  $b$  was not evicted by the best-effort memory accesses. Also note that while realizing  
 133 the same number of hits in the example compared to static locking, two important remarks  
 134 are required. First, the figure does not include the time spent to prefetch and lock the  $b$  and  
 135  $e$  blocks. Second, static locking causes additional misses for non-locked blocks compared to  
 136 DM-LRU. This exemplifies the on-demand nature of DM-LRU, which is able to retain in  
 137 cache blocks as they become needed during a task's execution. We discuss the analogies and  
 138 differences between DM-LRU and static/dynamic locking more extensively in Section 6.

139 Intuitively, it is thanks to the on-demand allocation and differential treatment of DM  
 140 memory blocks that DM-LRU enables tighter worst-case cache timing analysis, as we show  
 141 in the rest of the paper.

### 142 2.3 Cache Analysis via Abstract Interpretation

143 In this work, we extend abstract interpretation-based analysis to reason on the hit/miss  
 144 classification of memory accesses when a DM-LRU cache controller is implemented in  
 145 hardware. Analysis via abstract interpretation was originally proposed for LRU caches [11]  
 146 and better formalized and extended to FIFO and Pseudo-LRU in [41, 14]. An excellent  
 147 survey on the topic was proposed in [32]. We reuse the notation in [14, 32], while some details  
 148 are omitted due to space constraints. Since this work focuses on a DM-aware extension of  
 149 LRU, we introduce some of the background related to abstract interpretation-based LRU  
 150 analysis.

151 Imagine taking a snapshot of the cache state at a given point in time. In this case, one  
 152 could highlight the *state* of the cache in terms of: (i) which blocks are currently in cache,  
 153 and (ii) what is the *age* of each block. In LRU, the age of a block, say block  $a$ , captures  
 154 the number of memory accesses (to other blocks than  $a$ ) that were performed since the last  
 155 access to  $a$ . For instance, in the six steps in Figure 2, the LRU age for  $a$  is in the following  
 156 sequence: 0, 1, 2, 3, 0, 1, and 2. If  $a$  has an LRU age greater than or equal to the number of  
 157 ways (4 in our example), then  $a$  is not cached.

158 If the ages of all the cached blocks are known, the cache is in a **concrete state**. From a

159 concrete state, it is possible to produce a new concrete state that follows each new memory  
 160 access (state update), as shown in Figure 2. In a typical program, however, execution may  
 161 follow different paths. This means that at a given point in time, multiple concrete states are  
 162 possible, depending on the execution path taken by the program in its control-flow graph  
 163 (CFG).

164 Instead of keeping track of all the possible concrete states at any point of the CFG,  
 165 abstract interpretation keeps track of two main pieces of information: (i) the upper-bound  
 166 and (ii) the lower-bound on the age of any memory block among all the possible con-  
 167 crete states. Analysis on the age upper-bound and lower-bound is carried on separately.  
 168 The former is referred to as *must-analysis*, while the latter goes under the name of *may-*  
 169 *analysis*. A state that summarizes the upper-bound (resp., lower-bound) of each block  
 170 in a set of possible concrete states is called an **abstract state**. For instance, consider a  
 171 *must-analysis* abstract state of the form:  $\bar{q} = [\{\}, \{a, b\}, \{\}, \{d, e\}]$ . This corresponds to  
 172 all the concrete states where blocks  $a, b$  have age at most 1, and  $d, e$  at most 3. The full  
 173 *concretization* of  $\bar{q}$  is the set:  $\{[a, b, d, e], [b, a, d, e], [a, b, e, d], [b, a, e, d]\}$ . Similarly, consider  
 174 the *may-analysis* abstract state  $\underline{q} = [\{\}, \{\}, \{a, b\}, \{\}]$ . A concretization of  $\underline{q}$  is the set  
 175  $\{[\perp, \perp, \perp, \perp], [\perp, \perp, a, \perp], [\perp, \perp, \perp, a], [\perp, \perp, b, \perp], [\perp, \perp, \perp, b], [\perp, \perp, a, b], [\perp, \perp, b, a]\}$ , where  
 176  $\perp$  is a generic unknown block.

177 Given a *must-analysis* abstract state, it is possible to determine —i.e., *classify*— a  
 178 memory access as *always-hit* (H). These are accesses that result in hits regardless of the path  
 179 taken in the CFG. Similarly, given a *may-analysis* abstract state, it is possible to perform  
 180 classification of *always-miss* memory accesses. If neither classification applies, the block is  
 181 simply non-classified (NC). NC, often indicated as  $\top$ , represents the case in which some  
 182 execution paths lead to a miss while others lead to a hit for the same memory access.

183 Note that for architectures without timing anomalies [31, 20], *must-analysis* is sufficient  
 184 to safely compute the WCET of an application. In fact in this case NC accesses can be  
 185 simply treated as misses. We developed and implemented both *must-* and *may-*analysis for  
 186 DM-LRU, but we hereby focus in greater detail on *must-analysis*. Additional details about  
 187 *may-analysis* are provided in the appendix.

### 188 3 Cache Model and Terminology

189 In this section we discuss the cache model adopted to represent the behavior of DM-LRU, and  
 190 we introduce key concepts required to follow the proposed abstract interpretation analysis.

#### 191 3.1 DM-LRU Model

192 Algorithm 1 shows the full pseudo-code of the DM-LRU cache replacement algorithm. The  
 193 algorithm is defined for a generic  $A$ -way set-associative cache with  $S$  sets. The index of a set  
 194 is indicated with  $s \in \{0, \dots, S - 1\}$ . In the algorithm,  $DetMask_s$  denotes the bitmask of the  
 195 set  $s$ 's cache lines that contain deterministic memory. Consider a DM request ( $DM = 1$ ) that  
 196 resulted in a cache miss—see step 1 or 6 in Figure 2. The algorithm first tries to evict a BE  
 197 cache line, if such a line exists (Line 3-4). This also causes an additional bit to be asserted  
 198 in the  $DetMask_s$  bitmap. If no BE can be evicted (i.e., all lines are deterministic ones), it  
 199 chooses one of the deterministic lines (the older one in the LRU stack) as the victim (Line  
 200 6). On the other hand, consider the case where a BE memory block is requested ( $DM \neq 1$ ),  
 201 resulting in a miss—steps 1 and 2 in Figure 2. DM-LRU evicts one of the best-effort cache  
 202 lines, but not any of the deterministic cache lines (Line 9).

203 We assume a single-core, single-level set-associative cache. We indicate with  $A$  the  
 204 associativity of the cache. Since DM-LRU operates independently on each set, it is possible  
 205 to describe our analysis on a single set without loss of generality. Hereafter, we consider

```

Input :  $DetMask_s$  - deterministic ways of Set  $s$ 
Input :  $A$  - cache associativity
Output :  $victim$  - the victim way to be replaced or NULL if no replacement possible
1 if  $DM == 1$  then
2   if  $(\neg DetMask_s) \neq NULL$  then
3     // evict a best-effort line first
4      $victim = LRU(\neg DetMask_s)$ 
5      $DetMask_s | = 1 \ll victim$ 
6   else
7     // evict a deterministic line
8      $victim = LRU(DetMask_s)$ 
9   end
10 else
11   if  $(\neg DetMask_s) \neq NULL$  then
12     // evict a best-effort line
13      $victim = LRU(\neg DetMask)$ 
14   else
15     // no BE line can be allocated
16      $victim = NULL$ 
17   end
18 end
19 return  $victim$ 

```

**Algorithm 1:** Deterministic memory-aware cache line replacement algorithm.

206 a single cache set. At any point in time,  $D$  is the number of cache lines allocated to DM  
 207 memory blocks for the considered cache set.  $D$  is the number of bits set to “1” in the  
 208  $DetMask_s$  for the set under analysis. We indicate with  $B$  the number of lines that have not  
 209 been allocated for DM memory. It holds that  $D + B = A$ . Note that if  $D < A$ , and a DM  
 210 line that is currently not cached as a DM line is accessed, then the new DM line is allocated  
 211 and  $D$  is increased by one. This may trigger the eviction of the least recently used BE block,  
 212 as per Algorithm 1.

### 213 3.2 Terminology and notations

214 We indicate with  $\mathcal{B}$  the set of memory blocks that map to the cache set under analysis.  
 215 A generic memory block  $b^{CL} \in \mathcal{B}$  is comprised of an address  $b$  and an *eviction class*  
 216  $CL = \{DM, BE\}$ . The set of all the possible concrete states of a DM-LRU cache is denoted  
 217 as  $Q_{DM-LRU_A}$ , where each state  $q \in Q_{DM-LRU_A}$  is defined as follows:

$$218 \quad q := \{D, [b_0^{DM}, \dots, b_{D-1}^{DM}], [b_D^{CL}, \dots, b_{A-1}^{CL}]\}, \quad (1)$$

219 where  $D \in [0, A]$  and  $b_i^{CL} \in \mathcal{B}$ . Note that the first  $D$  cache lines are allocated as DM  
 220 cache lines, hence these are necessarily DM memory blocks. The remaining  $A - D$  blocks are  
 221 currently allocated BE memory blocks. Throughout this paper we will use the shorthand  
 222 notation  $b_i \in \mathcal{B}$  for blocks whose eviction class is obvious from context or unimportant. For  
 223 blocks allocated as BE, we assume BE class unless specified otherwise.

224 An important concept is the **age** of a memory block under DM-LRU, defined as follows.

225 ► **Definition 1 (DM-LRU Age).** *The age of a DM memory block  $a^{DM}$  is defined as the*  
 226 *number of distinct DM blocks accessed since the last access to  $a^{DM}$ ; the age of a BE memory*  
 227 *block  $b$  is set to the current value of  $D$  whenever  $b^{BE}$  is accessed. It is then defined as  $D + K$ ,*  
 228 *where  $K$  is the number of misses to DM blocks, or accesses to distinct BE blocks since the*  
 229 *last access to  $b^{BE}$ .*

230 Following Definition 1, the index of a given block  $b_i^{CL} \in q$  is also the age of the block  
 231 in DM-LRU. The age of a block  $b_i^{DM}$  allocated as DM can increase if: (1) a new DM line  
 232 is allocated (with age 0); or (2) a line  $b_j^{DM}$  already allocated as DM with age greater than  
 233  $b_i$  is accessed. Conversely, the age of a BE block  $b_i^{BE}$  can increase if: (1) a new DM line is  
 234 allocated (with age 0); (2) a new BE line is allocated (with age  $D$ ); or (3) a line  $b_j^{BE}$  already  
 235 in cache with age greater than  $b_i^{BE}$  is accessed.

236 Also note that Definition 1 remains consistent for the case in which a block  $b^{BE}$  is  
 237 accessed but cannot be allocated because all the sets have been reserved for DM lines. This  
 238 phenomenon goes under the name of *DM takeover*, and can be resolved by imposing a hard  
 239 cap on the maximum number of DM lines that can be allocated. The analysis for a DM-LRU  
 240 with an allocation cap is almost identical to an unrestricted DM-LRU, and only introduces  
 241 uninteresting subcases. For simplicity, we hereby focus on the analysis for unrestricted  
 242 DM-LRU. We demonstrate that preventing DM takeover is indeed necessary and beneficial  
 243 in Section 7.

## 244 4 DM-LRU Analysis

245 In this section we detail our abstract interpretation-based analysis [14, 32] for DM-LRU,  
 246 i.e. when the cache controller implements the policy defined in Algorithm 1. We discuss  
 247 *must*-analysis in detail. As previously mentioned, *may*-analysis is not strictly required for  
 248 architectures without timing anomalies. As such we only provide the intuition behind it and  
 249 defer the details to the appendix. We do not provide a persistence analysis for DM-LRU.  
 250 Persistence analysis is useful to determine if memory accesses inside loops can result in hits  
 251 after the first iteration. Instead, for our evaluations, we unroll the first iteration of each loop,  
 252 i.e., we perform *virtual unrolling*, *virtual inlining* (VIVU) [34, 32].

### 253 4.1 Must-analysis

254 *Must-analysis* is performed considering abstract cache states. In this case, *must-analysis*  
 255 keeps track of the upper bound on the number of allocated DM blocks indicated with  
 256  $D \in \{0, \dots, A\}$ , and the upper-bound on the DM-LRU age of each addressable memory block  
 257  $b \in \mathcal{B}$ . The abstract domain  $DMLru_A^{\square}$  is defined as:

$$258 \quad DMLru_A^{\square} := \{0, \dots, A\} \times \mathcal{B} \rightarrow \{0, \dots, A - 1, \infty\}. \quad (2)$$

259 Intuitively, the domain associates a *current eviction class* (DM or BE) and an age upper  
 260 bound ( $0, \dots, A$  or  $\infty$ ) to a memory block  $b \in \mathcal{B}$  mapping to the set under analysis. We use  
 261 the notation  $\bar{q}(b)$  to indicate the upper-bound on the age of  $b$  in  $\bar{q}$ . To represent a generic  
 262 abstract state  $\bar{q} \in DMLru_A^{\square}$  we use a compact notation that highlights the distinction  
 263 between DM and BE allocations. For instance, the notation

$$264 \quad \bar{q} = [\{\}, \{a, b\}], [\{c\}, \{d\}] \in DMLru_A^{\square} \quad (3)$$

265 denotes an abstract state  $\bar{q}$  where  $D \leq 2, B \geq 2, A = 4$ . Hence, blocks  $a$  and  $b$  have  
 266 upper-bound  $\bar{q}(a) = \bar{q}(b) = 1$  on their DM-LRU age. Similarly,  $c, d$  are BE blocks with  
 267  $\bar{q}(c) = 2$  and  $\bar{q}(d) = 3$ , respectively.

268 Given an abstract state  $\bar{q} \in DMLru_A^{\square}$ , the Boolean operator  $DM^{\square}(\bar{q}, b)$  returns *true*  
 269 only if the block  $b \in \mathcal{B}$  must exist as a DM-allocated block in  $\bar{q}$ . Formally

$$270 \quad DM^{\square}(\bar{q}, b^{CL}) := \begin{cases} true & \text{if } CL = DM \wedge \bar{q}(b) < \infty \\ false & \text{otherwise.} \end{cases} \quad (4)$$

271  
 272  
 273 For instance, considering  $\bar{q}$  defined as in Equation 3, we obtain  $DM^{\square}(\bar{q}, a) = true$ ,  
 274  $DM^{\square}(\bar{q}, d) = false$ , and so on. We use the simpler notation  $DM^{\square}(b)$  when the state is  
 275 implicit. The operator  $BE^{\square}(\bar{q}, b)$  is simply defined as  $BE^{\square}(\bar{q}, b) := \neg DM^{\square}(\bar{q}, b)$ . To prevent  
 276 additional clutter in our notation,  $DM^{\square}(\bar{q}, b^{DM})$  evaluates to *true* if and only if the DM  
 277 block  $b^{DM}$  *must* be allocated in cache in  $\bar{q}$ . As such, if the generic DM block  $b^{DM}$  has an  
 278 upper-bound on its DM-LRU age greater than  $A - 1$ , then  $BE^{\square}(\bar{q}, b^{DM}) = true$ .

279 An *abstract state transformer* for the  $DMLru_A^{\square}$  domain is an operator that takes in input  
 280 an abstract state  $\bar{q} \in DMLru_A^{\square}$  and any number of additional parameters, and returns in

281 output a transformed state  $\bar{q}' \in DM Lru_A^{\square}$ . We consider and define two abstract transformers  
 282 for  $DM Lru_A^{\square}$ : an update transformer  $U^{\square}(\bar{q}, a)$ , and a join transformer  $J^{\square}(\bar{q}, \bar{p})$ . We use the  
 283 operator  $\lambda b.$  to represent an age update operation carried on each  $b \in \mathcal{B}$  when considering a  
 284 transformation from state  $\bar{q}$  to  $\bar{q}'$ . This operator can be formally defined as:

$$285 \quad \lambda b. f(\bar{q}(b)) := \forall b \in \mathcal{B}, \bar{q}'(b) \leftarrow f(\bar{q}(b)) \quad (5)$$

## 286 Must-analysis Update

287 The update abstract transformer for the *must*-analysis  $U^{\square}(\bar{q}, a)$  is used to go from an initial  
 288 abstract state, to a new abstract state after a new memory access has been performed.  
 289  $U^{\square}(\bar{q}, a)$  takes in input an initial abstract state  $\bar{q}$  and a memory block  $a \in \mathcal{B}$ , and returns  
 290 the abstract state that results from accessing  $a$ . For ease of notation, we split the definition  
 291 of  $U^{\square}$  in two parts: the logic that corresponds to the update operation when a DM block  
 292  $a^{DM}$  is accessed, indicated with  $U_D^{\square}$ ; and the update transformation when a BE block  $a^{BE}$   
 293 is accessed, namely  $U_B^{\square}$ .  $U_D^{\square}$  is defined in Equation 6.

$$294 \quad U_D^{\square}(\bar{q}, a^{DM}) :=$$

$$295 \quad D' \leftarrow \begin{cases} D + 1 & \text{if } D < A \wedge BE^{\square}(a) \\ D & \text{if } D = A \vee DM^{\square}(a) \end{cases} \quad \begin{matrix} \text{(a.1)} \\ \text{(a.2)} \end{matrix}$$

$$\lambda b. \begin{cases} 0 & \text{if } b = a & \text{(b)} \\ \bar{q}(b) & \text{if } b \neq a \wedge \left\| \begin{array}{l} BE^{\square}(b) \wedge DM^{\square}(a) \\ DM^{\square}(b) \wedge \bar{q}(a) \leq \bar{q}(b) \\ BE^{\square}(b) \wedge BE^{\square}(a) \wedge \bar{q}(a) \leq \bar{q}(b) \end{array} \right. & \begin{matrix} \text{(c.1)} \\ \text{(c.2)} \\ \text{(c.3)} \end{matrix} \\ \bar{q}(b) + 1 & \text{if } b \neq a \wedge \bar{q}(a) > \bar{q}(b) \wedge \left\| \begin{array}{l} DM^{\square}(b) \wedge \bar{q}(b) < D' - 1 \\ BE^{\square}(b) \wedge BE^{\square}(a) \wedge \bar{q}(b) < A - 1 \end{array} \right. & \begin{matrix} \text{(d.1)} \\ \text{(d.2)} \end{matrix} \\ \infty & \text{if } b \neq a \wedge \bar{q}(a) > \bar{q}(b) \wedge \left\| \begin{array}{l} DM^{\square}(b) \wedge \bar{q}(b) \geq D' - 1 \\ BE^{\square}(b) \wedge BE^{\square}(a) \wedge \bar{q}(b) \geq A - 1 \end{array} \right. & \begin{matrix} \text{(e.1)} \\ \text{(e.2)} \end{matrix} \end{cases} \quad (6)$$

297 Here,  $D'$  ( $B'$ , resp.) is the new value of  $D$  ( $B$ , resp.) after the update. The conditions  
 298 following the  $\|$  operator are to be considered in logical “or” with each other.

299 The update abstract transformer  $U_B^{\square}$  for a best-effort memory access  $a$  can be defined as  
 300 follows:

$$303 \quad U_B^{\square}(\bar{q}, a^{BE}) :=$$

$$\lambda b. \begin{cases} D & \text{if } b = a \wedge D < A & \text{(a)} \\ \bar{q}(b) & \text{if } b \neq a \wedge \left\| \begin{array}{l} DM^{\square}(b) \\ BE^{\square}(b) \wedge \bar{q}(a) \leq \bar{q}(b) \end{array} \right. & \text{(b)} \\ \bar{q}(b) + 1 & \text{if } b \neq a \wedge BE^{\square}(b) \wedge \bar{q}(a) > \bar{q}(b) \wedge \bar{q}(b) < A - 1 & \text{(c)} \\ \infty & \text{if } \left\| \begin{array}{l} b = a \wedge D \geq A \\ b \neq a \wedge BE^{\square}(b) \wedge \bar{q}(a) > \bar{q}(b) \wedge \bar{q}(b) \geq A - 1 \end{array} \right. & \text{(d)} \end{cases} \quad (7)$$

304 To clarify the update operation, consider the abstract state  $\bar{q} = [\{\}, \{b, f\}], [\{c\}, \{d\}]$ ,  
 305 where  $D = 2$ . Assume that deterministic block  $a^{DM}$  is accessed, which has age upper-  
 306 bound  $\infty$  in  $\bar{q}$ , to obtain  $\bar{q}' = U^{\square}(\bar{q}, a) = U_D^{\square}(\bar{q}, a)$ . First, the value of  $D'$  is computed as

310  $D' = D + 1 = 3$  (a.1); next,  $b, f$  both satisfy the condition  $\bar{q}(a) > \bar{q}(b) = \bar{q}(f) = 1$ . Moreover,  
 311 we have that  $DM^{\square}(b) = DM^{\square}(f) = true$ , and that  $\bar{q}(b) = \bar{q}(f) = 1 < D' - 1 = 2$ . This  
 312 corresponds to condition (d.1) in Equation 6. Hence, the age of  $b, f$  in the resulting state is  
 313  $\bar{q}'(b) = \bar{q}'(f) = 2$ . Similarly, block  $c$  and  $d$  satisfy condition (d.2) and (e.2), respectively. The  
 314 resulting updated abstract state is:  $\bar{q}' = [\{a\}, \{\}, \{b, f\}], [\{c\}]$ .

315 An example for the abstract transformer  $U_{\bar{B}}^{\square}$  defined in Equation 7 is provided in Section 5.

316 ► **Theorem 2** (Correctness of *must*-analysis update). *Consider a generic abstract state*  
 317  $\bar{p} = U^{\square}(\bar{q}, a^{CL})$  *obtained from the must-analysis update state transformer when accessing*  
 318 *a generic block  $a^{CL}$  from an initial abstract state  $\bar{q}$ . Then for any block  $b \in \mathcal{B}$ ,  $\bar{p}(b)$  is an*  
 319 *upper-bound on the DM-LRU age of  $b$ .*

320 **Proof Sketch.** A proof can be constructed by considering two main sub-cases: (1) when  
 321  $CL = DM$  for the block being accessed; and (2) the case when  $CL = BE$ . Due to space  
 322 constraints, we provide an intuition for the former case, as the latter follows from the same  
 323 reasoning. When considering  $CL = DM$ , the new state  $\bar{p}$  is obtained as  $\bar{p} = U_{\bar{D}}^{\square}(\bar{q}, a^{DM})$ , as  
 324 per Equation 6.

325 First let us consider the rule on the update of  $D$ . If  $\bar{q}(a) = \infty$  then  $a$  is not necessarily  
 326 in cache and accessing  $a$  increases the upper-bound on the number of allocated DM blocks,  
 327 as long as the associativity  $A$  has not been exceeded, i.e.  $D < A$ . In this case, note that  
 328  $BE^{\square}(\bar{q}, a) = true$  and condition Equation 6 (a.1) applies.  $D$  does not change in any other  
 329 case (a.2). After the update, block  $a$  will have age upper-bound equal to 0 (b).

330 Next, consider all the blocks  $b \neq a$  that had age upper-bound of infinity in  $\bar{q}$  — i.e.  
 331  $\bar{q}(b) = \infty$ , and  $BE^{\square}(\bar{q}, b) = true$ . When  $a$  is accessed, their age upper-bound should not  
 332 change. If  $\bar{q}(a) = \infty$  then condition (c.3) applies. If  $\bar{q}(a) \neq \infty$  then  $DM^{\square}(\bar{q}, a) = true$  and  
 333 condition (c.1) applies.

334 Furthermore, consider all the blocks  $b^{DM}, b \neq a$  that must be allocated as DM blocks in  
 335  $\bar{q}$ , i.e. such that  $DM^{\square}(\bar{q}, b) = true$ . If  $\bar{q}(a) = \infty$ , the upper-bound on their DM-LRU age  
 336 will have to increase by 1 (d.1). If however the value of  $\bar{q}(b) + 1$  exceeds the updated value of  
 337  $D$ , namely  $D'$ , then the block may be evicted and the new upper-bound on its DM-LRU age  
 338  $\bar{p}(b) = \infty$  (e.1). The same cases apply when  $\bar{q}(a) < \infty$  and  $\bar{q}(a) > \bar{q}(b)$ .

339 On the other hand, if  $a$  has an age upper-bound that is same as or lower than  $b$ 's, i.e.  
 340  $\bar{q}(a) \leq \bar{q}(b)$ , then a concrete state where DM-LRU age of  $a$  is strictly larger than that of  $b$   
 341 cannot exist. As such, the upper-bound on the DM-LRU age of  $b$  will not change, as per  
 342 condition (c.2).

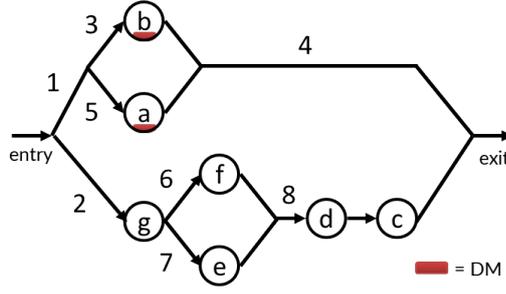
343 Lastly, consider all the blocks  $b^{BE}, b \neq a$  that must be allocated as BE blocks in  $\bar{q}$ , i.e.  
 344 such that  $BE^{\square}(\bar{q}, b) = true$  and  $\bar{q}(b) < \infty$ . The only case in which  $\bar{q}(a) > \bar{q}(b)$  is if  $\bar{q}(a) = \infty$ .  
 345 When  $a$  is accessed, the upper-bound on the age of  $b$  will have to increase by 1 (d.2), unless  
 346 by doing so the associativity  $A$  is exceeded. In the latter case,  $\bar{p}(b) = \infty$  (e.2). ◀

## 347 Must-analysis Join

348 The join abstract transformer  $J^{\square}(\bar{q}, \bar{p})$  is used to compute a new abstract state at the merging  
 349 point of two or more execution paths. There are strong similarities between the transformer  
 350 defined hereby and what used in traditional LRU *must*-analysis [14]. At a high level, the  
 351 joined state will consider as *must*-cached only those blocks in the intersection of the joining  
 352 states, each with the maximum age in any of the two states. For the new state,  $D$  is taken  
 353 as the maximum between the value of  $D$  in the joining states. Equation 8 formalizes the  
 354  $J^{\square}(\bar{q}, \bar{p})$  abstract transformer:

$$355 \quad J^{\square}(\bar{q}, \bar{p}) := D \leftarrow \max\{D_{\bar{q}}, D_{\bar{p}}\}, \lambda b. \max\{\bar{q}(b), \bar{p}(b)\} \quad (8)$$

356 If we were to join  $\bar{q} = [\{\}, \{b, f\}], [\{c\}, \{d\}]$  with  $\bar{q}' = [\{a\}, \{\}, \{b, f\}], [\{c\}]$ , the resulting  
 357 state would be  $\bar{q}'' = J^{\square}(\bar{q}, \bar{q}') = [\{\}, \{\}, \{b, f\}], [\{c\}]$ .



■ **Figure 3** Fragment of process CFG. At the end of the fragment, all the cache blocks in the figure may be cached.

358 ► **Theorem 3** (Correctness of *must*-analysis join). Consider an abstract state  $\bar{s} = J^\sqsubseteq(\bar{q}, \bar{p})$   
 359 obtained from the *must*-analysis join state transformer from two initial abstract states  $\bar{q}$  and  
 360  $\bar{p}$ . Then for any block  $b \in \mathcal{B}$ ,  $\bar{s}(b)$  is an upper-bound on the DM-LRU age of  $b$ .

361 **Proof Sketch.** A proof can simply follow from the definition of the  $J^\sqsubseteq$  operator in Equation 8.  
 362 By hypothesis  $\bar{q}$  and  $\bar{p}$  carry the upper-bound on the age of a generic block  $b$  along two disjoint  
 363 execution sub-paths. After the two sub-paths join, the maximum between  $\bar{q}(b)$  and  $\bar{p}(b)$  is a  
 364 safe upper-bound on the DM-LRU age of  $b$  in the resulting abstract state  $\bar{s}$ . Moreover, an  
 365 upper-bound on the number of allocated DM blocks in  $\bar{s}$  is the maximum between  $D_{\bar{q}}$  and  
 366  $D_{\bar{p}}$ . ◀

### 367 Must-analysis Classification

368 Every time an access is performed, it is possible to classify a memory access using a  
 369 classification function that will either return  $M$  for cache miss,  $H$  for cache hit, or  $\top$  in case  
 370 neither  $M$  nor  $H$  classification can be made given the current abstract state. In order to  
 371 classify memory accesses, for a given  $\bar{q}$  abstract state we define two helper sets  $\bar{\mathcal{D}}$  and  $\bar{\mathcal{B}}$   
 372 representing the deterministic and best-effort memory blocks that have finite upper bound  
 373 on their DM-LRU age:

$$\begin{aligned}
 374 \quad \bar{\mathcal{D}} &:= \{b^{CL} \in \mathcal{B} \mid CL = DM \wedge \bar{q}(b) < \infty\} \\
 375 \quad \bar{\mathcal{B}} &:= \{b^{CL} \in \mathcal{B} \mid CL = BE \wedge \bar{q}(b) < \infty\}
 \end{aligned} \tag{9}$$

377 The classification function of the *must* analysis is defined as:

$$C^\sqsubseteq(\bar{q}, a^{CL}) := \begin{cases} H & \text{if } \bar{q}(a) < \infty & \text{(a)} \\ M & \text{if } \begin{cases} CL = DM \wedge a \notin \bar{\mathcal{D}} \wedge |\bar{\mathcal{D}}| = D \\ CL = BE \wedge a \notin \bar{\mathcal{B}} \wedge |\bar{\mathcal{B}}| = B \end{cases} & \text{(b)} \\ \top & \text{otherwise} & \text{(c)} \end{cases} \tag{10}$$

379  
 380 We provide a complete step-by-step example on how *must*-analysis can be applied to an  
 381 application's CFG in Section 5.

## 382 4.2 May-analysis

383 The complete *may*-analysis is provided in the appendix (Section 10). We hereby provide a  
 384 sketch of the approach followed in the analysis.

385 The goal of *may*-analysis is to track the lower-bound on the age of memory blocks. Given  
 386 a *may*-analysis abstract state it is possible to classify a memory access as always leading to a  
 387 miss. Let us consider the example in Figure 3 and reason on the lower bound on the age of

each block for a 4-way fully associative cache. For block  $a^{DM}$ , the best case is represented by the execution pattern 1-5-4. In this case, the block has DM-LRU age 0. A similar situation occurs for block  $b^{DM}$  and path 1-3-4. For blocks  $f$  and  $g$ , the best-case is represented by the paths 2-6-8, and 2-7-8, respectively. This leads the two blocks to have a lower-bound of 2 on their DM-LRU age. Similarly, blocks  $c$ ,  $d$ , and  $g$  have lower-bound 0, 1, and 3, respectively.

We can represent the resulting *may*-analysis state obtained following the derivation above as:  $[\{a, b\}], [\{c\}, \{d\}, \{f, e\}, \{g\}]$ . What happens if another access to  $a$  occurs after path 4 and 8 join? Then the best-case for block  $b$  is still 1-3-4, but its age lower-bound will be 1. At the same time, because at least one DM block was allocated regardless of the taken path, the minimum lower-bound on the age of any BE block has to be 1. Also note that regardless of the execution path taken, block  $g$  will be evicted. The result is the following *may*-analysis abstract state:  $[\{a\}, \{b\}], [\{\}, \{c\}, \{d\}, \{f, e\}]$ .

## 5 Analysis Example

In this section we provide a description of how DM-LRU *must-analysis* can be applied to a CFG once the target of each memory access is known. The original CFG of the considered C program code generated by the Heptane tool is shown in Figure 4. The program consists of a single loop with four iterations, where the first iteration has been unrolled. The program accesses 7 memory locations. These are  $\mathcal{B} = \{a, b, c, d, e, f, g\}$  and are visible in the various basic blocks as operands of load/store instructions.

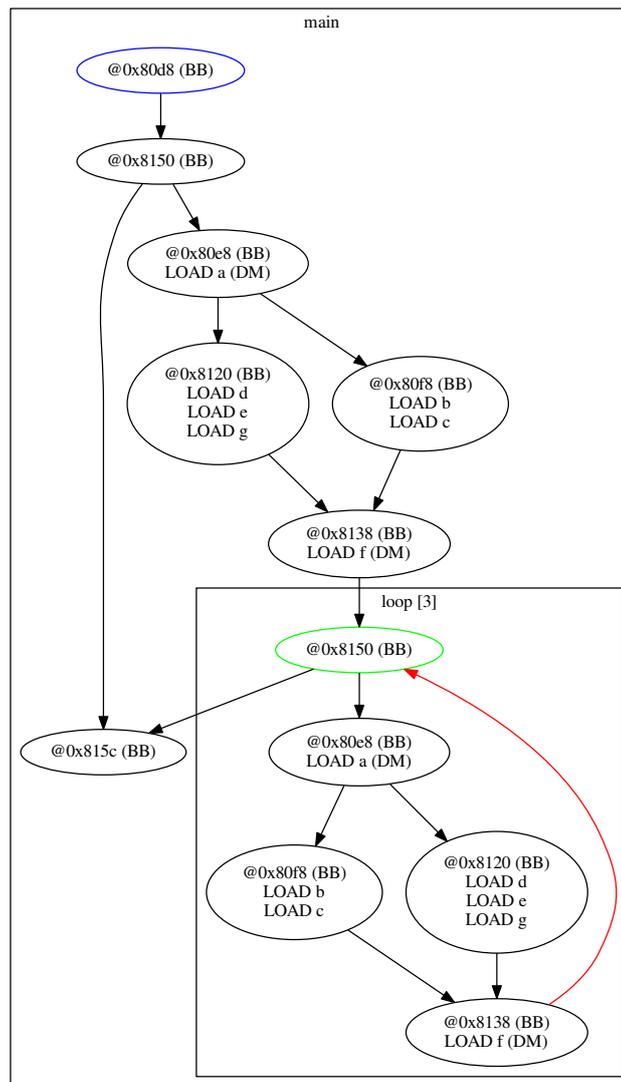
Figure 5 shows the same CFG as in Figure 4, but where only basic blocks in which memory accesses are being performed are kept. Moreover, basic blocks with multiple memory accesses are depicted as sequences of blocks, each with a single memory access. The nodes are annotated with their corresponding abstract states. We apply *must-analysis* starting from the entry node  $a$ . We compare the behavior of traditional LRU analysis and DM-LRU when blocks  $a$  and  $f$  have been declared as DM. We consider a fully-associative cache with 4 ways. For DM-LRU analysis, the cache state before the first access  $\bar{q}_0 = \square, [\{\}, \{\}, \{\}, \{\}]$  ( $D_0 = 0$ ); for LRU analysis it is  $[\{\}, \{\}, \{\}, \{\}]$ . Under DM-LRU, when block  $a^{DM}$  is accessed, the performed operation is  $\bar{q}_1 = U^{\square}(\bar{q}_0, a_{DM}) = U_D^{\square}(\bar{q}_0, a)$ . Following Equation 6, we have  $D_1 = D_0 + 1$ , then condition (a) is satisfied by  $a$ , all the other blocks  $b \in \mathcal{B}$  satisfy condition (d.2). As such, we have  $\bar{q}_1 = [\{a\}], [\{\}, \{\}, \{\}]$ , as reported in the figure.

Let us now follow the upper branch with access sequence  $d \rightarrow e \rightarrow g$  (all of them are best-effort memory accesses). For each memory access, we apply Equation 7 to obtain a new abstract state. After accessing  $e$ , the resulting abstract state is:  $\bar{q}_2 = [\{a\}][\{e\}, \{d\}, \{\}]$ . Let us now show more clearly how we obtain  $\bar{q}_3 = U^{\square}(\bar{q}_2, g^{BE}) = U_B^{\square}(\bar{q}_2, g)$ , when we next access  $g$ . Considering all blocks in  $\mathcal{B}$  and using Equation 7 we know: block  $a$  satisfies condition (b.1) and its age remains the same;  $b$  and  $c$  satisfy (d.2) and their age remains  $\infty$ ; block  $e$  satisfies (c) and its age increases by 1, from 1 to 2; the age of block  $d$  increases from 2 to 3; and finally, block  $g$  (being accessed) satisfies condition (a) and its age is set to  $D_2 = 1$ . The final state is  $\bar{q}_3 = [\{a\}], [\{g\}, \{e\}, \{d\}]$ , as shown in the figure above node  $g$ . The same procedure applies to the lower branch of the CFG, and we obtain the state  $\bar{q}_4 = [\{a\}], [\{c\}, \{b\}, \{\}]$ , after we access  $c$ .

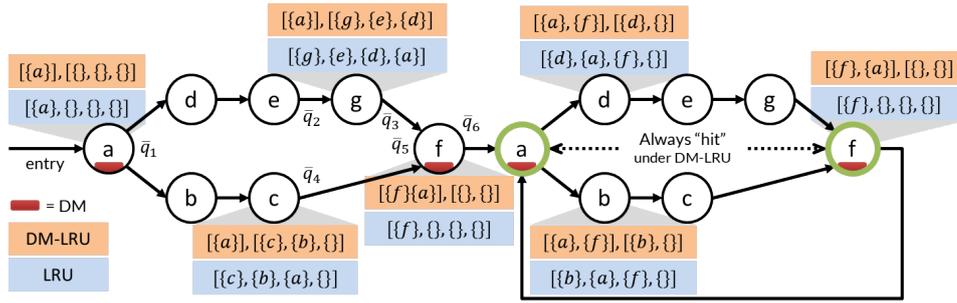
Before accessing  $f$ , we need to join states  $\bar{q}_3$  and  $\bar{q}_4$  derived above. In this case, we apply Equation 8 to obtain  $\bar{q}_5 = J^{\square}(\bar{q}_3, \bar{q}_4)$ . It follows that  $D_5 = 1$ . Moreover, the only block present in both states  $\bar{q}_3$  and  $\bar{q}_4$  is  $a$ . All other blocks in  $\mathcal{B}$  will have age  $\infty$  in  $\bar{q}_5$ . As such we have  $\bar{q}_5 = [\{a\}], [\{\}, \{\}, \{\}]$ . Next, accessing  $f^{DM}$  yields  $\bar{q}_6 = U_D^{\square}(\bar{q}_5, f) = [\{f\}, \{a\}], [\{\}, \{\}]$ , as shown in the figure. This is because  $D_6 = D_5 + 1$ , and because  $a$  satisfies condition (c.1) in Equation 6. The same reasoning can be applied to obtain the remaining states depicted in the figure.

Consider now the state  $\bar{q}_6$  and apply the *must-analysis* classifier before accessing  $a^{DM}$ ,

## 20:12 Impact of DM-LRU on WCET: a Static Analysis Approach



■ **Figure 4** Original CFG of considered example as rendered by the Heptane tool, with annotated memory accesses (LOAD). Note that VIVU has been performed on the loop.



**Figure 5** An example of *must* analysis under DM-LRU (orange states), compared to traditional LRU (blue states). If  $a$  and  $f$  are marked as DM, their accesses inside the loop can be classified as always hits.

437 i.e. compute  $C^{\sqsubseteq}(\bar{q}_6, a)$  as in Equation 10. First, the sets  $\bar{D}_6$  and  $\bar{B}_6$  can be computed using  
 438 Equation 9 as  $\bar{D}_6 = \{a, f\}$ , and  $\bar{B}_6 = \{\}$ . Hence condition (a) is satisfied and access to  $a$  is  
 439 classified as  $H$  (hit). Conversely, no access can be classified as hit under LRU.

## 6 Analogies and Differences with Cache Locking

441 Cache locking refers to a technique where cache blocks that are deemed important for  
 442 an application's timing are *pinned* (locked) in cache. Similar to DM-LRU, cache locking  
 443 represents a way to partially override the best-effort replacement strategy offered by the  
 444 hardware. And like DM-LRU, specialized hardware support is required to perform locking.  
 445 With respect to WCET analysis, the big advantage provided by cache locking is that all those  
 446 accesses for locked cache blocks can be immediately classified as hits. While cache locking  
 447 was commonly supported in previous-generation embedded systems, the current trend in  
 448 embedded SoCs is toward cache controllers that offer little or no management primitives.

449 Despite the strong similarities, some profound differences exist between cache locking  
 450 and DM-LRU. Leveraging cache locking implies injection of additional logic—in either the  
 451 application, the compiler, and/or the OS—to perform a series of prefetch&lock operations.  
 452 On the contrary, a system featuring a DM-LRU cache only requires that memory blocks are  
 453 tagged appropriately at task load time.

454 In case of static locking, prefetch&lock can be performed at initialization time. As such,  
 455 the extra logic required to perform locking does not impact the task's WCET. Conversely,  
 456 with dynamic cache locking, the locked cache content is modified at runtime. Depending on  
 457 the available hardware support, this operation may not be directly possible in user-space,  
 458 requiring instead a costly switch to kernel-space. Regardless, an online prefetch&lock routine  
 459 can pollute the rest of the cache, resulting in overheads that may largely offset any benefit.  
 460 In other words, additional system-level assumptions are required to make a meaningful  
 461 comparison with dynamic locking. For this reason, we do not compare DM-LRU against  
 462 dynamic locking.

463 Interestingly enough, however, the proposed DM-LRU analysis can be re-used to analyze  
 464 dynamic locking schemes if additional system parameters are available. In a nutshell, consider  
 465 a 4-way fully associative cache. Next, assume that the locked content is switched whenever  
 466 a given branch in the CFG is taken. Then, consider the case where the new content to  
 467 be locked is comprised of blocks  $a, b, c$ . A special node on the branch can be added with  
 468 associated a modified update abstract transformer  $Lock^{\sqsubseteq}$ . This is such that the resulting  
 469 *must*-analysis abstract state  $\bar{q}$  after the update is:  $\bar{q} = Lock^{\sqsubseteq}(\{a, b, c\}) = [\{a\}, \{b\}, \{c\}][\{\}$ .

## 470 **7** Evaluation

471 The DM-LRU analysis presented in the previous sections provides a way to understand how  
 472 the WCET of applications varies as memory blocks addressed in applications are declared  
 473 as DM. We now apply DM-LRU analysis to a set of realistic embedded benchmarks. In  
 474 this section, we first briefly describe our implementation. Next, we investigate three main  
 475 aspects: (1) what is the degree of WCET improvement that can be achieved via DM-LRU  
 476 when compared to LRU? (2) Is there an advantage in imposing a limit to the number of DM  
 477 lines that can be simultaneously allocated, i.e. in preventing DM takeover? (3) how does  
 478 DM-LRU compares to static cache locking?

479 In our evaluation we focus on the degree of WCET improvement that DM-LRU can  
 480 provide compared to LRU. However, because supporting DM-LRU implies changes to the  
 481 hardware cache memory and controller, it is also important to determine if a DM-LRU  
 482 implementation can be efficiently carried out. In short, only one additional bit to distinguish  
 483 between DM and BE lines is required per cache line. Additionally, compact changes<sup>1</sup> are  
 484 required to the cache controller to restrict victim selection based on the classification (DM  
 485 or BE) of a new line being allocated. No additional logic is required to appropriately set the  
 486 DM bits at line fetch. Additional considerations on the incurred hardware costs to support  
 487 DM memory are also provided in [10].

### 488 **7.1** Implementation

489 We have implemented support for DM-LRU inside a state of the art static WCET analysis tool,  
 490 namely Heptane [23]. Heptane implements Implicit Path Enumeration Technique (IPET) [29]  
 491 and performs analysis for many cache architectures: e.g., LRU, FIFO, Pseudo-LRU, multi-  
 492 level non-inclusive caches, and shared caches. In our setup, we consider a single-level of  
 493 cache, divided into an instruction (I) cache, and a data (D) cache. For simplicity, we assume  
 494 in all our experiments that both caches are selected to have the same number of sets and  
 495 ways. Heptane supports two architectures: ARMv7 and MIPS. The ARMv7 target was used  
 496 for this paper.

497 We have modified the Heptane tool to support two variants of DM-LRU, as well as  
 498 static locking. Most importantly, we have extended the support for abstract interpretation-  
 499 based cache analysis to implement the *must*- and *may*-analysis presented in the previous  
 500 sections. The performed changes allow backward compatibility with the original set of  
 501 policies supported by the tool. Next, we have integrated the logic to differentiate between BE  
 502 memory and DM memory. For this purpose, we have added a table of DM addresses—the  
 503 DM Table—that can be specified by an external tool, mimicking the selection of DM blocks  
 504 by the OS at binary load time. Furthermore, we have added appropriate logic in Heptane  
 505 to output per-memory-block statistics in terms of references, hits, and misses, as computed  
 506 during WCET analysis. These statistics are then used to build a DM-block selection heuristic.  
 507 Finally, we have modified Heptane’s CFG extraction routines to perform VIVU—i.e., to  
 508 recursively peel the first iteration of every loop.

509 We have developed and employed a simple heuristic to determine which memory block-  
 510 s/addresses should be marked as DM and inserted into the DM Table. The heuristic initially  
 511 performs WCET analysis without selecting any DM line. Next, it analyzes the per-memory-  
 512 block statistics and selects as DM the block with the largest number of misses. At this point,  
 513 WCET analysis is performed again with the new DM Table containing a single entry. Using

---

<sup>1</sup> Whenever a line eviction has to occur, the DM/BE bits of all the lines in the considered set form a bitmask. Victim selection for a BE access is then performed by excluding all those lines that have a bit set to 1 in the DM bitmask.

514 the per-memory-block statistics of the latest run, a new DM block is selected in addition to  
 515 the previously selected block. The same steps are performed until no more addresses can be  
 516 selected as DM. Note that when no lines are selected as DM, the behavior of the cache is  
 517 indistinguishable from vanilla LRU. Similarly, when the entire memory of an applications is  
 518 selected as DM, no differences exist with LRU. In practice, however, we saw no differences  
 519 between DM-LRU and LRU when more than  $3 \times S \times A$  lines are selected as DM, where  $S$   
 520 and  $A$  is the number of sets and ways of the cache, respectively. In our experiments, we  
 521 acquired analytical results for a number of DM lines in the range  $[1, 3 \times S \times A]$ .

## 522 7.2 Setup

523 We compare two variants of DM-LRU and static cache locking against LRU. A description  
 524 of the three scenarios follows.

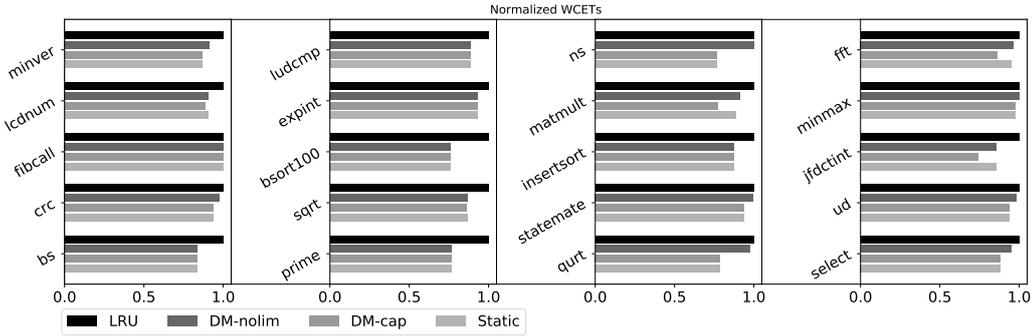
- 525 **1. Unrestricted DM-LRU ("DM-nolim"):** in this variant, no restriction is imposed on  
 526 the maximum number of cache sets that can be reserved for DM lines. It follows that  
 527 the only constraint for the allocation of DM lines is the cache associativity itself. The  
 528 analysis for unrestricted DM-LRU is the one presented in the previous sections.
- 529 **2. Limited DM-LRU ("DM-cap"):** in this variant, a hard cap in the maximum number  
 530 of ways is imposed on the expansion of DM lines. This represents a solution to the  
 531 aforementioned problem of *DM takeover*. Imposing a cap of 0 makes DM-cap to be  
 532 identical to vanilla LRU. Similarly, imposing a *cap* of  $A$  makes DM-cap to be identical  
 533 to DM-nolim. In our experiments, we explore all the possible values of *cap* in the range  
 534  $[1, A]$ .
- 535 **3. Static locking ("Static"):** this case is used to draw a comparison between the considered  
 536 DM-LRU variants and static locking. In case of static locking, selection of lines to statically  
 537 allocate is performed following the same heuristic used for DM lines selection. Similar to  
 538 DM-cap, we impose how many ways can be dedicated to statically locked content (locked  
 539 ways). The maximum number of allocated line is then  $S \times \textit{locked}$ . Note that the main  
 540 performance difference between DM-cap and Static lies in the additional flexibility that  
 541 DM-cap provides. In DM-cap, in fact, more lines than  $S \times \textit{cap}$  can be selected, while it is  
 542 not allowed in static locking.

543 For all the considered variants, we explore a number of cache configurations. Specifically,  
 544 we vary the associativity  $A$  of the I/D caches in the set  $\{2, 4, 8, 16\}$ . We vary the number  
 545 of cache sets  $S$  such that  $S \in \{2, 4, 8, 16, 32\}$ . As previously mentioned, for DM-nolim and  
 546 DM-cap, we progressively select up to  $3 \times S \times A$  DM lines following the heuristic described  
 547 above. In each system instance, we perform WCET analysis using the modified Heptane  
 548 tool. Then, we keep track of which configuration— $S$ ,  $A$ , DM-lim *cap*, *locked* ways, number  
 549 of DM lines—for each of the three scenarios leads to the best reduction in WCET compared  
 550 to the vanilla LRU case.

551 For our benchmarks, we use a subset of realistic benchmarks from the Mälardalen suite [19].  
 552 Unfortunately, vanilla HEPTANE is not able to perform WCET analysis for some of the  
 553 benchmarks in the suite. As such, our evaluation only includes those benchmarks that  
 554 are correctly handled by HEPTANE. Notably, the aforementioned changes to implement  
 555 DM-LRU analysis did not impact the set of benchmarks that can be correctly analyzed by  
 556 the tool.

## 557 7.3 Results

558 Figure 6 provides an overview of the obtained results. A cluster of bars is provided for each of  
 559 the considered benchmarks. Reading the plot from top to bottom, the first bar corresponds  
 560 to the WCET under LRU. All the results in the figure are normalized to the LRU case. The



■ **Figure 6** Computed WCETs for vanilla LRU (LRU), unrestricted DM (DM-nolim), DM limited to a subset of ways (DM-cap), and static locking (Static).

561 second bar represents the best WCET improvement that was observed under DM-nolim. The  
 562 WCET improvement is calculated as:  $\frac{WCET_{DM-nolim}}{WCET_{LRU}}$ , where the WCETs under DM-nolim  
 563 and under LRU are obtained in the same system configuration. A similar calculation was  
 564 performed to derive the remaining two bars, i.e. for the DM-lim and Static cases.

565 What emerges from the plot is that in 16 out of 20 cases, DM-nolim is able to achieve  
 566 WCET reductions compared to vanilla LRU. Notably, in case of `bsort100` and `prime`, it is  
 567 possible to achieve a WCET reduction of around 23.73% and 23.47%, respectively. It can  
 568 also be noted that DM-cap outperforms DM-nolim. Moreover, DM-cap performs generally  
 569 better than static locking. For instance, the best WCET reduction achieved via static locking  
 570 for the `jfdctint` benchmark is 26.09% under DM-cap (with a  $S = 4$ ,  $A = 8$ , 19 DM lines,  
 571  $cap = 1$ ). But the best WCET reduction under static locking is only 14.64%, which is  
 572 achieved for a cache with parameters  $S = 2$ ,  $A = 16$  and 15 ways occupied by statically  
 573 locked lines. Similar results can be observed for the benchmarks `matmult` and `fft`.

574 The reason for the performance improvement that can be obtained with DM-cap is  
 575 twofold. On the one hand, the problem of DM takeover is solved. This prevents the case that  
 576 all the accesses to BE lines result in misses. On the other hand, for applications that exhibit  
 577 changes in working sets, static locking can be sub-optimal. Conversely, under DM-cap, is  
 578 possible to mark lines belonging to different working sets as DM. In this case, at working  
 579 set changes over time, those DM lines belonging to a previous working set will be naturally  
 580 evicted, without suffering pollution from BE lines.

581 A more detailed overview of the obtained experimental results is provided in Table 1. In  
 582 the table, the first column reports the name of the benchmark under analysis. If multiple  
 583 configurations are of interest, multiple rows are shown for a given benchmark. The second  
 584 column reports the cache configuration in terms of sets  $S$  and ways  $W$  for the results on  
 585 each row. Next, the WCET obtained with LRU is reported in the following column, followed  
 586 by the best WCET obtained for the same configuration under DM-nolim and the relative  
 587 improvement (due to the space limitation, the number of DM lines that were selected has  
 588 been omitted in the table.) Similarly, the best result obtained under DM-cap is reported  
 589 next, and the value of  $cap$  under which the result was achieved is reported in the adjacent  
 590 column. Finally, the last two columns report the WCET (and the relative improvement) for  
 591 static locking with the given cache configuration and number of locked ways reported in the  
 592 last column.

## 593 8 Related Work

594 **Memory Tagging and Hardware Support.** In this work, we assume that hardware

| Benchmark  | $S \times A$ | LRU      | DM-nolim          | DM-cap            | cap | Static            | locked |
|------------|--------------|----------|-------------------|-------------------|-----|-------------------|--------|
| bs         | 2×2          | 6613     | 5513 (-16.63%)    | 5513 (-16.63%)    | 2   | 5513 (-16.63%)    | 2      |
| crc        | 4×2          | 2492320  | 2425920 (-2.66%)  | 2330620 (-6.49%)  | 1   | 2330620 (-6.49%)  | 1      |
| fibcall    | 2×2          | 14191    | 14191 (-0.00%)    | 14191 (-0.00%)    | 1   | 14191 (-0.00%)    | 1      |
| lcdnum     | 4×2          | 16291    | 14791 (-9.21%)    | 14791 (-9.21%)    | 2   | 14791 (-9.21%)    | 2      |
|            | 2×4          | 16191    | 16191 (-0.00%)    | 14391 (-11.12%)   | 2   | 15291 (-5.56%)    | 2      |
| minver     | 4×2          | 126558   | 115758 (-8.53%)   | 109958 (-13.12%)  | 1   | 109958 (-13.12%)  | 1      |
| prime      | 2×4          | 611425   | 467925 (-23.47%)  | 467925 (-23.47%)  | 3   | 467925 (-23.47%)  | 3      |
|            | 2×4          | 54983    | 47552 (-13.52%)   | 47252 (-14.06%)   | 3   | 52552 (-4.42%)    | 2      |
| sqrt       | 2×4          | 54983    | 47552 (-13.52%)   | 47252 (-14.06%)   | 3   | 47583 (-13.30%)   | 6      |
|            | 2×2          | 12434700 | 9484580 (-23.72%) | 9484580 (-23.72%) | 1   | 9484580 (-23.72%) | 1      |
| bsort100   | 2×2          | 12434700 | 9484580 (-23.72%) | 9484580 (-23.72%) | 1   | 9484580 (-23.72%) | 1      |
| expint     | 2×4          | 759551   | 709651 (-6.57%)   | 709651 (-6.57%)   | 4   | 709651 (-6.57%)   | 4      |
| ludcmp     | 16×2         | 638233   | 564633 (-11.53%)  | 564633 (-11.53%)  | 2   | 564633 (-11.53%)  | 2      |
| qurt       | 2×8          | 217555   | 212160 (-2.48%)   | 173755 (-20.13%)  | 6   | 173755 (-20.13%)  | 6      |
|            | 4×4          | 217555   | 220355 (-1.29%)   | 171155 (-21.33%)  | 3   | 171155 (-21.33%)  | 3      |
| statemate  | 2×2          | 616218   | 612918 (-0.54%)   | 576718 (-6.41%)   | 1   | 576718 (-6.41%)   | 1      |
|            | 8×8          | 383718   | 382818 (-0.23%)   | 359118 (-6.41%)   | 6   | 359118 (-6.41%)   | 6      |
| insertsort | 2×2          | 80126    | 70126 (-12.48%)   | 70126 (-12.48%)   | 1   | 70126 (-12.48%)   | 1      |
| matmult    | 2×2          | 7191620  | 6568220 (-8.67%)  | 5555520 (-22.75%) | 1   | 6391620 (-11.12%) | 1      |
|            | 4×2          | 193481   | 193481 (-0.00%)   | 193481 (-0.00%)   | 1   | 193481 (-0.00%)   | 1      |
| ns         | 2×2          | 530781   | 534781 (-0.75%)   | 406681 (-23.38%)  | 1   | 406681 (-23.38%)  | 1      |
|            | 4×2          | 170766   | 162266 (-4.98%)   | 157966 (-7.50%)   | 1   | 157966 (-7.50%)   | 1      |
| select     | 2×4          | 170766   | 162866 (-4.63%)   | 150966 (-11.59%)  | 3   | 150966 (-11.59%)  | 3      |
|            | 4×2          | 226843   | 223243 (-1.59%)   | 223243 (-1.59%)   | 2   | 225243 (-0.71%)   | 2      |
| ud         | 2×2          | 302443   | 354143 (-17.09%)  | 283843 (-6.15%)   | 1   | 283843 (-6.15%)   | 1      |
|            | 2×16         | 150234   | 128734 (-14.31%)  | 111034 (-26.09%)  | 2   | 128234 (-14.64%)  | 15     |
| jfdctint   | 4×8          | 150234   | 130334 (-13.25%)  | 111134 (-26.03%)  | 1   | 147134 (-2.06%)   | 1      |
|            | 4×8          | 150234   | 130334 (-13.25%)  | 111134 (-26.03%)  | 1   | 130334 (-13.25%)  | 7      |
| minmax     | 2×2          | 4034     | 4034 (-0.00%)     | 4034 (-0.00%)     | 1   | 4034 (-0.00%)     | 1      |
|            | 2×4          | 4034     | 4034 (-0.00%)     | 3934 (-2.48%)     | 1   | 4034 (-0.00%)     | 1      |
| fft        | 32×2         | 1683830  | 1623930 (-3.56%)  | 1623430 (-3.59%)  | 1   | 1623430 (-3.59%)  | 1      |
|            | 4×4          | 2488230  | 2494360 (-0.25%)  | 2140830 (-13.96%) | 1   | 2443230 (-1.81%)  | 1      |
|            | 4×4          | 2488230  | 2494360 (-0.25%)  | 2140830 (-13.96%) | 1   | 1716630 (-4.62%)  | 2      |

■ **Table 1** Summary of notable experimental results under four strategies: (1) vanilla LRU ("LRU"); (2) unrestricted DM-LRU ("DM-nolim"); (3) restricted DM-LRU ("DM-cap"); and (4) static locking ("Static").

allows us to encode (tag) extra information (e.g., importance) on memory locations at a fine-granularity. The basic idea of memory tagging has first explored in the security community, to prevent memory corruption (e.g., buffer overflow) [6, 38] and to enforce data flow integrity [45] and capability protection [51]. Efficient hardware designs for word-granularity single-bit and multi-bit memory tagging have been investigated [24] and several real SoC prototypes have been built [45, 4], demonstrating the feasibility. In the real-time systems community, several works explored the use physical memory address based differentiated hardware designs (mostly cache) in a more coarse-grained manner (i.e., memory segments, page, and task granularity). Kumar et al, proposed a criticality-aware cache design, called Least Critical (LC), which includes a memory criticality-aware extension to LRU replacement policy [27]. The LC cache's replacement policy is similar to the replacement policy we assumed in this work (Algorithm 1), while its memory tagging mechanism, which uses a fixed number of specialized range registers, does not allow flexible and fine-grained memory tagging. Therefore, our static analysis method can be directly applicable to analyze the LC cache. PRETI [28] also proposes a criticality-aware cache design but it focuses on shared cache for SMT hardware, while we focus on private caches. More recently, OS-level page-granularity memory tagging and supporting multicore architecture designs (including a new cache design) have been explored to provide efficient hardware isolation (incl. cache isolation) in multicore [10].

613 **Static Cache Analysis.** There exists a broad literature on static cache analysis [32, 50].  
 614 With respect to existing literature, this work is closely related to approaches that propose  
 615 abstract interpretation-based cache analysis. This approach was initially proposed in [1, 12].  
 616 These works illustrate LRU analysis and hit/miss classification using *may*- and *must*-analysis.  
 617 The work in [12] also proposes a persistence analysis based on abstract states, which was found  
 618 to be unsafe and for which a fix was proposed in [8, 25]. We base our DM-LRU extension on  
 619 the *may*- and *must*-analysis proposed in [12], but use the improved formalization in [14]. In  
 620 order to perform access classification in case of loops we use an approach similar to *virtual*  
 621 *inlining & virtual unrolling* (VIVU) originally proposed in [34]. A large body of works has  
 622 considered cache replacement policies other than LRU. These include FIFO [14, 15, 18],  
 623 MRU [17], Pseudo-LRU [16]. Comparatively less work has been produced to analyze non-  
 624 inclusive [36, 21] as well as inclusive [22] multi-level caches. With respect to these works,  
 625 the proposed methodology set itself apart because it focuses on the impact on the WCET  
 626 of designer-driven selection of frequently accessed memory blocks. In this sense, proposed  
 627 approach can be used to analyze caches that support the definition of touch-and-lock cache  
 628 lines, under the assumption that no more than  $A$  blocks are simultaneously locked on any  
 629 set, where  $A$  is the associativity of the cache.

630 **Cache Locking and Scratchpad Memory.** Some COTS cache designs [2, 7, 13]  
 631 support selective cache locking, which prevents evictions for certain programmer selected  
 632 cache-lines. Exploiting the feature, various static and dynamic cache locking schemes for both  
 633 instruction and data caches have been investigated [5, 40, 39, 48, 35]. In [47, 48], for instance,  
 634 cache locking statements are inserted in the task's execution flow at compilation time, when  
 635 the uncertainty about the memory locations being accessed negatively impacts the static  
 636 WCET analysis. Some recent works combined cache locking with cache partitioning to  
 637 improve task WCET in multicore [30, 43, 33]. As an alternative to cache, scratchpad memory  
 638 has received significant attention in the real-time systems community for its predictability  
 639 benefits [46, 9, 49, 44]. More recently, a technique called invalidation-driven allocation  
 640 (IDA) [26] was proposed to achieve the same level of determinism of a locked cache in  
 641 spite of lack of hardware-assisted locking primitives. IDA can be used as long as precise  
 642 invariants on the size of an application's working set hold. To overcome its high programming  
 643 complexity, however, many researchers proposed various compiler-based techniques. In [44],  
 644 for instance, a sophisticated compiler-based technique is proposed to break each task into  
 645 intervals and at the beginning of each interval, the required memory blocks of the interval  
 646 are prefetched onto a scratchpad memory via a DMA controller without blocking the task  
 647 execution. Dividing a task into a sequence of well-defined memory and computation phases  
 648 was originally proposed in [37, 52]. In both cache locking and scratchpad memory based  
 649 techniques, a common limitation is the overhead of explicitly executing additional instructions  
 650 (prefetch, lock/unlock, or data movement to/from scratchpad). Furthermore, these additional  
 651 instructions are context sensitive in the sense that they must be executed before actual accesses  
 652 occur, and if they are executed too early, they can negatively impact both performance and  
 653 WCET. In contrast, our approach is context insensitive in the sense that, once DM blocks  
 654 are flagged, actual allocation and replacement are automatically performed by the hardware  
 655 (cache controller) without additional instruction execution overhead.

## 656 **9 Conclusion**

657 In this paper, we presented the DM-LRU cache replacement policy and proposed an abstract  
 658 interpretation-based analysis for DM-LRU. We implemented the proposed analysis and DM-  
 659 LRU support in the Heptane static WCET analysis tool. Using the Heptane, we evaluated  
 660 the WCET impacts of our DM-LRU based approach on a number of benchmarks. The results  
 661 show that our DM-LRU approach can provide lower task WCETs with less performance

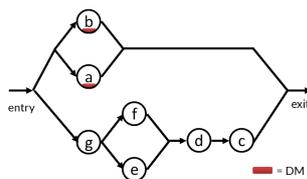


Figure 7 Fragment of process CFG that leads to abstract DM-LRU state  $\underline{q} = [\{a, b\}], [\{c\}, \{d\}, \{e, f\}, \{g\}]$ .

overhead and programming complexity, compared to the standard LRU and cache locking based approaches.

## Acknowledgements

We are especially grateful to Daniel Grund for making his research thesis [14] promptly available to us. This research is supported in part by NSF CNS 1718880. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the authors and do not necessarily reflect the views of the NSF.

## 10 Appendix: May Analysis

In the DM-LRU analysis framework, *may*-analysis is once again performed by considering abstract cache states. Recall that *may*-analysis keeps track of the lower-bound on the age of each addressable memory block. There are a number of differences compared to the analytical tools used for *must* analysis. In *may*-analysis it is necessary to keep track of both  $D \in \{0, \dots, A\}$  and  $B \in \{0, \dots, A\}$ . Here, the meaning of  $D$  and  $B$  changes. In this case,  $D$  represents the maximum lower-bound of any possibly cached DM block. Conversely,  $B$  captures the minimum lower-bound on the DM-LRU age of any BE block. It may be the case that  $B + D > A$  in order to correctly abstract the age lower-bound resulting from multiple execution paths. It must hold however that  $A \leq D + B \leq 2A$ . It follows that the abstract domain for *may*-analysis  $DMLru_A^\exists$  is defined as:

$$DMLru_A^\exists := \{0, \dots, A\} \times \{0, \dots, A\} \times \mathcal{B} \rightarrow \{0, \dots, A - 1, A\}. \quad (11)$$

An abstract state  $\underline{q} \in DMLru_A^\exists$  is then represented as two sets of memory blocks, for instance:  $\underline{q} = [\{a, b\}], [\{c\}, \{d\}, \{e, f\}, \{g\}] \in DMLru_A^\exists$ . In this example, we have  $D = 1, B = 4, A = 4$ . It follows that the upper-bound on the number of DM memory blocks is 1, and that blocks  $a$  and  $b$  have at least DM-LRU age 0, and *may* be marked as deterministic blocks. On the other hand,  $c$  is a best-effort memory block with DM-LRU age at least 0. It should not come as a surprise that in some states  $D + B > A$ . Consider the execution depicted in Figure 7 that produces  $\underline{q}$ . When execution reaches the end of the figure, there could be 0 or 1 DM blocks allocated in cache. Hence the upper-bound on the number of DM blocks has to be  $D = 1$ . On the other hand, the upper bound on the number of BE blocks is  $B = 4$ .

The operator  $DM^\exists(\underline{q}, a)$  takes an abstract state  $\underline{q}$  and a block  $a$ , and returns *true* if *may* be allocated as a DM block in  $\underline{q}$ . For ease of notation, we simply use  $DM^\exists(a)$  when the considered abstract state is obvious. We define the operator  $DM^\exists(\underline{q}, a)$  as follows:

$$DM^\exists(\underline{q}, b^{CL}) := \begin{cases} true & \text{if } CL = DM \wedge \underline{q}(b) < A \\ false & \text{otherwise.} \end{cases} \quad (12)$$

The update abstract transformer  $U_D^\exists$  for a DM memory access  $a$  can be defined as follows:

$$U_D^\exists(\underline{q}, a) :=$$

$$697 \quad D' \leftarrow \begin{cases} D+1 & \text{if } D < A \wedge BE^\sqsupset(a) \\ D+1 & \text{if } D < A \wedge \exists x^{DM} \neq a : \underline{q}(x) = \underline{q}(a) = D-1, \\ D & \text{otherwise} \end{cases} \quad (13)$$

$$698 \quad B' \leftarrow \begin{cases} B-1 & \text{if } B > 0 \wedge \underline{q}(a) \geq A-B \\ B & \text{if } B = 0 \vee \underline{q}(a) < A-B \end{cases} \quad (14)$$

$$699 \quad \lambda b. \begin{cases} 0 & \text{if } b = a & (a) \\ \underline{q}(b) & \text{if } b \neq a \wedge \\ & \left\| \begin{array}{l} DM^\sqsupset(b) \wedge \underline{q}(a) < \underline{q}(b) \\ BE^\sqsupset(b) \wedge \underline{q}(a) < A-B \end{array} \right. & (b) \\ \underline{q}(b)+1 & \text{if } b \neq a \wedge \\ & \left\| \begin{array}{l} DM^\sqsupset(b) \wedge \underline{q}(a) \geq \underline{q}(b) \wedge \underline{q}(b) < D'-1 \\ BE^\sqsupset(b) \wedge \underline{q}(a) \geq A-B \wedge \underline{q}(b) < A-1 \end{array} \right. & (c) \\ A & \text{if } b \neq a \wedge \\ & \left\| \begin{array}{l} DM^\sqsupset(b) \wedge \underline{q}(a) \geq \underline{q}(b) \wedge \underline{q}(b) \geq D'-1 \\ BE^\sqsupset(b) \wedge \underline{q}(a) \geq A-B \wedge \underline{q}(b) \geq A-1 \end{array} \right. & (d) \end{cases} \quad (15)$$

700  
701  
702 Where  $D'$  ( $B'$ , resp.) is the new value of  $D$  ( $B$ , resp.) after the update. Similarly, the  
703 update abstract transformer  $U_{\bar{B}}^\sqsupset$  for a best-effort memory access  $a$  can be defined as follows:

$$704 \quad U_{\bar{B}}^\sqsupset(\underline{q}, a) := \quad (16)$$

$$705 \quad \lambda b. \begin{cases} A-B & \text{if } b = a & (a) \\ \underline{q}(b) & \text{if } b \neq a \wedge \left\| \begin{array}{l} DM^\sqsupset(b) \\ BE^\sqsupset(b) \wedge \underline{q}(a) < \underline{q}(b) \end{array} \right. & (b) \\ \underline{q}(b)+1 & \text{if } b \neq a \wedge BE^\sqsupset(b) \wedge \underline{q}(a) \geq \underline{q}(b) \wedge \underline{q}(b) < A-1 & (c) \\ A & \text{if } b \neq a \wedge BE^\sqsupset(b) \wedge \underline{q}(a) \geq \underline{q}(b) \wedge \underline{q}(b) \geq A-1 & (d) \end{cases} \quad (17)$$

706  
707  
708 To clarify how the  $U^\sqsupset$  operation transforms a given state, consider the abstract state  
709  $\underline{q} = [\{a, b\}], [\{c\}, \{d\}], \{e, f\}, \{g\}]$ , where  $D = 1, B = 4$ . Assume that DM block  $h$  is accessed,  
710 whose DM-LRU age is currently  $A$  or higher. First, the value of  $D'$  ( $B'$ , resp.) is computed  
711 as  $D' = D + 1$  ( $B' = B - 1$ , resp.); next,  $\{a, b\}$  both satisfy the fourth condition in  $U_{\bar{D}}^\sqsupset$   
712 —Equation 15, first case of (c); block  $c, d, e$  and  $f$  satisfy the fifth condition;  $g$  the seventh.  
713 The resulting updated abstract state is:  $\underline{q}' = [\{h\}, \{a, b\}], [\{\}, \{c\}, \{d\}], \{e, f\}$ . Note that in  
714 the resulting state  $B = 3$ , hence the least lower-bound on any BE block is  $A - B = 1$ .

715 **May-analysis Join:** The join abstract transformer for DM-LRU *may*-analysis is sym-  
716 metric to the join abstract transformer used for DM-LRU *must*-analysis. The joined state  
717 will contain all the blocks in the union of the joining states, each with the minimum age  
718 in any of the two states. Furthermore,  $D$  is taken as the maximum between the value of  
719  $D$  in the joining states. Similarly,  $B$  is taken as the maximum between the value of  $B$  in  
720 the joining states. As such, after a join, it always holds that  $D + B \leq 2A$ . Equation 18  
721 formalizes the  $J^\sqsupset(\bar{q}, \bar{p})$  abstract transformer:

$$722 \quad J^\sqsupset(\bar{q}, \bar{p}) := D \leftarrow \max\{D_{\bar{q}}, D_{\bar{p}}\}, B \leftarrow \max\{B_{\bar{q}}, B_{\bar{p}}\}, \lambda b. \min\{\bar{q}(b), \bar{p}(b)\}. \quad (18)$$

723 To clarify the join operation, consider the state  $\underline{q} = [\{a, b\}], [\{c\}, \{d\}], \{e, f\}, \{g\}]$  obtained  
724 in Figure 7, and the state  $\underline{q}' = [\{h\}, \{a, b\}], [\{\}, \{c\}, \{d\}], \{e, f\}]$  obtained as  $\underline{q}' = U^\sqsupset(\underline{q}, h^{DM})$   
725 (i.e. by accessing the DM block  $h$ ). If we were to join  $\underline{q}$  with  $\underline{q}'$ , the resulting state would be  
726  $\underline{q}'' = [\{a, b, h\}], [\{c\}, \{d\}], \{e, f\}, \{g\}]$ .

727 **May-analysis Classification:** It is possible to classify a memory access using a classi-  
728 fication function that will either return  $M$  for cache miss, or  $\top$  in case access to a memory  
729 block cannot be guaranteed to be a miss given the current abstract state. The classification  
730 function of the *may* analysis is defined as:

$$731 \quad C^\sqsupset(\underline{q}, a^{CL}) := \begin{cases} M & \text{if } \underline{q}(a) = A \\ \top & \text{otherwise.} \end{cases} \quad (19)$$

## 733 — References —

- 734 **1** M. Alt, C. Ferdinand, F. Martin, and R. Wilhelm. Cache behavior prediction by abstract  
735 interpretation. In *Proceedings of the Third International Symposium on Static Analysis, SAS*  
736 '96, pages 52–66, Berlin, Heidelberg, 1996. Springer-Verlag.
- 737 **2** ARM. *PL310 Cache Controller Technical Reference Manual, Rev: r0p0*, 2007.
- 738 **3** R. Banakar, S. Steinke, B. Lee, M. Balakrishnan, and P. Marwedel. Scratchpad memory: design  
739 alternative for cache on-chip memory in embedded systems. In *Int. Symp. Hardware/Software*  
740 *Codesign (CODES+ISSS)*, pages 73–78. ACM, 2002.
- 741 **4** A. Bradbury, G. Ferris, and R. Mullins. Tagged memory and minion cores in the lowrisc soc.  
742 *Memo, University of Cambridge*, 2014.
- 743 **5** M. Campoy, A. P. Ivars, and J. Busquets-Mataix. Static use of locking caches in multitask  
744 preemptive real-time systems. In *Real-Time Embedded Systems Workshop (Satellite of the*  
745 *IEEE Real-Time Systems Symposium)*, pages 1–6, 2001.
- 746 **6** S. Chen, J. Xu, N. Nakka, Z. Kalbarczyk, and R. K. Iyer. Defeating memory corruption  
747 attacks via pointer taintedness detection. In *Dependable Systems and Networks (DSN)*, pages  
748 378–387. IEEE, 2005.
- 749 **7** N. Corp. Variable SMP – A Multi-Core CPU Architecture for Low Power and High Performance.  
750 Technical report, Nvidia, 2011.
- 751 **8** C. Cullmann. Cache persistence analysis: Theory and practice. *ACM Trans. Embed. Comput.*  
752 *Syst.*, 12(1s):40:1–40:25, Mar. 2013.
- 753 **9** J.-F. Deverge and I. Puaut. WCET-directed dynamic scratchpad memory allocation of data.  
754 In *Euromicro Conference on Real-Time Systems (ECRTS)*, pages 179–190. IEEE, 2007.
- 755 **10** F. Farshchi, P. K. Valsan, R. Mancuso, and H. Yun. Deterministic memory abstraction and  
756 supporting multicore system architecture. In *Euromicro Conf. Real-Time Syst. (ECRTS)*.  
757 IEEE, 2018.
- 758 **11** C. Ferdinand, F. Martin, R. Wilhelm, and M. Alt. Cache behavior prediction by abstract  
759 interpretation. *Sci. Comput. Program.*, 35(2-3):163–189, Nov. 1999.
- 760 **12** C. Ferdinand and R. Wilhelm. Efficient and precise cache behavior prediction for real-  
761 timesystems. *Real-Time Syst.*, 17(2-3):131–181, Dec. 1999.
- 762 **13** Freescale. *e500mc Core Reference Manual*, 2012.
- 763 **14** D. Grund. *Static Cache Analysis for Real-Time Systems: LRU, FIFO, PLRU*. epubli, 2012.
- 764 **15** D. Grund and J. Reineke. Precise and Efficient FIFO-Replacement Analysis Based on Static  
765 Phase Detection. In *Euromicro Conference on Real-Time Systems (ECRTS)*, pages 155–164,  
766 July 2010.
- 767 **16** D. Grund and J. Reineke. Toward precise PLRU cache analysis. In *International Workshop*  
768 *on Worst-Case Execution Time Analysis (WCET)*, pages 23–35, 2010.
- 769 **17** N. Guan, M. Lv, W. Yi, and G. Yu. WCET Analysis with MRU Caches: Challenging LRU for  
770 Predictability. In *Real Time and Embedded Technology and Applications Symposium (RTAS)*,  
771 pages 55–64, April 2012.
- 772 **18** N. Guan, X. Yang, M. Lv, and W. Yi. FIFO Cache Analysis for WCET Estimation: A  
773 Quantitative Approach. In *Design, Automation and Test in Europe (DATE)*, pages 296–301,  
774 San Jose, CA, USA, 2013. EDA Consortium.
- 775 **19** J. Gustafsson, A. Betts, A. Ermedahl, and B. Lisper. The Mälardalen WCET benchmarks –  
776 past, present and future. In B. Lisper, editor, *Proceedings of the 10th International Workshop*  
777 *on Worst-Case Execution Time Analysis (WCET'2010)*, pages 137–147, Brussels, Belgium,  
778 July 2010. OCG.
- 779 **20** S. Hahn and J. Reineke. Design and analysis of sic: A provably timing-predictable pipelined  
780 processor core. In *Real-Time Systems Symposium (RTSS)*, pages 469–481. IEEE, 2018.
- 781 **21** D. Hardy and I. Puaut. WCET Analysis of Multi-level Non-inclusive Set-Associative Instruction  
782 Caches. In *Real-Time Systems Symposium (RTSS)*, pages 456–466, Nov 2008.
- 783 **22** D. Hardy and I. Puaut. Wcet analysis of instruction cache hierarchies. *J. Syst. Archit.*,  
784 57(7):677–694, Aug. 2011.

- 785 **23** D. Hardy, B. Rouxel, and I. Puaut. The Heptane Static Worst-Case Execution Time Estimation  
786 Tool. In *17th International Workshop on Worst-Case Execution Time Analysis (WCET*  
787 *2017)*, volume 8 of *International Workshop on Worst-Case Execution Time Analysis*, page 12,  
788 Dubrovnik, Croatia, June 2017.
- 789 **24** A. Joannou, J. Woodruff, R. Kovacsics, S. W. Moore, A. Bradbury, H. Xia, R. N. Watson,  
790 D. Chisnall, M. Roe, B. Davis, et al. Efficient tagged memory. In *International Conference on*  
791 *Computer Design (ICCD)*, pages 641–648. IEEE, 2017.
- 792 **25** L. Ju, S. Chakraborty, and A. Roychoudhury. Accounting for cache-related preemption delay in  
793 dynamic priority schedulability analysis. In *Design, Automation and Test in Europe (DATE)*,  
794 pages 1623–1628, San Jose, CA, USA, 2007. EDA Consortium.
- 795 **26** T. Kloda, M. Solieri, R. Mancuso, N. Capodici, P. Valente, and M. Bertogna. Deterministic  
796 memory hierarchy and virtualization for modern multi-core embedded systems. In *2019 IEEE*  
797 *Real-Time and Embedded Technology and Applications Symposium (RTAS)*, Montreal, Canada,  
798 April 2019.
- 799 **27** N. C. Kumar, S. Vyas, R. K. Cytron, C. D. Gill, J. Zambreno, and P. H. Jones. Cache design  
800 for mixed criticality real-time systems. In *Computer Design (ICCD)*, pages 513–516. IEEE,  
801 2014.
- 802 **28** B. Lesage, I. Puaut, and A. Sez nec. Preti: Partitioned real-time shared cache for mixed-  
803 criticality real-time systems. In *Real-Time and Network Systems (RTNS)*, pages 171–180.  
804 ACM, 2012.
- 805 **29** Y. S. Li and S. Malik. Performance analysis of embedded software using implicit path  
806 enumeration. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and*  
807 *Systems*, 16(12):1477–1487, Dec 1997.
- 808 **30** T. Liu, Y. Zhao, M. Li, and C. J. Xue. Task assignment with cache partitioning and locking  
809 for wcet minimization on mp soc. In *2010 39th Int. Conf. Parallel Processing*, pages 573–582,  
810 Sept 2010.
- 811 **31** T. Lundqvist and P. Stenström. Timing anomalies in dynamically scheduled microprocessors. In  
812 *Proceedings of the 20th IEEE Real-Time Systems Symposium*, RTSS '99, pages 12–, Washington,  
813 DC, USA, 1999. IEEE Computer Society.
- 814 **32** M. Lv, N. Guan, J. Reineke, R. Wilhelm, and W. Yi. A survey on static cache analysis for  
815 real-time systems. *Leibniz Transactions on Embedded Systems*, 3(1):05–1, 2016.
- 816 **33** R. Mancuso, R. Dudko, E. Betti, M. Cesati, M. Caccamo, and R. Pellizzoni. Real-time cache  
817 management framework for multi-core architectures. In *Real-Time and Embedded Technology*  
818 *and Applicat. Symp. (RTAS)*. IEEE, 2013.
- 819 **34** F. Martin, M. Alt, R. Wilhelm, and C. Ferdinand. Analysis of loops. In *International Conference*  
820 *on Compiler Construction (CC)*, pages 80–94, London, UK, UK, 1998. Springer-Verlag.
- 821 **35** S. Mittal. A survey of techniques for cache locking. *Transactions on Design Automation of*  
822 *Electronic Systems (TODAES)*, 21(3):49:1–49:24, May 2016.
- 823 **36** F. Mueller. Timing predictions for multi-level caches. In *In ACM SIGPLAN Workshop on*  
824 *Language, Compiler, and Tool Support for Real-Time Systems*, pages 29–36, 1997.
- 825 **37** R. Pellizzoni, E. Betti, S. Bak, G. Yao, J. Criswell, M. Caccamo, and R. Kegley. A predictable  
826 execution model for COTS-based embedded systems. In *Real-Time and Embedded Technology*  
827 *and Applicat. Symp. (RTAS)*, pages 269–279. IEEE, 2011.
- 828 **38** K. Piromsopa and R. J. Enbody. Secure bit: Transparent, hardware buffer-overflow protection.  
829 *Transactions on Dependable and Secure Computing*, 3(4):365–376, 2006.
- 830 **39** I. Puaut. Wcet-centric software-controlled instruction caches for hard real-time systems. In  
831 *Euromicro Conference on Real-Time Systems (ECRTS)*, pages 10–pp. IEEE, 2006.
- 832 **40** I. Puaut and D. Decotigny. Low-complexity algorithms for static cache locking in multitasking  
833 hard real-time systems. In *Real-Time Systems Symposium (RTSS)*, pages 114–123. IEEE,  
834 2002.
- 835 **41** J. Reineke. *Caches in WCET analysis: predictability, competitiveness, sensitivity*. epubli,  
836 2008.
- 837 **42** J. Reineke, D. Grund, C. Berg, and R. Wilhelm. Timing predictability of cache replacement  
838 policies. *Real-Time Syst.*, 37(2):99–122, November 2007.

- 839 **43** A. Sarkar, F. Mueller, and H. Ramaprasad. Static task partitioning for locked caches in  
840 multicore real-time systems. *ACM Trans. Embed. Comput. Syst.*, 14(1):4:1–4:30, Jan. 2015.
- 841 **44** M. R. Soliman and R. Pellizzoni. WCET-Driven dynamic data scratchpad management with  
842 compiler-directed prefetching. In *Euromicro Conference on Real-Time Systems (ECRTS)*,  
843 volume 76, pages 24:1–24:23, 2017.
- 844 **45** C. Song, H. Moon, M. Alam, I. Yun, B. Lee, T. Kim, W. Lee, and Y. Paek. HDFI: hardware-  
845 assisted data-flow isolation. In *Symposium on Security and Privacy (SP)*, pages 1–17. IEEE,  
846 2016.
- 847 **46** V. Suhendra, T. Mitra, A. Roychoudhury, and T. Chen. WCET centric data allocation to  
848 scratchpad memory. In *Real-Time Systems Symposium (RTSS)*, pages 10–pp. IEEE, 2005.
- 849 **47** X. Vera, B. Lisper, and J. Xue. Data cache locking for higher program predictability. *SIG-*  
850 *METRICS Perform. Eval. Rev.*, 31(1):272–282, June 2003.
- 851 **48** X. Vera, B. Lisper, and J. Xue. Data cache locking for tight timing calculations. *ACM Trans.*  
852 *Embed. Comput. Syst.*, 7(1):4:1–4:38, Dec. 2007.
- 853 **49** J. Whitham and N. Audsley. Studying the applicability of the scratchpad memory management  
854 unit. In *Real-Time and Embedded Technology and Applications Symposium (RTAS)*, pages  
855 205–214. IEEE, 2010.
- 856 **50** R. Wilhelm, J. Engblom, A. Ermedahl, N. Holsti, S. Thesing, D. Whalley, G. Bernat, C. Fer-  
857 dinand, R. Heckmann, T. Mitra, F. Mueller, I. Puaut, P. Puschner, J. Staschulat, and  
858 P. Stenström. The worst-case execution-time problem - overview of methods and survey of  
859 tools. *ACM Trans. Embedded Comput. Syst. (TECS)*, 7(3), 2008.
- 860 **51** J. Woodruff, R. N. Watson, D. Chisnall, S. W. Moore, J. Anderson, B. Davis, B. Laurie, P. G.  
861 Neumann, R. Norton, and M. Roe. The cheri capability model: Revisiting risc in an age of  
862 risk. In *International Symposium on Computer Architecture (ISCA)*, 2014.
- 863 **52** G. Yao, R. Pellizzoni, S. Bak, E. Betti, and M. Caccamo. Memory-centric scheduling for  
864 multicore hard real-time systems. *Real-Time Syst.*, 48(6):681–715, 2012.