

# Reverse Engineering of Cache Replacement Policies in Intel Microprocessors and Their Evaluation

Andreas Abel and Jan Reineke  
Department of Computer Science  
Saarland University  
Saarbrücken, Germany  
Email: {abel, reineke}@cs.uni-saarland.de

**Abstract**—Performance modeling techniques need accurate cache models to produce useful estimates. However, properties required for building such models, like the replacement policy, are often not documented. In this paper, using a set of carefully designed microbenchmarks, we reverse engineer a precise model of caches found in recent Intel processors that enables accurate prediction of their cache performance by simulation. In particular, we identify two variants of pseudo-LRU that, unlike previously documented policies, employ randomization. We evaluate their performance and demonstrate that it differs significantly from known pseudo-LRU variants on some benchmarks.

## I. INTRODUCTION

To bridge the increasing latency gap between the processor and main memory, modern microarchitectures employ memory hierarchies with multiple levels of cache memory. These caches are small but fast memories that make use of temporal and spatial locality. Typically, they have a big impact on the execution time of computer programs.

In recent years, different approaches have been proposed to estimate the performance of software systems. This includes analytical modeling, as well as simulation-based and profile-based prediction techniques. All these approaches need sufficiently detailed models of the cache hierarchy in order to produce useful estimates. Similarly, such models are an essential part of worst-case execution time (WCET) analyzers for real-time systems [1]. Furthermore, information on cache properties is also required by self-optimizing software systems, as well as platform-aware compilers.

Unfortunately, documentation of relevant properties at the required level of detail is often not available, or may be misleading. One such property is the cache replacement policy. In this paper, we develop a novel set of microbenchmarks to reverse engineer replacement policies used in recent Intel processors. Based on the results we obtain from these microbenchmarks, we then propose models for the replacement policies that are detailed enough to precisely predict the performance of applications by simulation. Finally, we compare these policies to well-known existing ones by evaluating their performance on the PARSEC benchmark suite.

## II. PROBLEM DESCRIPTION

CPU Caches are structured as follows. They consist of a number of *cache sets*, each of which can store  $k$  memory blocks from the main memory ( $k$  is also called the *associativity* of the cache). A specific memory block can only be stored in

one cache set; this set is determined by a part of the block’s memory address. Upon a cache miss, a so called *replacement policy* must decide which memory block of the corresponding set to replace. One popular strategy is to replace the least-recently used (LRU) block. As the cost of implementing this policy is rather high for larger associativities, processors often use a tree-based approximation to LRU, called *pseudo-LRU* or PLRU (for details we refer to [2]).

In [3] we observed that several Intel Core 2 Duo CPUs appear to use different replacement policies for their L2 caches. According to Intel [4], these CPUs “use some variation of a pseudo LRU replacement algorithm”. However, while we could verify that the Core 2 Duo E6300 (2MB, 8-way set-associative cache) uses the tree-based PLRU policy, the cache behavior of the E6750 (4MB, 16-way set-associative) and the E8400 (6MB, 24-way set-associative) was found to be different from previously documented PLRU variants. The following experiment reveals these differences:

- 1) Clear the cache.
- 2) Access one block in all cache sets.
- 3) Access  $n$  different blocks in all cache sets.
- 4) Access the blocks from 2) again and measure the misses.

Figure 1 shows the result of running this experiment with different values for  $n$  on the CPUs mentioned above. Note that all of those CPUs have 4096 cache sets. For the tree-based PLRU policy, we would expect to get 0 misses if  $n$  is smaller than the associativity, and 4096 misses otherwise, as is the case for the Core 2 Duo E6300. The goal of our work is to develop techniques to build a precise model of the policies used by the other two Core 2 Duo processors.

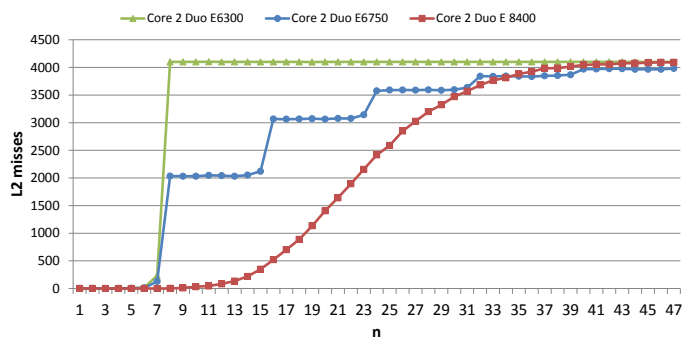


Fig. 1: Experimental analysis of the L2 cache behavior of the Intel Core 2 Duo E6300, E6750, and E8400.

