

Project 3 Report: Cache Replacement Policies on AI Workloads

Tai Duc Nguyen
ECEC 412: Modern Processor Design

November 10, 2019

Abstract

Cache is a solution to the problem of slow memory access in modern computer architecture. From the dawn of the 21st century, many CPU manufacturer, especially Intel and IBM has implemented combinations of advanced dynamic caching techniques, which can be divided into two categories: hardware and software. Some famous hardware techniques, which tackle the issues of maintaining efficient power management while providing fast cache access, are: multi-tiered cache levels and set-associative caching. However, software techniques provide the ability to increase cache hit rate. Some examples include way-predictions methods and cache replacement policies. In this project, different cache replacement policies: Least Recently Used (LRU) and Least Frequently Used (LFU) are evaluated performance wise with different hardware configurations: block size, cache size and number of associative sets.

1 Introduction

In order to achieve the objective of evaluating the performance of different cache replacement policies, a simulation was written in C by Shihao Song and Adarsha Balaji (available here: https://github.com/Shihao-Song/Computer-Architecture-Teaching/tree/master/C621/Cache_Policy). Using the framework established, a mechanism to test all pre-defined hardware configurations at once is written on top. After automating the testing of different configs, the LFU replacement policy is implemented in a similar fashion to the already existed code for the LRU replacement policy.

2 Simulation and Results

From the procedures in Section 1, a sample output from the output is shown in Figure 1 below:

The columns in the sample are respectively:

- Type of replacement policy
- Block size (in bytes)
- Cache size (in KB)
- Number of associative sets
- Number of cache blocks
- Cache hit rate in percentage

```

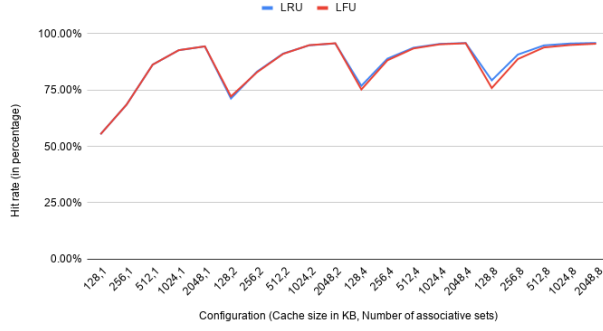
→ Cache_Policy git:(master) X ./Main 541.leela_r_llc.mem_trace
lru, 64, 128, 1, 2048, 38.069094%
lru, 64, 256, 1, 4096, 57.662301%
lru, 64, 512, 1, 8192, 73.476087%
lru, 64, 1024, 1, 16384, 82.539810%
lru, 64, 2048, 1, 32768, 88.332566%
lru, 64, 128, 2, 2048, 39.318178%
lru, 64, 256, 2, 4096, 67.056948%
lru, 64, 512, 2, 8192, 85.927612%
lru, 64, 1024, 2, 16384, 92.713996%
lru, 64, 2048, 2, 32768, 96.241536%
lru, 64, 128, 4, 2048, 42.554451%
lru, 64, 256, 4, 4096, 70.743799%
lru, 64, 512, 4, 8192, 92.363097%
lru, 64, 1024, 4, 16384, 98.220996%
lru, 64, 2048, 4, 32768, 99.455119%
lru, 64, 128, 8, 2048, 42.586477%
lru, 64, 256, 8, 4096, 75.736557%
lru, 64, 512, 8, 8192, 94.443418%
lru, 64, 1024, 8, 16384, 98.867637%
lru, 64, 2048, 8, 32768, 99.597940%
lfu, 64, 128, 1, 2048, 38.069094%
lfu, 64, 256, 1, 4096, 57.662301%
lfu, 64, 512, 1, 8192, 73.476087%
lfu, 64, 1024, 1, 16384, 82.539810%
lfu, 64, 2048, 1, 32768, 88.332566%
lfu, 64, 128, 2, 2048, 41.786363%
lfu, 64, 256, 2, 4096, 67.800881%
lfu, 64, 512, 2, 8192, 85.255685%
lfu, 64, 1024, 2, 16384, 92.605316%
lfu, 64, 2048, 2, 32768, 96.448646%
lfu, 64, 128, 4, 2048, 45.554564%
lfu, 64, 256, 4, 4096, 69.851638%
lfu, 64, 512, 4, 8192, 89.236583%
lfu, 64, 1024, 4, 16384, 96.634206%
lfu, 64, 2048, 4, 32768, 99.132972%
lfu, 64, 128, 8, 2048, 47.165296%
lfu, 64, 256, 8, 4096, 71.792928%
lfu, 64, 512, 8, 8192, 87.736296%
lfu, 64, 1024, 8, 16384, 96.410854%
lfu, 64, 2048, 8, 32768, 99.282395%

```

Figure 1: Sample output of the simulation

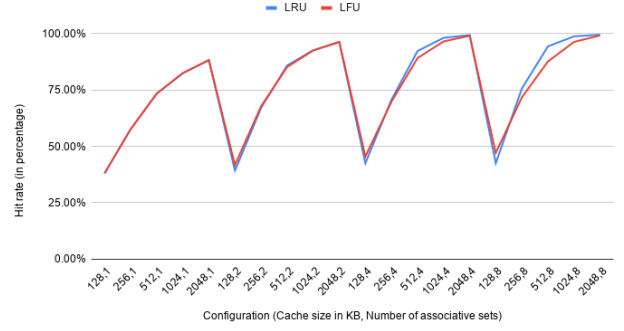
Performing this simulation over all 3 memory traces (each trace includes many Load and Store instructions), the results are shown in graphs (Figure 2 2a, 2b, 2c, and 2d).

Hardware Configuration vs. Hit rate for 531.deepsjeng memory trace



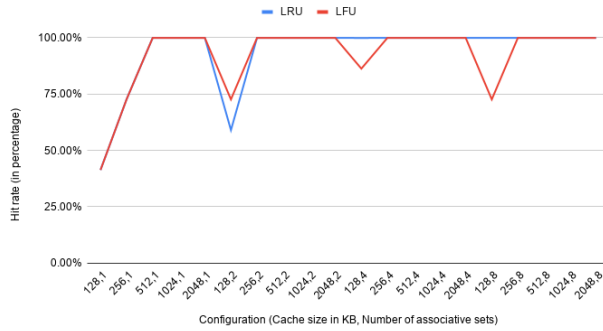
(a) Results for 531.deepsjeng memory trace

Hardware Configuration vs. Hit rate for 541.leela memory trace



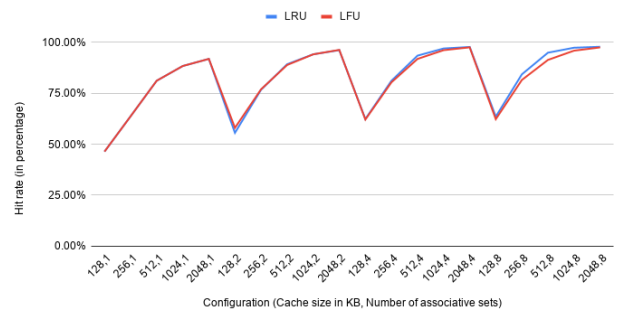
(b) Results for 541.leela memory trace

Hardware Configuration vs. Hit rate for 548.exchange2 memory trace



(c) Results for 548.exchange2 memory trace

Hardware Configuration vs. Hit rate as a weighted average over all 3 traces



(d) Overall results as weighted averages of 3 traces

Figure 2: Simulation results. The x-axis have different hardware configurations with the first number being the cache size in KB and the second number being the number of associative sets. The y-axis is the cache hit rate in percentage.

From all the figures above, the general trend is that the hit rate increase when:

- The cache size increases
- The number of associative cache stores increases

Regarding the 2 different replacement policies, LFU does better than LRU only when the number of associative cache stores is small (1 or 2). LRU has higher hit rate than LFU with larger number of associative cache stores (4 or 8). This behavior is seen most clear from the results for the 548.exchange2 memory trace. It is noted that different memory traces have different number of instructions, hence, the overall performance of each policy is decided as a weighted average and shown in Figure 2d.

3 Discussion

From the results above, the fact that cache hit rate is directly proportional to cache size and number of associative cache stores is consistent with other literature on this subject. Also, from the data, it can also be shown that the increase in performance is not linear, but rather curving out, which means that there will be diminishing returns for large cache size and high number of associative cache due to limitation in hardware area and power consumption.