

# HW3 Cache Optimization

郭宗信, 111550105

**Abstract**—This study focuses on optimizing the L1 data cache (D\$) of the Aquila System-on-Chip (SoC) to improve the efficiency of  $\pi$  computation using Machin's formula. The baseline design, a 4-way set associative cache, was thoroughly analyzed to assess its hit/miss rates, latency, and resource utilization. The optimization process involved adjustments to cache associativity and replacement policies, with a particular emphasis on implementing a random and LRU replacement policy. This approach outperformed traditional FIFO strategy in terms of overall performance. However, reducing associativity was found to have a limited impact on the computation time for  $\pi$ . Future work will explore the integration of prefetching mechanisms and dynamic cache designs to further enhance performance and adaptability for a wider range of workloads.

**Keywords**—Data Cache Optimization; Cache Associativity; Replacement Policies; FPGA-based SoC; Random Replacement Policy;

## I. INTRODUCTION

The performance of FPGA-based System-on-Chip (SoC) architectures is closely tied to the efficiency of their memory hierarchies.

This work focuses on optimizing the L1 data cache for  $\pi$  computation. The baseline 4-way set associative cache was analyzed for hit/miss rates, latency, and resource utilization to identify optimization opportunities. Key contributions include the implementation and evaluation of replacement policies, specifically LRU and Random, with the latter demonstrating superior adaptability and performance compared to traditional FIFO strategy. Experiments on varying associativity levels (2-way, 4-way, and 8-way) revealed that changes in associativity had limited impact on computation time for  $\pi$ , while replacement policies significantly influenced cache performance. The LRU replacement policy optimized hit rates and reduced latency more effectively than other strategies. These findings highlight the potential of replacement policy modifications in enhancing cache performance, establishing a basis for future exploration into advanced optimization techniques such as prefetching and dynamic replacement policies to improve SoC adaptability and efficiency.

## II. SYSTEM DESCRIPTION

### A. Description

The implementation leverages a modular Verilog cache design (dcache.v) and a C program (pi.c) to generate the workload.

This study utilizes the Aquila System-on-Chip (SoC) as the experimental platform, featuring a configurable L1 data cache (D\$). The cache design serves as a critical component for

evaluating the performance of memory-intensive computations, such as  $\pi$  calculation using Machin's formula.

### B. Cache Configuration

The baseline cache design includes the following key characteristics:

- **Cache Size:** 4 KB
- **Associativity:** 4-way set associative
- **Replacement Policy:** FIFO
- **Data Block Size:** 32 bytes

### C. Workload

The workload for this study involves calculating  $\pi$  to 5000 decimal places using Machin's formula. The computation introduces a high volume of repeated memory accesses (read:36,282,622 write: 26,147,435).

## III. IMPLEMENTATION

This part focused on three main aspects of implementation: data analysis, replacement policy modification, and associativity adjustment. Each aspect contributed to evaluating and optimizing the L1 data cache for better performance under specific workloads.

### A. Data Analysis

To facilitate performance analysis, key cache-related variables were monitored using (`*mark_debug = "true"`), enabling observation in the Integrated Logic Analyzer (ILA). Metrics such as hit/miss counts, hit/miss rates, and latency were computed for both read and write operations. The analysis process is summarized in the flowchart (Fig. 1 and 2). By focusing on these metrics, the implementation provided a detailed view of cache performance and identified potential bottlenecks for further optimization.

Fig. 1. cache\_statistics\_flowchart)

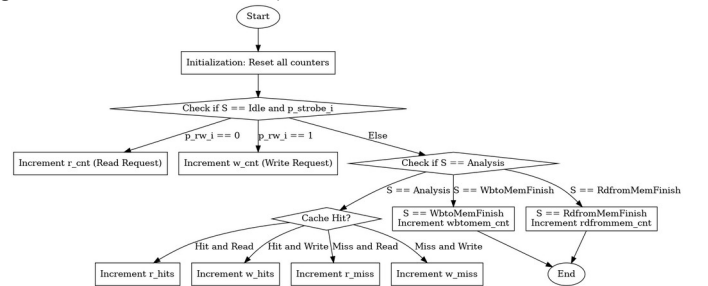


Fig. 2. latency\_tracking\_flowchart



TABLE II. LATENCY OF DIFFERENT SITUATION

Policy	Dirty Miss Latency (Cycles)	Clean Miss Latency (Cycles)
4-way FIFO	51	31
4-way LRU	50	31
4-way Random	50	30

We can find that :

1. The latency overhead of dirty misses significantly impacts cache performance, with a single dirty miss averaging 51 cycles.

2. Both LRU and Random policies slightly reduced average latency compared to FIFO by better managing replacements and minimizing unnecessary write-backs. But this may due to the analysis residual.

### C. Execution Time Comparison

Execution time was measured for  $\pi$  computations to 5,000 decimal places across different configurations. Table III highlights the execution times:

TABLE III. EXECUTION TIME OF DIFFERENT POLICY AND ASSOCIATIVITY

Policy	Execution Time (ms)
4-way FIFO	30,524
4-way LRU	28,983
4-way Random	29,702
8-way FIFO	30,517
2-way FIFO	30,516

LRU achieved the best execution time, improving by approximately 5% over FIFO due to fewer dirty misses and improved replacement decisions.

The Random policy followed closely, demonstrating its ability to adapt to workloads with unpredictable access patterns.

Execution times for 4-way and 8-way configurations were comparable, suggesting diminishing returns for increased

associativity under the tested workload or just the acceptable tolerance.

## V. FUTURE WORK

Future work can focus on exploring alternative strategies for improving the efficiency of  $\pi$  computation and optimizing the D-cache design. For instance, modifying the program logic to reduce memory access intensity or reorganizing data flow in the  $\pi$  computation process could minimize cache dependency. Using alternative methods for calculating  $\pi$ , such as the Gauss-Legendre algorithm or Monte Carlo simulations, may reveal how different computational patterns interact with cache policies and associativity levels.

Also, from a hardware perspective, implementing a write-back buffer could mitigate the latency caused by dirty block writes. By temporarily storing write-back data, the processor could continue executing subsequent instructions without waiting for memory operations to complete. This approach could significantly improve throughput, particularly for write-heavy workloads.

Further more, adopting advanced replacement policies like LFU (Least Frequently Used) or dynamic replacement strategies that adjust based on workload patterns could further enhance cache efficiency. For example, a dynamic policy that transitions between random and LRU depending on access behavior may optimize hit rates while reducing resource overhead.

Last, we can focus on dynamic associativity, where the cache adjusts the number of active ways based on the workload, could also be explored. This method would allow the cache to balance between resource usage and performance, activating more ways for complex workloads and reducing them for simpler tasks.

In summary, these potential advancements in both software and hardware approaches could pave the way for more efficient  $\pi$  computation and optimized cache performance, addressing current limitations while setting a foundation for future exploration.