

UCHICAGO CMSC-22200 COMPUTER ARCHITECTURE, AUTUMN 2021

LAB 4: SIMULATING CACHES

Instructor: Prof. Frederic T. Chong

TAs: Joshua Vízslai, Siddharth Dangwal

Assigned: Wednesday, November 16th, 2022

Due: **11:59 pm, Friday, December 2nd, 2022**(non-graduating students)

Introduction

In this lab, you will extend the 5-stage pipelined ARM machine that you have implemented in previous labs with level 1 (L1) instruction and data caches. We will fully specify the microarchitecture of the caches (see “Microarchitectural Specification” section). Skeleton source code files for this lab can be downloaded from the course webpage. In the `src/` directory, we provide two files (`shell.c` and `shell.h`) that already implement the shell and another six files (`pipe.h`, `pipe.c`, `bp.h`, `bp.c`, `cache.c`, and `cache.h`) that you are allowed to change. You have already implemented the pipelined simulator and branch predictor in `pipe.h`, `pipe.c`, `bp.h`, and `bp.c`. You should implement caches (`cache.c` and `cache.h`) and employ the caches in your pipelined simulator (`pipe.h` and `pipe.c`).

Microarchitectural Specification

Instruction Cache

The *instruction cache* is accessed every cycle in the fetch stage (unless the pipeline is stalled).

Organization. It is a **four-way** set-associative cache that is **8 KB** in size with **32 byte** blocks (this implies that the cache has **64** sets). When accessing the cache, the set index is calculated using bits [10:5] of the PC.

Miss Timing. When the fetch stage *misses* in the instruction cache, the block must be retrieved from main memory. An access to main memory takes **50** cycles. On the 50th cycle, the new block is inserted into the cache. In total, an instruction cache miss stalls the pipeline for 50 cycles. In the 51st cycle, the data is returned to the processor, so the Instruction Fetch stage is completed at the 51st cycle.

Replacement. When a new block is retrieved from main memory, it is inserted into the appropriate set within the instruction cache. If any block within the set are empty, the new block is simply inserted into an empty block. However, if none of the blocks in the set are empty, the new block *replaces* the *least-recently-used* block. For both cases, the new block becomes the *most-recently-used* block.

Control-Flow. While the fetch stage is stalled due to a miss in the instruction cache, a control-flow instruction further down the pipeline may *redirect* the PC. As a result, the pending miss may turn out to be unnecessary: it is retrieving the wrong block from main memory. In this case, the pending miss is *anceled*: the block that is eventually returned by main memory is *not* inserted into the cache. (A test input of this case is provided in `inputs/cancel_req.s`.) Finally, note that a redirection that accesses the *same* block as a pending miss does *not* cancel the pending miss.

Data Cache

The *data cache* is accessed whenever a load or store instruction is in the memory stage.

Organization. It is an **eight-way** set-associative cache that is **64 KB** in size with **32 byte** blocks (this implies that the cache has **256** sets). When accessing the cache, the set index is calculated using bits [12:5] of the data address that is being loaded/stored. The L1 data cache is a write-through, allocate-on-write cache.

Miss Timing & Replacement. Miss timing and replacement of the data cache are identical to the instruction cache.

Handling Load/Store. Both load and store misses stall the pipeline for 50 cycles. They both retrieve a new block from main memory and insert it into the cache. The Memory stage of the load/store instruction

will take 51 cycles when miss, and 1 cycle when hit. For write miss, memory should be updated in the 51st cycle.

Assumptions & Initial States

- Assume that both caches are initially empty.
- Assume that a program that runs on the processor *never* modifies its own code (referred to as self-modifying code), and that a given block *cannot* reside in both caches.
- Assume that if we hit either cache, the data is returned immediately (in the Fetch stage for I cache, and Memory stage for D cache).
- Initialize all data structures related to the caches (e.g., cache blocks, tag/valid-bit arrays, etc.) to 0.
- As stated in lab1, for all load/store instructions, the address accessed is aligned to the size of the data element being accessed. This implies that every memory access will involve one and only one cache line.

Testing Your Code

Use the simulator that you developed in lab 3 as a starting point to implement the pipelined simulator with L1 caches. Your simulator in lab3 can be used as a reference to verify that your simulator (with cache) functions correctly.

We are also providing you with a reference simulator and a set of test files, which you can also use to verify the functionality of your cache implementation. Note that the reference simulator provides hit and miss information for every memory access.

Handin

If you are working in groups of 2, only one copy of the code should be submitted. **Please make sure you write down both names and CNET IDs at the top of the pipe.c and cache.c files.**

You should electronically hand in your code (all files in the `src/` directory) within a `lab4` folder of your own SVN repository. Note that we will check out your code automatically and look for the following folder:

`Your_SVN_ROOT_Directory/lab4/src/`

So, if your code is not in the correct directory or the name of the folders are not the same as above (e.g. “Lab4” instead of “lab4”), you will see delay in grading your lab. Also, please do not insert large files in that folder.

Contact the TA if your handin SVN repository does not exist. Your code should be readable and well-documented.

You should stick with the same partner as Lab3.

Graduating students cannot use late days for this lab.

Key Reminders

1. You should write more input files in order to be confident that your simulator is correct.
2. We will test your simulator with many input programs (some provided with the problem, some not). We will check your code cycle by cycle for each test input (using `rdump` and `mdump`) and your output (the contents of the `dumpsim` file) must match the reference simulator output EXACTLY.
3. In this lab, you need to implement the correct number of **retired instructions** in your `rdump` output.