# Homework #7 – Caching and Virtual Memory

## Due Wednesday, Dec 4 at 5:00pm

Directions:

- For short-answer questions, submit your answers in PDF format to the GradeScope assignment "Homework 7 written".
- For the programming question, submit your source file using the filename specified in the question to the GradeScope assignment "Homework 7 code".
  - Programs that show good faith effort will receive a minimum of 25% credit.
- **You must do all work individually, and you must submit your work electronically via GradeScope.**
  - All submitted code will be tested for suspicious similarities to other code, and the test will uncover cheating, even if it is "hidden" (by reordering code, by renaming variables, etc.).

## Q1. Cache policies

[5 points] Why are write-back caches usually also write-allocate?

## Q2. Cache performance

[5] Your L1 data cache has an access latency of 2ns, and your L2 cache has an access latency of 8ns. Assume that 80% of your L1 accesses are hits, and assume that 100% of your L2 accesses are hits. What is the average memory latency as seen by the processor core? **You must show your work to receive full credit.**

## Q3. Virtual memory layout

[20] You have a 64-bit machine and you bought 32GB of physical memory. Pages are 256KB. For all calculations, **you must show your work to receive full credit**.

(a) [1] How many virtual pages do you have per process?

(b) [1] How many bits are needed for the Virtual Page Number (VPN)? (Hint: use part (a))

(c) [1] How many physical pages do you have?

(d) [1] How many bits are needed for the Physical Page Number (PPN)? (Hint: use part (c))

(e) [1] How big in _bytes_ does a page table entry (PTE) need to be to hold a single PPN plus a valid bit? Round up to the nearest power-of-2 number of bytes.

(f) [1] How big would a flat page table be for a single process, assuming PTEs are the size computed in part (e)?

(g) [10] Why does the answer above suggest that a "flat page table" isn't going to work for a 64-bit system like this? Research the concept of a _multi-level page table_, and briefly define it here. Why could such a data structure be much smaller than a flat page table?

(h) [4] Does a TLB miss always lead to a page fault? Why or why not?

# Q4. Cache simulator program

[70] In Java or C, write a simulator of a single-level cache and the memory underneath it.[1]

The simulator, called `cachesim`, takes the following input parameters on the command line: name of the file holding the loads and stores, cache size (not including tags or valid bits) in kB, associativity, and the block size in bytes. The replacement policy is always FIFO. FIFO stands for "first in, first out" and it means exactly that. Whatever block is first brought into a set will be the first one evicted. Consider a 4-way set-associative cache in which blocks A, B, C, D, E, and F all map to a given set. Assume these six blocks are brought into the cache in that order. Thus, when E is brought in, A is evicted to make room. When F is brought in, B is evicted. (It does not matter if A or B was accesses more recently.)

For example, "`cachesim tracefile 1024 4 32`" should simulate a cache that is 1024kB (=1MB), 4-way set-associative, has 32-byte blocks, and uses FIFO replacement. This cache will be processing the loads and stores in the file called tracefile.

**Important Assumptions:** Addresses are 24-bits (3 bytes), and thus addresses range from 0 to $2^{24}$-1 (i.e., there is 16MB of address space). The machine is byte-addressed and big-endian. The cache size, associativity, block size, and access size will all be powers of 2. Cache size will be no larger than 2MB, block size will be no larger than 1024 bytes, and no access will be larger than the block size. No cache access will span multiple blocks (i.e., each cache access fits within a single block).

All cache blocks are initially invalid. All cache misses are satisfied by the main memory (and you must track the values written through to memory in case they are subsequently loaded). If a block has never been written before, then its value in main memory is zero. The cache is **write-back** (and thus each block needs a "dirty" bit) and **write-allocate**. This means your program will need to store both the state of cache and the entire content of simulated memory; the memory part can be represented as an array of 16M bytes.

If you have any known bugs, please include those in a README file to help the grader give partial credit.

## Calling syntax

The program will be called `cachesim`, and will have the following calling syntax:

```
./cachesim <trace-file> <cache-size-kB> <associativity> <block-size>
```

Arguments:

- `<trace-file>`: Filename of the memory access trace file.
- `<cache-size-kB>`: Total capacity of the cache, <u>kilobytes (kB)</u>. A power of two between 1 and 2048.
- `<associativity>`: The set associativity of the cache, AKA the number of ways. A power of two.
- `<block-size>`: The size of the cache blocks, in <u>bytes</u>. A power of two between 2 and 1024.

All numeric arguments are in base-10 format.

---

[1] Note: a cache can be simulated as a 2D array, and dealing with 2D arrays can be trickier in C than in Java.

## Trace file format

The trace file will be in the following format. There will be some number of lines. Each line will specify a single load or store, the 24-bit address that is being accessed (in base-16), the size of the access in bytes, and the value to be written if the access is a store (in base-16). For example:

```
store 0xd53170 4 7d2f13ac
load  0xd53172 1
store 0xd53170 2 f0b1
store 0x1a25bb 2 c77a
load  0xd53170 4
load  0x12 2
store 0x23 8 d687eb9f1bc687ec
```

As can be seen, leading 0 bits will not be in addresses in the trace file. Also, as viewed in the store commands, values following the access size in bytes will be the correct size. Accesses will be no larger than 8 bytes at a time. Because all parts of the file are whitespace-delimited tokens, `fscanf` will be your friend.

## File reading helpers

To help with reading the trace file, C and Java versions of the following three functions are provided in the starter GitLab repo:

-   `void traceInit(char* filename)`: Opens the trace file, given its name
-   `bool traceFinished()`: True if all accesses have already been read with `traceNextAccess`
-   `cacheAccess traceNextAccess()`: Gets the next access from the trace file, returning a C struct (or Java class) with the following fields:
    -   `bool isStore`: true for stores, false for loads
    -   `int address`: the address being accessed as an integer
    -   `int accessSize`: the size of the access in bytes
    -   `byte[8] data`: for stores, this contains the data being stored

After forking and cloning the GitLab repo, these functions will be available in the `cachesim.c` and `cachesim.java` files. You can use them from your `main` method to help process the trace file.

## Program output

Your simulator must produce the following output. For every access, it must print out what kind of access it is (load or store), what address it's accessing (in base-16), and whether it is a hit or a miss. For each load, it must print out the value that is loaded (possibly after satisfying the miss from memory). If an access requires a valid block to first be evicted, your simulator must print out the address of the evicted block (i.e., the address of byte 0 in that block) and whether it is dirty or clean (i.e., not dirty).

The output format must be as follows and may not be graded if format is ignored. Below is output for the example input file shown above with a 1MB 4-way cache with 32-byte blocks. (Thanks for reading thoroughly; put a picture of a basketball in your Q1 answer for extra credit.) This means the cache has

1MB/32 = 32768 frames spread among 4 ways per set, so 32768/4 = 8192 sets. Each line of output is annotated with an explanation:

```
$ ./cachesim traces/example.txt 1024 4 32
store 0xd53170 miss              (First seen, and write-allocate means we cache it now)
load 0xd53172 hit 13             (We just cached this guy, so load hit)
store 0xd53170 hit               (Store hit)
store 0x1a25bb miss              (First seen, and write-allocate means we cache it now)
load 0xd53170 hit f0b113ac       (Another hit, and see how the data reflects the stores above)
load 0x12 miss 0000              (First seen, so we cache it now; data is the starting 00's)
store 0x23 miss                  (First seen also, as this is in a different cache block than 0x12)
```

In a longer example with cache replacements, a replacement of a dirty block would be output like this:

```
replacement 0x96 dirty
```

And a replacement of a clean block would be output like this:

```
replacement 0x128 clean
```

The "0x" before each address is required; leading 0s are not printed.

Always print out the exact number of hex digits corresponding with the number of bytes missed. Memory defaults to all zeroes on boot.

## Restriction

In this assignment, **you may NOT use the modulus (%) operator** and **you may NOT include math.h**. You must use bitwise operations to decompose the address.  Penalty: 50% off overall score.

Additionally, your program should exit with a status of 0 (EXIT_SUCCESS). Penalty: 25% off score.

If you program in C, you do NOT have to worry about memory leaks on this assignment. However, valgrind is still a useful tool for detecting misuse of memory, and may help you find hidden bugs!

## Building and testing

You can simply build your program as per usual with g++ (`javac` for Java):

```
g++ -g -o cachesim cachesim.c
```

A number of input files are in the `traces` subdirectory; these are described by `traces/INFO.txt`. As with prior assignments, a suite of tests is provided in the `tests` subdirectory and an automated testing tool, `hwtest.py`, has been provided to automate testing.

If you're using Java, there's a `settings-java.json` file in the `tests` directory which configures the tester to run Java instead of C. To switch to that settings file, move it with:

```
mv tests/settings-java.json tests/settings.json
```

## Tips relevant to Q4/Q5

- To manipulate bit fields, use bitwise operators (`&`, `|`, `~`), shifts (`<<`, `>>`) and masks.
- You can compute $2^N$ with the expression: `(1<<N)`
- You can get a bit string of N ones with the expression: `((1<<N)-1)`
- Here's a simple implementation of base-2 log using only integer math, that way you don't have to mess with the math library:

```
int log2(int n) {
    int r=0;
    while (n>>=1) r++;
    return r;
}
```

- Parsing/printing tips for `cachesim`:
    - Use the provided functions to read the trace file passed in the command line
    - You can print a zero-padded two-digit hex representation of a byte with the "`%02hhx`" `printf` specifier.
    - See the documentation for the provided functions and the manpages for `printf` for more information.