

ENCM 369 Winter 2023 Lab 10 for the Week of March 27

Steve Norman
Department of Electrical & Software Engineering
University of Calgary

March 2023

Administrative details

You may work with a partner on this assignment

If you choose to work with a partner, please make sure that both partners fully understand all parts of your assignment submission, and please follow the instructions regarding submission of your PDF document.

Partners must be in the same lab section. The reason for this rule to keep teaching assistant workloads balanced and to make it as easy as possible for teaching assistants to record marks.

Due Dates

The Due Date for this assignment is 6:00pm Thursday, April 6.

The Late Due Date is 6:00pm Tuesday, April 11.

The penalty for handing in an assignment after the Due Date but before the Late Due Date is 3 marks. In other words, X/Y becomes $(X-3)/Y$ if the assignment is late. There will be no credit for assignments turned in after the Late Due Date; they will not be marked.

How to package and hand in your assignments

You must submit your work as a *single PDF file* to the D2L dropbox that will be set up for this assignment. The dropbox will be configured to accept only file per student, but you may replace that file as many times as you like before the due date.

See the Lab 1 instructions for more information about preparing and uploading your PDF file.

Important update for those working with a partner: Please submit only one PDF file for both students. On the cover page, put lab section, name and ICID information in this format:

Group Submission for [lab section]

Submitted by: [submitter's name]

UCID: [submitter's UCID]

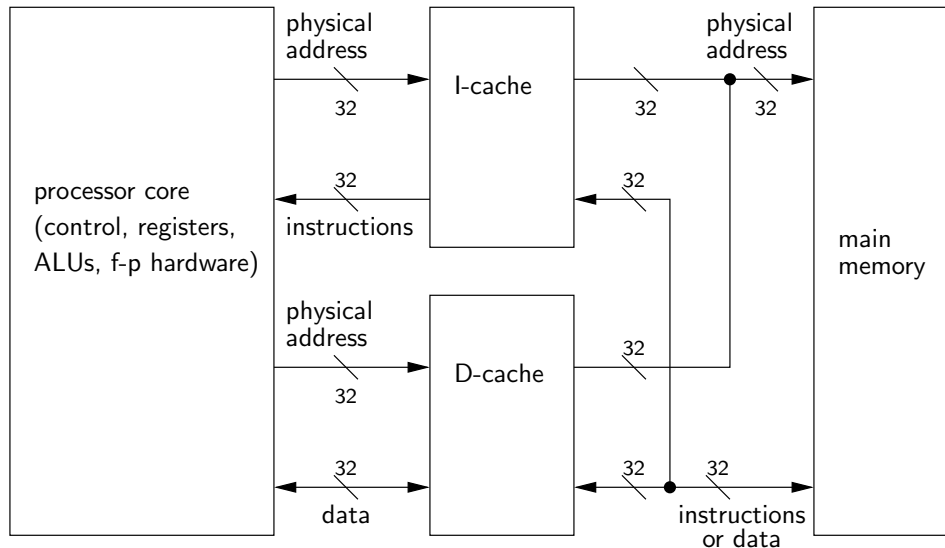
Partner: [partner's name]

UCID: [partner's UCID]

Marking scheme

A	5 marks
B	5 marks
C	4 marks
D	3 marks
TOTAL	17 marks

Figure 1: Computer with one level of cache, and no address translation.



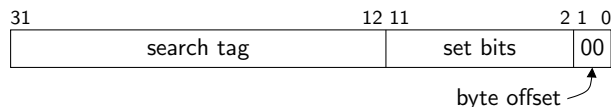
Exercise A: Tracing behaviour of an I-cache

Read This First

In learning about caches, it is useful to trace through all of the copying of bit patterns that occurs in a sequence of memory accesses.

What to Do

The table in Figure 3 shows some of the main memory contents for a program running on a computer like the one shown in Figure 1, running the RISC-V instruction set. Note that the first sequence of instructions is a complete procedure, but the second sequence is only part of a procedure. The second sequence includes a loop; within that loop there is a call to the procedure comprised of the first sequence of instructions. You will be asked to trace the interaction between the I-cache and the main memory starting with `PC = 0x0040_5b50`, at the moment in time just before the `beq` instruction is fetched.

Figure 2: How instruction addresses are split for access to the cache of Exercise A.**Figure 3:** Small fragments of main memory contents for Exercise A.

address	instruction at address	disassembly of instruction
0x0040_4b58	0x0004_2283	P1: lw t0, (a0)
0x0040_4b5c	0xfc02_8313	addi t1, t0, -64
0x0040_4b60	0x0065_2023	sw t1, (a0)
0x0040_4b64	0x0000_8067	jr ra
⋮	⋮	⋮
0x0040_5b50	0x015a_0a63	beq s4, s5, L2
0x0040_5b54	0x000a_0533	L1: add a0, s4, zero
0x0040_5b58	0x004a_0a13	jal P1
0x0040_5b5c	0x800f_f0ef	addi s4, s4, 4
0x0040_5b60	0xff5a_1ae3	bne s4, s5, L1
0x0040_5b64	0x0000_0a33	L2: add s4, zero, zero

Figure 4: Part of the initial state of the I-cache for Exercise A. (Only sets 724–729 in the cache are shown.)

set	valid	tag	instruction
724	1	0x00405	0x015a_0a63
725	0	0x00000	0x0000_0000
726	1	0x00404	0x0004_2283
727	1	0x00404	0xfc02_8313
728	1	0x00404	0x0065_2023
729	1	0x00404	0x0000_8067

The I-cache for this computer is direct-mapped with 1024 sets, and the block in each set contains one instruction. This structure exactly matches an example that has been presented in a lecture. For this cache, main memory addresses are split as shown in Figure 2.

At the moment in time mentioned above, the state of sets 724–729 of the cache is shown in Figure 4. Use 0x1000_1230 as the initial value for `s4` and 0x1000_1238 as the initial value for `s5`. Trace all the instruction fetches until after the instruction at address 0x0040_5b64 has been fetched. Record your answer in tabular form, using this trace of the first two instructions as a model:

address	tag	set	action
0x0040_5b50	0x00405	724	I-cache hit—no I-cache update
0x0040_5b54	0x00405	725	I-cache miss—instruction 0x000a_0533 is copied into instruction field in set 725, V-bit in that set is changed to 1, tag to 0x00405

Hints: (1) Including the two instruction fetches given as examples, there will be a

total of 18 instruction fetches. (2) There is a useful C program in `encm369w23lab10/exA`

What to Include in Your PDF Submission

Include your completed table.

Exercise B: Analysis of direct-mapped caches

Read This First

The base-two logarithm

The *base-two logarithm* is a simple and useful concept that has many useful applications in computer systems, one of which is describing dimensions in caches. The \log_2 function is simply the inverse of the function $f(x) = 2^x$, in the same way that the \ln function is the inverse of $f(x) = e^x$ and the \log_{10} function (often written as simply \log) is the inverse of $f(x) = 10^x$. For example,

$$\log_2 1 = 0, \log_2 4 = 2, \log_2 32 = 5, \text{ and } \log_2 65536 = 16.$$

Dimensions within direct-mapped caches

Direct-mapped lookup is the simplest practical way to organize a data or instruction cache. It is illustrated in Figures 8.7 and 8.12 in the course textbook. The general structure of all direct-mapped caches is shown in Figure 5.

To search for an instruction word or data word in a direct-mapped cache, the main memory address of that word is broken into three or four parts:

- *tag*: used to distinguish a block address among a group of block addresses that all generate the same set bits;
- *set bits*: used to select one set among the S sets in the cache;
- *block offset*: used to select a single word within a multi-word block, so not necessary if the cache has one-word blocks;
- *byte offset*: can be assumed to be zero when the cache access is to an entire data or instruction word, but would be important for access to individual bytes in instructions such as RISC-V `lb`, `lbu`, and `sb`.

How wide are each of these fields? Let's go right-to-left within an address. First,

$$\text{byte offset width} = \log_2(\text{number of bytes per word}).$$

So with 4-byte words, as in textbook and lecture examples, the byte offset width is $\log_2 4 = 2$, but with 8-byte words, the byte offset width would be $\log_2 8 = 3$.

Next,

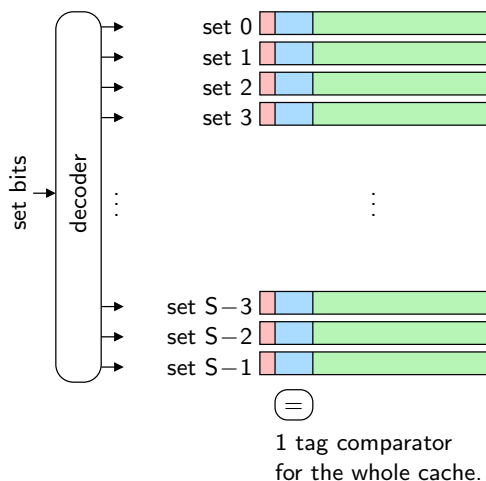
$$\text{block offset width} = \log_2(\text{number of words per block}).$$

For example, in textbook Figure 8.12, there are 4 words per block, so the block offset is $\log_2 4 = 2$ bits wide. Note that that gives you four different bit patterns to choose one of the four words within a block: 00, 01, 10, 11. Note also that $\log_2 1 = 0$, consistent with the idea that if the block size is one word, no bits from the address should be used as a block offset—that is what you see in textbook Figure 8.7.

Moving on, let's let S stand for the number of sets within the cache. Then

$$\text{number of set bits} = \log_2 S.$$

Figure 5: General organization of a direct-mapped cache. Wiring and some logic components have been left out to reduce clutter.



Key to colouring of storage cells

- status bit(s): 1 valid bit, plus, possibly, another bit to help with writes
- tag bits
- block of data or instruction words

Note: Relative sizes are *not* to scale. A block might be as large as 64 bytes (512 bits), which is difficult to describe graphically in proportion to a single status bit.

I hope you can see a general pattern here: If X is a power of two, then $\log_2 X$ bits are needed to select one of X things.

Finally, the tag is everything in an address that hasn't already been used:

tag width =

address width – number of set bits – block offset width – byte offset width

The *capacity* of a cache is usually defined as the maximum number of bytes of data or instructions that a cache can hold. Note that that definition *excludes* storage of status bits and tag bits. Let's let Bpl ("bytes per line") stand for the size of a block. (*Cache line* is a synonym for *cache block*, and unlike "block", "line" does not confusingly start with the letter b.) From Figure 5 it should be clear that for a direct-mapped cache

$$C = S \times \text{Bpl},$$

and if we want to think of block size measured in words instead of bytes,

$$C = S \times \text{words per block} \times \text{bytes per word}.$$

Example calculations

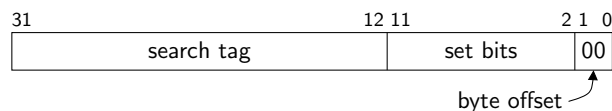
1. Suppose it has been decided that a direct-mapped cache should have 1024 entries, each with one 32-bit data word. How should addresses be split into parts? And what will the capacity of this cache be?

Byte offset: A 32-bit word is 4 bytes, so the width is $\log_2 4 = 2$.

Block offset: There is no block offset, because there is only one word per block.

Set bits: $S = 1024$, so we need $\log_2 1024 = 10$ set bits.

Tag: The tag is all the bits to the left of the set bits. So addresses should be split this way:



Capacity is

$$\begin{aligned} C &= 1024 \text{ blocks} \times 1 \frac{\text{word}}{\text{block}} \times 4 \frac{\text{bytes}}{\text{word}} \\ &= 4096 \text{ bytes} \\ &= 4 \text{ KiB.} \end{aligned}$$

Remark: This has been a review of the organization of the cache used in Exercise A of this lab.

2. Suppose it has been decided to build a direct-mapped cache with a (very tiny) capacity of 32 bytes, in which each block holds 4 32-bit words. What is S , the number of sets? And how should addresses be split into parts?

To find S , solve for it in this equation:

$$32 \text{ bytes} = 2^5 \text{ bytes} = S \times 2^2 \frac{\text{words}}{\text{block}} \times 2^2 \frac{\text{bytes}}{\text{word}}$$

That gives

$$S = \frac{2^5 \text{ bytes}}{2^2 \frac{\text{words}}{\text{block}} \times 2^2 \frac{\text{bytes}}{\text{word}}} = 2^{5-2-2} \text{ blocks} = 2^1 \text{ blocks} = 2 \text{ blocks.}$$

(The equation gives S as a number of blocks rather than a number of sets, but that's fine, because in a direct-mapped cache there is one block per set.)

The address split is: byte offset, $\log_2 4 = 2$ bits; block offset, $\log_2 4 = 2$ bits; set bits, $\log_2 2 = 1$ bit; tag, $32 - 1 - 2 - 2 = 27$ bits. Graphically, that's:



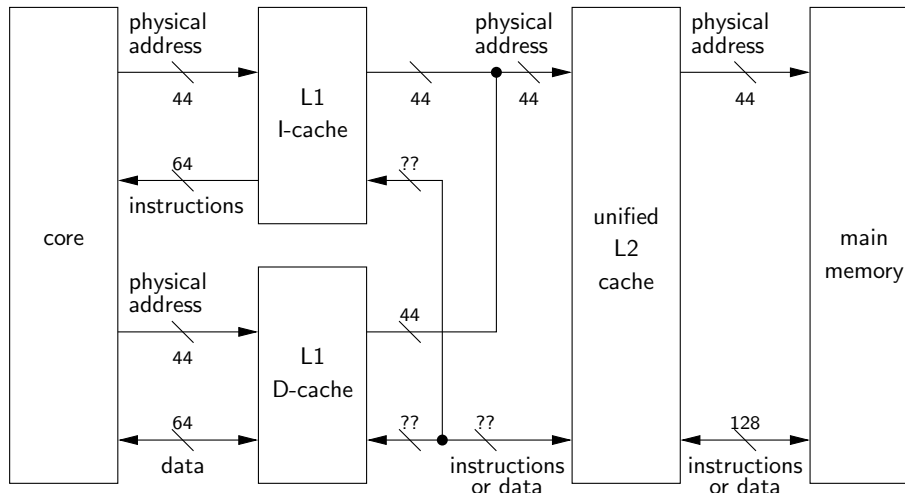
Remark: This example has been a review of the organization of the cache in textbook Figure 8.12.

What to Do

Write well-organized solutions to the following problems.

3. Suppose the specification for the cache of textbook Figure 8.12 is changed. The block size is now supposed to be 8 words, and the capacity has been increased to a much more practical size of 16 KiB.
- What is S , the number of sets?
 - How should addresses be split into parts? Show this with a diagram like the previously given examples—include bit numbers marking the boundaries between parts.

Figure 6: Simplified view of memory organization of a recent single-core 64-bit system. (The major simplification here is the omission of address translation.) The core can fetch two 32-bit instructions at once from the Level 1 I-cache. The width of the data/instruction bus between the Level 1 and Level 2 caches is unspecified but should be much wider than 64 bits to support speedy transfers of large blocks.



- (c) The cache will be built using SRAM cells, one SRAM cell for every V-bit, every bit within a tag, and every bit within a block of data or instruction words. How many SRAM cells are needed for the whole cache? Show your work carefully.
4. The term “64-bit processor” usually describes a processor in which general-purpose registers are 64 bits wide, and memory addresses *within the processor core and in pointer variables* are managed as 64-bit patterns. However, 2^{64} bytes of DRAM is enormously larger than the amount of DRAM that can practically be connected to a single processor chip, so a typical 64-bit design might use only the least significant 44 bits of an address to access caches and main memory. This is illustrated in Figure 6.

Do the following calculations for the computer of Figure 6. We’ll assume direct-mapped design for all three caches, and we’ll consider the word size to be 64 bits.

- The capacity of the L1 D-cache is 32 KiB. The block size is 64 bytes. Draw a diagram to show how a 44-bit address input to the cache would be split into these fields: tag, set bits, block offset, and byte offset. Indicate exactly how wide each field is.
- Repeat part (a) for the L2 cache, which has a capacity of 4 MiB. The block size is again 64 bytes.
- Use your results from (b) to determine how many one-bit SRAM cells are needed for all the V-bits, tags, and data/instruction blocks in the L2 cache.

What to Include in Your PDF Submission

Include solutions to the problems given in “What to Do”.

Exercise C: Simulating a Cache

Read This First

In this exercise you will work with a C program to simulate some aspects of the behaviour of data caches. For two different sequences of memory accesses you will determine which accesses are cache hits and which accesses are cache misses. (This is relatively easy to do; a complete simulation that modeled write-through or write-back to main memory would be a much more complex problem.)

The sequences of memory accesses were generated by determining the data memory accesses that would be made in sorting an array of 3000 integers using two different sort algorithms on a machine similar to a 32-bit RISC-V computer. Memory is stored in 32-bit words, and byte addresses are 32-bits in length. All accesses to data memory by the sort procedures are loads and stores of *words*, at addresses that are multiples of four.

These sequences are available in text files. Here are the first ten lines of one of the files:

```
r 804e6ac
r 804fe1c
w 804e6ac
r 804e6a8
r 804fe14
r 804fe18
w 804e6a8
w 804fe18
r 804e6a4
r 804fe0c
```

`r` means read (load) and `w` means write (store). The addresses are in hexadecimal notation, even though there is no leading `0x`.

Note that the actual data values are not included in the file, just the addresses used to access data. It turns out that to count cache hits and misses, the sequence of addresses used is all that really matters.

What to Do, Part I

Inspect the files in `encm369w23lab10/exC`

Read the C source file `sim1.c` carefully. The program does a simulation of a data cache with 1024 words in 1-word blocks. Build an executable and run it with both data files, following the instructions in a comment near the top of the source file.

Here is some information about two memory access traces:

- In `heapsort_trace.txt` all data memory accesses are to the array elements. Each of the 3000 elements is read at least once, and is likely to be read and written many more times.
- In `mergesort_trace.txt` the data memory accesses are to all of the array elements, to a 1500-word array of temporary storage needed by the mergesort algorithm, and to 72 words of stack used to manage a sequence of recursive function calls. So the total number of different words accessed is 4572.

In both cases the 1024-word cache is much too small to hold all of the different data memory words being accessed. If you naïvely guess that access to memory words is truly random, you would expect very high miss rates, such as $(3000 - 1024)/3000 = 65.9\%$ or worse for the heapsort run, and $(4572 - 1024)/4572 = 77.6\%$ or worse for

the mergesort run. However, you should see much lower miss rates due to *locality of reference* in the memory accesses.

What to Do, Part II

Let's examine the effect of changing the block size of the cache of Part I, while maintaining capacity.

Make a copy of `sim1.c` called `sim2.c`, and edit it to simulate a direct-mapped cache with 128 eight-word blocks. You will not have to edit many lines of the C code, but to do it correctly you will have to do some calculations like the ones in Exercise B, then think carefully about how to isolate the correct set and tag bits.

Run the program using the two given input files. By screenshotting, copy-pasting text, or whatever other method is convenient, make records of your program source code and program outputs that you can include in your lab report.

Also, answer this question:

Compare results obtained in this part with results from Part I. Do they suggest that there is significant *spatial* locality of reference in the memory accesses done by the heapsort and mergesort algorithms? Give a brief explanation.

What to Do, Part III

Now let's consider a direct-mapped cache with the same eight-word block size as in Part II, but with the cache capacity doubled to 8 KiB.

Make another copy of `sim1.c`; call this one `sim3.c`. Edit it to simulate a direct-mapped cache with the appropriate number of eight-word blocks.

Run the program using the two given input files. As in Part II, make records of your source code and program outputs suitable for inclusion in your lab report.

Final Note

Cache-friendliness (a tendency towards low miss rates) is an important factor in performance of algorithms on modern computer hardware.

But don't use this exercise to draw any firm conclusions about which of the two sort algorithms is more cache-friendly. (My own suspicion is that mergesort may often be significantly more cache-friendly than heapsort.) The caches being simulated are unrealistically small, and using one run of each algorithm for a single array hardly constitutes enough input data for a good experiment.

Researchers doing simulation experiments to compare cache designs will use traces of trillions of memory accesses from many different programs in order to ensure that performance is being measured for many different patterns of memory access.

What to Include in Your PDF Submission

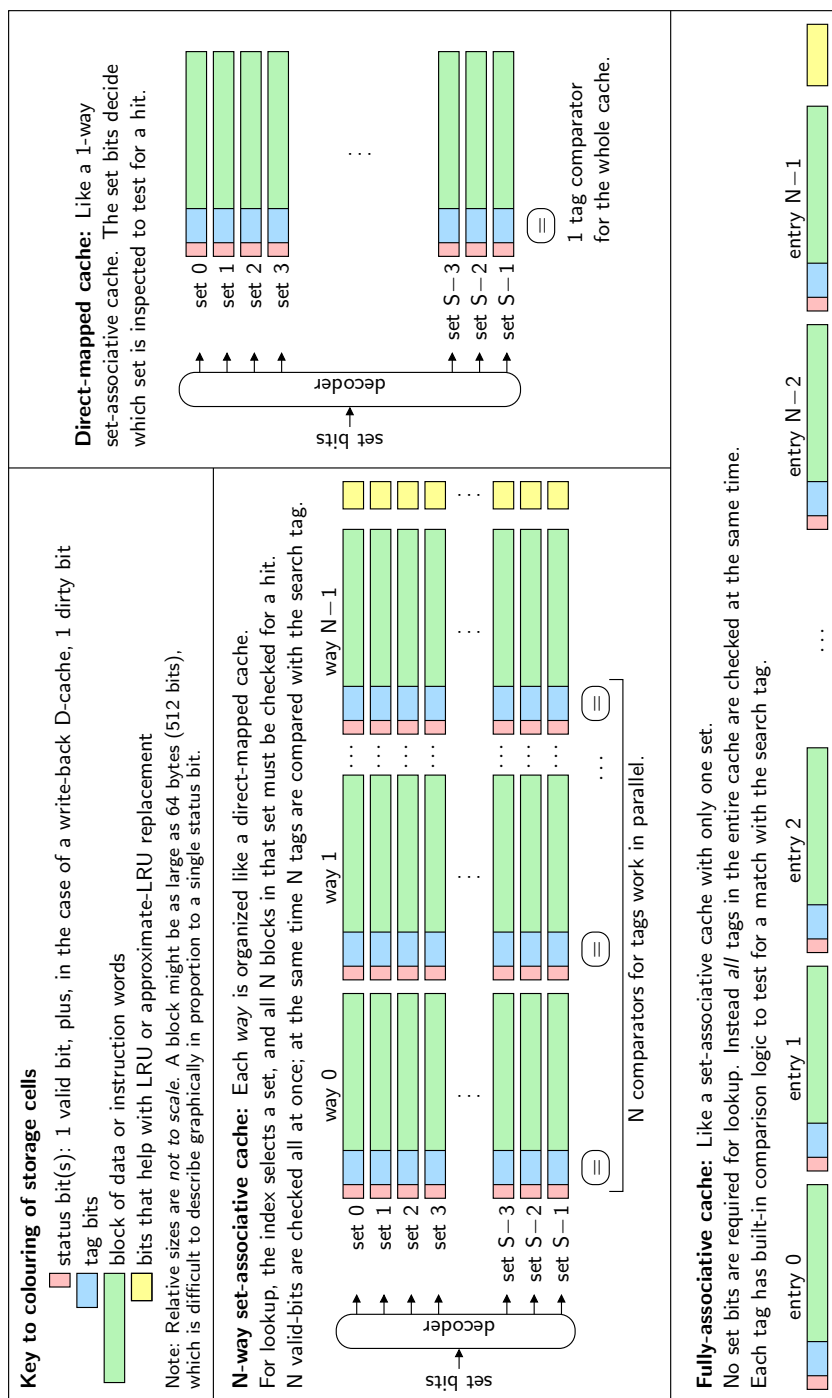
Include listings of code and program output from Parts II and III, and your answer to the question in Part II. Please label the listings clearly so your TA's don't have to guess which listing is from which part.

Exercise D: Analysis of set-associative caches

Read This First

As explained in lectures, set-associative caches provide a good (but not perfect) solution to the problem of set-bit conflicts in direct-mapped caches. The difference

Figure 7: Graphical comparison of set-associative caches, direct-mapped caches, and fully-associative caches. *After reading this caption you should probably use the “Rotate Right” tool in your PDF reader to view the diagram.* The term “search tag” refers to the tag extracted from the address of the instruction or data item the processor core is seeking. The diagrams are intended to show conceptual structure, and leave out wires and some important logic elements (e.g., multiplexers). Layout in the diagram may or may not reflect actual physical layout within integrated circuits.



between direct-mapped and set-associative structure is illustrated in Figure 7.

Just as for a direct-mapped cache, in a set-associative cache an address must be split into parts: tag, set bits, block offset (if there are multi-word blocks), and byte offset. The formulas for the widths of these address parts are exactly the same as given for direct-mapped caches in Exercise A of this lab.

However, the equations for the *capacity* of a set-associative cache are slightly different from those for the capacity of a direct-mapped cache. If C is capacity in bytes, S is the number of sets, N is the number of ways, and Bpl (“bytes per line”) is the size in bytes of a block, then

$$C = S \times N \times \text{Bpl}.$$

To use a block size given in words instead of bytes:

$$C = S \times N \times \text{words per block} \times \text{bytes per word}.$$

For units to make sense in these equations, it turns that appropriate units for N are *blocks per set*.

Here are some example calculations:

- What is the capacity of the cache of textbook Figure 8.9?

Obtaining all the dimensions we need from Figure 8.9,

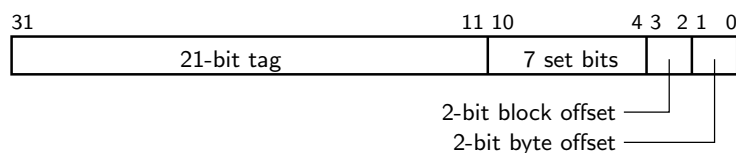
$$\begin{aligned} C &= S \times N \times \text{words per block} \times \text{bytes per word} \\ &= 4 \text{ sets} \times \frac{2 \text{ blocks}}{\text{set}} \times \frac{1 \text{ word}}{\text{block}} \times \frac{4 \text{ bytes}}{\text{word}} \\ &= 32 \text{ bytes.} \end{aligned}$$

- Suppose we modify the design of the cache of Exercise A this lab, so that it is two-way set-associative instead of direct-mapped, and has four-word blocks instead of one-word blocks. The capacity will remain 1024 words, that is 4096 = 2^{12} bytes. How will addresses be split into parts to access this new design?

We need to know S to determine the number of set bits.

$$\begin{aligned} S &= \frac{C}{N \times \text{words per block} \times \text{bytes per word}} \\ &= \frac{4 \text{ KiB}}{\frac{2 \text{ blocks}}{\text{set}} \times \frac{4 \text{ words}}{\text{block}} \times \frac{4 \text{ bytes}}{\text{word}}} \\ &= \frac{2^2 \times 2^{10} \text{ bytes}}{2^1 \times 2^2 \times 2^2 \text{ bytes/set}} \\ &= 128 \text{ sets.} \end{aligned}$$

The number of set bits is $\log_2 S = \log_2 128 = 7$. The word size is 4 bytes, so there is a 2-bit byte offset. The block size is 4 words, so the block offset is also 2 bits wide. Here is the address split:



What to Do

Refer to Problem 4 in Exercise B of this lab. Do parts (a), (b), and (c) again, with the following assumptions:

- no change in word size;
- no change in block size;
- no changes to the capacities of the caches;
- L1 D-cache changed to 8-way set-associative;
- L2 cache changed to 16-way set-associative.

What to Include in Your PDF Submission

Include solutions to the problems given in “What to Do”.