

Computer Architecture I
Spring, 2022
Homework 6

ShanghaiTech University

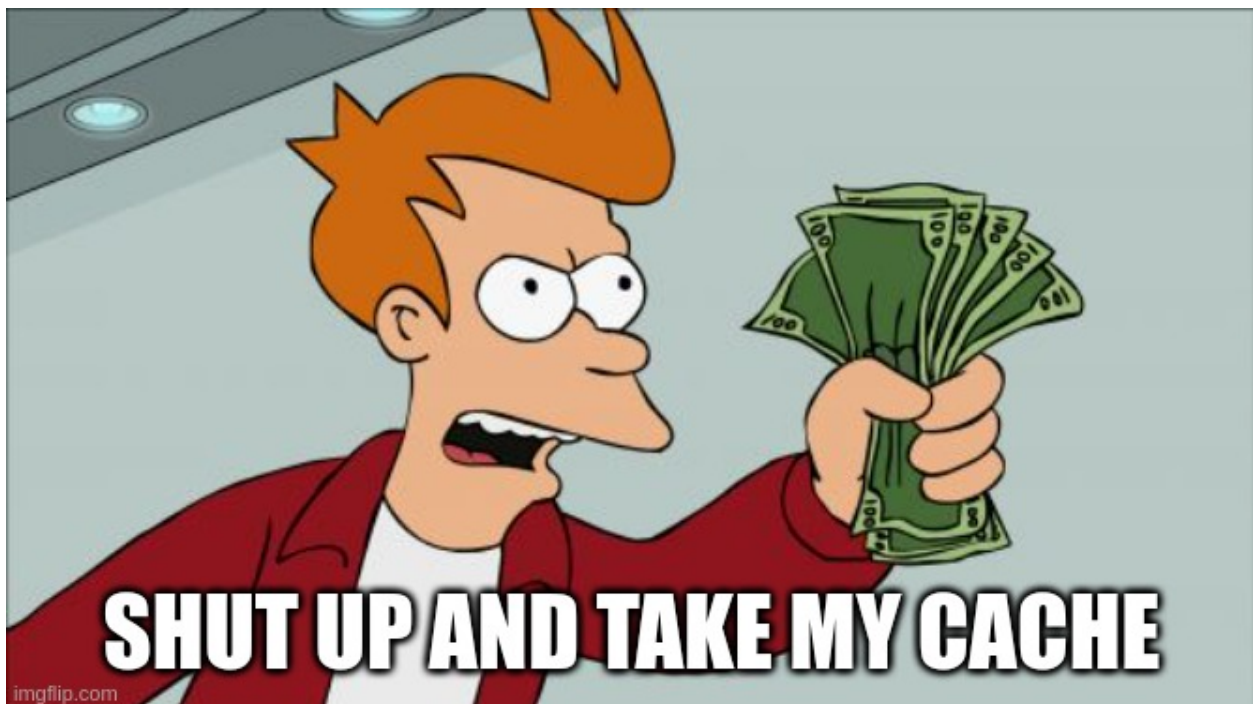
Due: May 7, 2022, 23:59

Instructions

This homework will help you review CPU caches. Remember to go through the textbook and all the lecture slides related.

Submission Guideline

- Both hand-writing and typesetting are accepted. You may find a LaTeX template helpful on our course website.
- Only PDF submissions to Gradescope will be counted. If you choose to write by hand, make sure it is transformed into a PDF document. Both scanning and taking photos are accepted.
- Please assign your answers properly on Gradescope. Any submission without proper assignment will result in a 25% point reduction.



1 Shut Up and Take My Cache!

In this section, we will review some basics of cache. A program is run on a byte-addressed system with a single-level cache. After a while, the entire cache has the state in Figure 1.

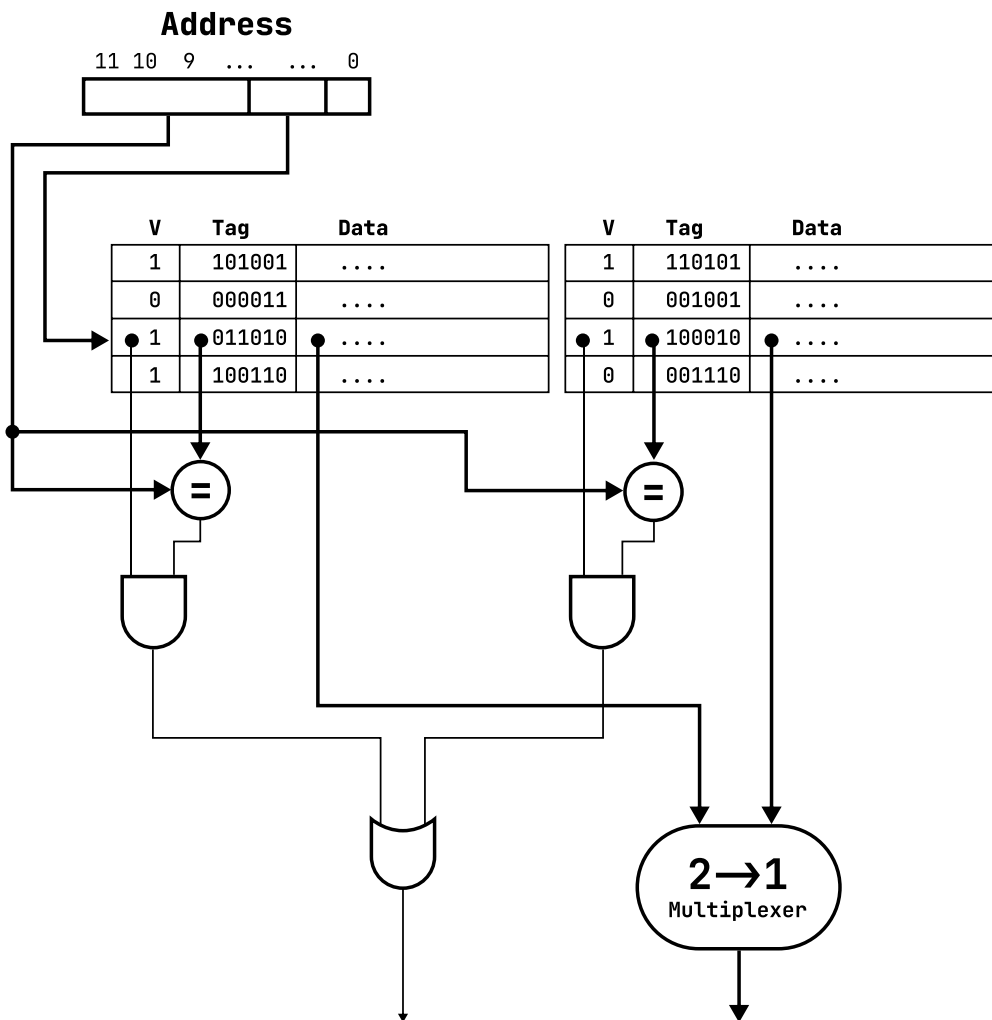


Figure 1: Layout for a cache implementation

\

1. (2 points) In Figure 1, what is a row stands for?

- ☒ A. One set ☐ B. One cache block. ☐ C. One entire cache. ☐ D. Not listed in choices.

Your answer: ☒ A ☐ B ☐ C ☐ D

2. (2 points) In Figure 1, what is V stands for?

- ☒ A. Valid bit ☐ B. Vacant bit. ☐ C. Visited bit.

Your answer: ☒ A ☐ B ☐ C

3. (2 points) What is the type of the cache described in Figure 1?

☐ A. Direct-Mapped cache.

☒ B. Set-Associative Cache (If you choose this, write down its associativity).

☐ C. Fully-Associative Cache (If you choose this, write down its associativity).

Your answer:

☐ A

☒ B (Associativity: 2)

☐ C (Associativity: _____)

4. (3 points) What is the (Tag : Set Index : Byte Offset) breakdown of memory addresses in Figure 1?

Your answer:

1. Tag: 6

2. Index: 2

3. Byte Offset: 4

5. (2 points) Is it TRUE that conflict misses cannot occur in fully-associated caches?

- ☐ A. Yes. ☐ B. No.

Your answer: ☒ A ☐ B

6. (3 points) Tell the difference(s) between Conflict Miss and Capacity Miss.

Your answer: *Conflict Miss occurs when the cache line has been replaced by others with the same index in direct map and n-way set associativity. (there is some space in other set.)*
Capacity Miss occurs when the limit capacity of cache cannot set more data, which means there are no enough space to contain all data. (cache is all full)

7. (12 points) For each of the following accesses to the cache described in Figure 1, determine if each access would be a hit or miss based on the cache state shown above. If it is a miss, classify the miss type(s). If multiple miss types may exist depending on prior memory accesses, select *all possible* miss types. For each access, you will get

- 2 points if you choose all correct choice(s),
- 1 point if you choose partial correct choice(s), and,
- 0 point if you give no choice(s) or wrong choice(s).

Each memory access should be considered *dependently*. In particular, *Do update* the cache status after each memory access. If a replacement happens, data in the first slot will always be evicted.

Order	Address	Access Outcome
1	0b 101001 00100	<u>1</u>
2	0b 011010 110100	<u>2</u>
3	0b 111110 101000	<u>3</u>
4	0b 000011 111100	<u>4</u>
5	0b 000011 011001	<u>5</u>
6	0b 100110 101100	<u>6</u>

1 tag 1

Your answer:

Note: Each access may have one or more correct choice(s).

1. ☒ Hit ☐ Compulsory Miss ☐ Conflict Miss
2. ☐ Hit ☒ Compulsory Miss ☐ Conflict Miss
3. ☐ Hit ☒ Compulsory Miss ☒ Conflict Miss
4. ☐ Hit ☒ Compulsory Miss ☒ Conflict Miss
5. ☐ Hit ☒ Compulsory Miss ☐ Conflict Miss
6. ☐ Hit ☒ Compulsory Miss ☒ Conflict Miss

8. (6 points) The specification sheet of the system is given below. What is the Average Memory Access Time (AMAT) of memory accesses in Question 7? What is the AMAT if we remove this cache? Please give your answer in nanosecond (ns). You shall have the formula, the unit of results and the deriving procedure presented in your answer. (Note: A 1 gigahertz (GHz) processor ticks a cycle for each 1 nanosecond)

System Frequency	2 GHz	/ns
Cache Access Latency	2 Cycles	
Main Memory Access Latency	600 Cycles	

Table 1: Specification Sheet

Your answer:

have cache:

$$AMAT = \text{Hit time} + \text{MissRate} \times \text{Miss Penalty}$$

$$= 2 \text{ cycles} \times \frac{1 \text{ ns}}{2 \text{ cycles}} + \frac{5}{6} \times 600 \text{ cycles} \times \frac{1 \text{ ns}}{2 \text{ cycles}}$$

$$= 25 \text{ ns}$$

remove cache:

$$AMAT = \text{Hit time} + \text{MissRate} \times \text{Miss Penalty}$$

$$= 0.4 \times 600 \text{ cycles} \times \frac{1 \text{ ns}}{2 \text{ cycles}}$$

$$= 300 \text{ ns}$$

2 Oops ... Too many bites (bytes)

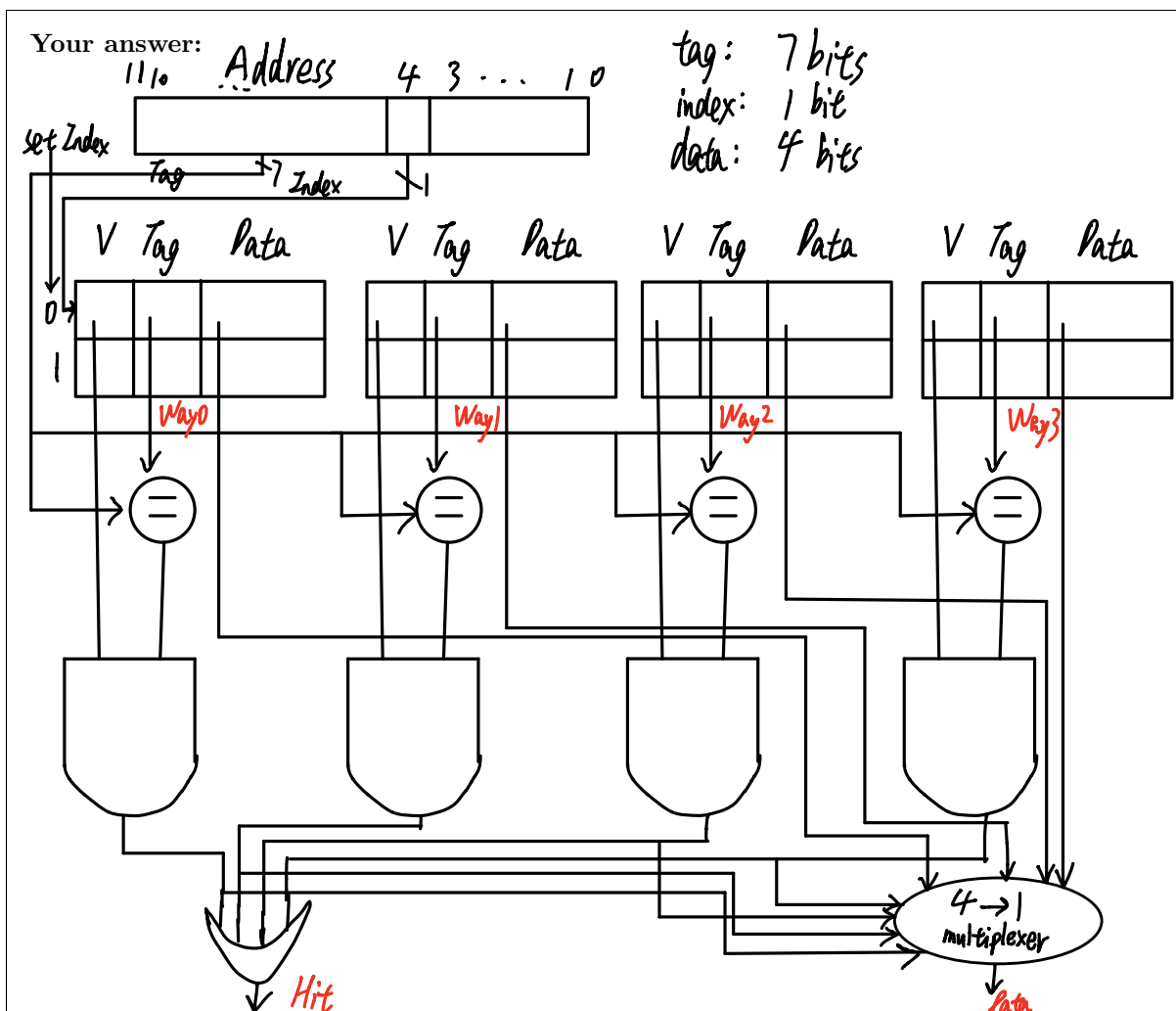
In this sections, we will review the implementation of caches and replacement policies.

1. (15 points) Let's Draw It Out! ⁴

Sketch the organization of a *four-way set associative* cache with a cache block size of 16 bytes and a total size of 128 bytes. Your sketch should have a style similar to Figure 1. The memory addresses are 12-bit long. ^{2⁷} ^{2⁴}

In your sketch, the following components should be also presented.

1. the width of set index, tag and data fields of memory addresses (3 points),
2. logical components used for comparison and selection (2 points),
3. the type of multiplexer (e.g. $2 \rightarrow 1$, $4 \rightarrow 1$, $8 \rightarrow 1$, $16 \rightarrow 1$, etc.) (1 point),
4. the number of sets (1 point) , and,
5. an implementation layout including wiring and placement of cache elements (8 points).



1 ① ② 2 ③ 3 4 ④ ⑤ ⑥
 0010 0011 0111 1100 0101 0000 0010 0001 1001 0011

In the following questions, we will examine how replacement policies affect miss rate.

After the system is cold start, the address sequence of warm-up accesses is:

0x2A, 0x3C, 0x7D, 0xC1, 0x5B, 0x01, 0x2C, 0x1D, 0x9B, 0x3E.

After all warm-up accesses are done, the address sequence of follow-up accesses is:

0x30, 0x40, 0x52, 0x44, 0x56, 0x48, 0x5A, 0x4C, 0x10, 0x3B, 0x5C, 0x30, 0x5E.

2. (2 points) Which locality(s) can you observe in the follow-up accesses?

Your answer: Note: There may exist(s) one or more correct choice(s).

☒ Spatial locality ☒ Temporal locality

3. (2 points) Assume *Least Recently Used* (LRU) replacement policy is applied. Circle out all access(es) with cache hit in the follow-up accesses.

Your answer: Circle the access(es) that meets a cache hit! (like this: 0xFF)

X
 0x30, 0x40, 0x52, 0x44, 0x56, 0x48, 0x5A, 0x4C, 0x10, 0x3B, 0x5C, 0x30, 0x5E.
 0100

4. (2 points) Assume *Most Recently Used* (MRU) replacement policy is applied. Circle out all access(es) with cache hit in the follow-up accesses.

Your answer: Circle the access(es) that meets a cache hit! (like this: 0xFF)

X
 0x30, 0x40, 0x52, 0x44, 0x56, 0x48, 0x5A, 0x4C, 0x10, 0x3B, 0x5C, 0x30, 0x5E.
 0100

5. (6 points) The specification sheet of the system is given below. What is the Average Memory Access Time (AMAT) of memory accesses in the question 3 and 4? Please give your answer in nanosecond (ns). You shall have the formula, the unit of results and the deriving procedure presented in your answer. (Note: A 1 gigahertz (GHz) processor ticks a cycle for each 1 nanosecond)

System Frequency	2 GHz
Cache Access Latency	2 Cycles
Main Memory Access Latency	130 Cycles

Table 2: Specification Sheet

Your answer: Miss Rate: Q3 $\frac{1}{13}$ Q4 $\frac{5}{13}$

Q3: $AMAT = Hit\ time + Miss\ Rate * Miss\ Penalty$

$$= 2\ Cycles * \frac{1\ ns}{2\ Cycles} + \frac{1}{13} * 130\ Cycles * \frac{1\ ns}{2\ Cycles}$$

$$= 6\ ns$$

$$Q4: AMAT = Hit\ time + Miss\ Rate * Miss\ Penalty$$

$$= 2\ Cycles * \frac{1\ ns}{2\ Cycles} + \frac{5}{13} * 130\ Cycles * \frac{1\ ns}{2\ Cycles}$$

$$= 26\ ns$$

3 Let's See Some Real World Example

Each time when you access the course website, your activity will be recorded into our web server logs! This is the definition of the web server log for our Computer Architecture course website. Assume our web server is a 32-bit machine. In this question, we will examine code optimizations to improve log processing speed. The data structure for the log is defined below.

```
struct log_entry {
    4int src_ip; /* Remote IP address */
    128char URL[128]; /* Request URL. You can consider 128 characters are enough. */
    8long reference_time; /* The time user referenced to our website. */
    64char browser[64]; /* Client browser name */
    4int status; /* HTTP response status code. (e.g. 404) */
} log[NUM_ENTRIES];
```

Assume the following processing function for the log. This function determines the most frequently observed source IPs during the given hour that succeed to connect our website.

```
topK_success_sourceIP (int hour);
```

1. (2 points) Which field(s) in a log entry will be accessed for the given log processing function?

Your answer: Note: There may exist(s) one or more correct choice(s).

☒ src_ip ☐ URL ☒ reference_time ☐ browser ☒ status

2. (1 point) Assuming 32-byte cache blocks and no prefetching, how many cache misses per entry does the given function incur on average?

Your answer: 2.5

3. (3 points) How can you reorder the data structure to improve cache utilization and access locality? Justify your modification.

Your answer:

```
struct log_entry {
    int status; /* HTTP response status code. (e.g. 404) */
    long reference_time; /* The time user referenced to our website. */
    int src_ip; /* Remote IP address */
    char browser[64]; /* Client browser name */
    char URL[128]; /* Request URL. You can consider 128 characters are enough. */
} log[NUM_ENTRIES];
```

If check status is right and reference_time is in range then this ip is the effect ip we need.
In this data structure, when we visit status, we will load reference_time and src_ip to cache by spatial locality.

4. (6 points) To mitigate the miss in the question 2, design a different data structure. How would you rewrite the program to improve the overall performance?

Your answer shall include:

- A new layout of data structure of our server logs.
- A description of how your function would improve the overall performance.
- How many cache misses per entries does your improved design incur on average?

Your answer:

```
struct log_entry {  
    int status; /* HTTP response status code. (e.g. 404) */  
    int src_ip; /* Remote IP address */  
    long reference_time; /* The time user referenced to our website. */  
} log[NUM_ENTRIES];
```

In this case, we only need 16 bytes to store one log-entry in cache, we have 32 bytes, by array and spatial locality, we can store next log-entry before we leave front one.

$\frac{1}{2} = 0.5$, only the first time before two log-entries will miss, it's a sequence
(miss, hit, hit, hit, hit, hit; miss, hit, hit, hit, hit, hit; ...)

One more thing...

Phew! That's all about cache! Tell us your feeling after finish this homework. Thank You! Don't worry if you did not assign this to Gradescope. This part is *optional*.

1. (0 points) How do you *feel* after you have done this homework?

Your answer:

- ☒ :) I feel good.
☐ :(I feel bad.

2. (0 points) Do you think this homework is hard for you?

Your answer:

- ☐ It's really difficult for me!
☒ I think it is a little bit challenging.
☐ I am okay with that.
☐ This is too easy for me.

3. (0 points) Your voice will be secretly heard by our team. Is there anything you want to feedback?

Your answer:

Thank you! Now you are clear with cache! :)