

CS152 Computer Architecture and Engineering
CS252 Graduate Computer Architecture
Spring 2021

Solution

Caches and the Memory Hierarchy

Assigned 02/10/2021

Problem Set #2

Due 02/24/2021

<http://inst.eecs.berkeley.edu/~cs152/sp21>

The problem sets are intended to help you learn the material, and we encourage you to collaborate with other students and to ask questions in discussion sections and office hours to understand the problems. However, each student must turn in their own solution to the problems.

The problem sets also provide essential background material for the exam and the midterms. The problem sets will be graded primarily on an effort basis, but if you do not work through the problem sets you are unlikely to succeed on the exam or midterms! We will distribute solutions to the problem sets on the day the problem sets are due to give you feedback. Homework assignments are due at the beginning of class on the due date, and all assignments are to be submitted through **Gradescope**. The CS152 Gradescope can be joined using the code: **X3GV6K**. Late homework will not be accepted, except for extreme circumstances and with prior arrangement.

Problem 1: Cache Access-Time & Performance

This problem requires the knowledge of Handout #2 and Lectures 6 & 7. Please, read these materials before answering the following questions.

Ben is trying to determine the best cache configuration for a new processor. He knows how to build two kinds of caches: direct-mapped caches and 4-way set-associative caches. The goal is to find the better cache configuration with the given building blocks. He wants to know how these two different configurations affect the clock speed and the cache miss-rate, and choose the one that provides better performance in terms of average latency for a load.

Problem 1.A

Access Time: Direct-Mapped

Now we want to compute the access time of a direct-mapped cache. We use the implementation shown in Figure H2-A in Handout #2. Assume a 128-KB cache with 8-word (32-byte) cache lines. The address is 32 bits and byte-addressed, so the two least significant bits of the address are ignored since a cache access is word-aligned. The data output is also 32 bits (1 word), and the MUX selects one word out of the eight words in a cache line. Using the delay equations given in Table 2.1-1, **fill in the column for the direct-mapped (DM) cache in the table.** *In the equation for the data output driver, 'associativity' refers to the associativity of the cache (1 for direct-mapped caches, A for A-way set-associative caches).*

Component	Delay equation (ps)		DM (ps)	SA (ps)
Decoder	$30 \times (\# \text{ of index bits}) + 100$	Tag	460	400
		Data	460	400
Memory array	$30 \times \log_2 (\# \text{ of rows}) + 30 \times \lceil \log_2 (\# \text{ of bits in a row}) \rceil + 100$	Tag	610	610
		Data	700	700
Comparator	$30 \times (\# \text{ of tag bits}) + 100$		550	610
N-to-1 MUX	$50 \times \log_2 N + 100$		250	250
Buffer driver	200			200
Data output driver	$50 \times (\text{associativity}) + 100$		150	300
Valid output driver	100		100	100

Table 2.1-1: Delay of each Cache Component

What is the critical path of this direct-mapped cache for a cache read? What is the access time of the cache (the delay of the critical path)? To compute the access time, assume that a 2-input gate (AND, OR) delay is 50 ps. If the CPU clock is 2.5 GHz, how many CPU cycles does a cache access take?

For the given cache structure which is byte addressable, we can know that the

$$\# \text{ of offset bits} = \log_2(\# \text{ of byte in a word line}) = 5 \text{ bits.}$$

The # of lines in the cache is

$$\# \text{ of lines} = \text{size} / \text{wordline size} = 2^{17} / 2^5 = 2^{12} \text{ lines}$$

Because the cache is direct map, then

$$\# \text{ of index bit} = \log_2(2^{12}) = 12 \text{ bits}$$

The total address bits is 32 bits, then

$$\# \text{ of tag bits} = 32 - 12 - 5 = 15 \text{ bits}$$

Applying all values we calculate above to the delay equations, we have:

$$\text{Decoder (tag)} = 30 * 12 + 100 = 460\text{ps}$$

$$\text{Decoder (data)} = 30 * 12 + 100 = 460\text{ps}$$

$$\text{Memory array (tag)} = 30 * \log_2(2^{12}) + 30 * \text{ceil}(\log_2(15+2)) + 100 = 610\text{ps}$$

$$\text{Memory array (data)} = 30 * \log_2(2^{12}) + 30 * \text{ceil}(\log_2(32*8)) + 100 = 700\text{ps}$$

$$\text{Comparator} = 30 * 15 + 100 = 550\text{ps}$$

$$\text{N-1 mux} = 50 * \log_2(8) + 100 = 250\text{ps}$$

$$\text{Data output driver} = 50 * 1 + 100 = 150\text{ps}$$

To determine the critical path for a cache read, we need to compute the time it takes to go through each path in hardware (tag check and data read). By taking the maximum delay of these two paths, we are left with the critical path.

Time to tag check valid driver from tag array

$$= \text{Decoder (tag)} + \text{Memory array (tag)} + \text{comparator} + \text{AND gate} + \text{valid output driver}$$

$$= 460 + 610 + 550 + 50 + 100 = 1770\text{ps}$$

Time to data output drive from data array

$$= \text{Decoder (data)} + \text{Memory array (data)} + \text{8-1 MUX} + \text{data output driver}$$

$$= 460 + 700 + 250 + 150 = 1560\text{ps}$$

From the above results, we can see that the critical path is tag check. **The access time is 1770ps.** At 2.5GHz, the cache access takes $(1770\text{ps} / (1/2.5\text{GHz})) = 4.4 \sim 5$ cycles. Here, **rounding up** to the nearest cycle is sensible, as this reflects how a synchronous system would work.

Problem 1.B**Access Time: Set-Associative**

We also want to investigate the access time of a set-associative cache using the 4-way set-associative cache in Figure H2-B in Handout #2. Assume the total cache size is still 128-KB (each way is 32KB), a 2-input gate delay is 50 ps, and all other parameters (such as the input address, cache line, etc.) are the same as part 2.1.A. **Compute the delay of each component and fill in the column for a 4-way set-associative cache in Table 2.1-1.**

What is the critical path of the 4-way set-associative cache? What is the access time of the cache (the delay of the critical path)? What is the main reason that the 4-way set-associative cache is slower than the direct-mapped cache? If the CPU clock is 2.5 GHz, how many CPU cycles does a cache access take?

For the given cache structure which is byte addressable, we know that the

$$\# \text{ of offset bits} = \log_2(\# \text{ of byte in a word line}) = 5 \text{ bits.}$$

The # of lines in a way of the cache

$$\# \text{ of lines} = (\$ \text{ size} / \text{wordline size}) / \text{nWays} = (2^{17}/2^5) / 4 = 2^{10} \text{ lines}$$

The number of index bits is then

$$\# \text{ of index bit} = \log_2(2^{10}) = 10 \text{ bits}$$

The total address bits is 32 bits, then

$$\# \text{ of tag bits} = 32 - 10 - 5 = 17 \text{ bits}$$

Applying all values we calculate above to the delay equations, we have:

$$\text{Decoder (tag)} = 30 * 10 + 100 = 400\text{ps}$$

$$\text{Decoder (data)} = 30 * 10 + 100 = 400\text{ps}$$

$$\text{Memory array (tag)} = 30 * \log_2(2^{10}) + 30 * \text{ceil}(\log_2((17+2)*4)) + 100 = 610\text{ps}$$

$$\text{Memory array (data)} = 30 * \log_2(2^{10}) + 30 * \text{ceil}(\log_2(32*8*4)) + 100 = 700\text{ps}$$

$$\text{Comparator} = 30 * 17 + 100 = 610\text{ps}$$

$$\text{N-1 mux} = 50 * \log_2(8) + 100 = 250\text{ps}$$

$$\text{Data output driver} = 50 * 4 + 100 = 300\text{ps}$$

There are three possible critical paths in an associative cache. The first two are the same as those in the direct mapped cache. The third one is the path through the tag array, the tag comparators, through the way-select mux, and through the data output driver.

Time to tag check valid driver

$$= \text{Decoder (tag)} + \text{Memory array (tag)} + \text{comparator} + \text{AND gate} + 2 * \text{OR gate} + \text{valid output driver}$$

$$= 400 + 610 + 610 + 50 + 100 + 100 = 1870\text{ps}$$

Time to data output drive:

= Decoder (data) + Memory array (data) + 8-1 MUX + data output driver
 = 400 + 700 + 250 + 300 = 1650ps

Time to tag valid check to output driver:

= Decoder (tag) + Memory array (tag) + comparator + AND gate + buffer driver + data output driver
 = 400 + 610 + 610 + 50 + 200 + 300 = 2170ps

From the above results, we can see that the critical path is tag valid check to output driver. The access time is 2170ps. At 2.5GHz, the cache access takes $(2170\text{ps}/(1/2.5\text{GHz})) = 5.4 \sim 6$ cycles. Here, rounding up to the nearest cycle is sensible, as this reflects how a synchronous system would work.

Problem 1.C

Miss-rate analysis

Now Ben is studying the effect of set-associativity on the cache performance. Since he now knows the access time of each configuration, he wants to know the miss-rate of each one. For the miss-rate analysis, Ben is considering two small caches: a direct-mapped cache with 8 lines with 32 bytes/line, and a 4-way set-associative cache of the same size and line size. For the set-associative cache, Ben tries out two replacement policies – least recently used (LRU) and round robin (FIFO).

Ben tests the cache by accessing the following sequence of hexadecimal byte addresses, starting with empty caches. For simplicity, assume that the addresses are only 12 bits. **Complete the following tables** for the direct-mapped cache and both types of 4-way set-associative caches showing the progression of cache contents as accesses occur (in the tables, ‘inv’ = invalid, and the column of a particular cache line contains the tag of that line). Also, for each address calculate the tag and index (which should help in filling out the table). *You only need to fill in elements in the table when a value changes.*

Address: 12 bits

Tag: 4 bits [11:8]

Index: 3 bits [7:5]

Offset: 5 bits [4:0]

D-map	Addresses and tags are in HEX								
	line in cache (tag)								hit?
	L0	L1	L2	L3	L4	L5	L6	L7	
11B	1	inv	inv	inv	inv	inv	inv	inv	no
134		1							no
20D	2								no
1A2						1			no
105	1								no
360				3					no

27D				2					no
121		1							yes
1A3						1			yes
17A				1					no
307	3								no
273				2					no
131		1							yes

	D-map
Total Misses	10
Total Accesses	13

Address: 12 bits
Tag: 6 bits [11:6]
Index: 1 bits [5:5]
Offset: 5 bits [4:0]

<u>4-way</u>	LRU -- addresses and tags are in HEX								
Address	line in cache								hit?
	Set 0				Set 1				
	way0	way1	Way2	way3	way0	way1	way2	way3	
11B	4	inv	inv	inv	inv	inv	inv	inv	no
134					4				no
20D		8							no
1A2						6			no
105	-								yes
360							D		no
27D								9	no
121					-				yes
1A3						-			yes
17A							5		no
307			C						no
273								-	yes
131					-				yes

	4-way LRU
Total Misses	8
Total Accesses	13

4-way Address	FIFO -- addresses and tags are in HEX								hit?
	line in cache (tag)								
	Set 0				Set 1				
	way0	way1	way2	way3	way0	way1	way2	way3	
11B	4	inv	inv	inv	inv	inv	inv	inv	no
134					4				no
20D		8							no
1A2						6			no
105	-								yes
360							D		no
27D								9	no
121					-				yes
1A3						-			yes
17A					5				no
307			C						no
273								-	yes
131						4			no

	4-way FIFO
Total Misses	9
Total Accesses	13

Problem 1.D

Average Latency

Assume that the results of the above analysis can represent the average miss-rates of the direct-mapped and the 4-way set-associative 128-KB caches studied in 1.A and 1.B. What would be the average memory access latency in CPU cycles for each cache (assume that the cache miss penalty is 30 cycles)? Which one is better? For the different replacement policies for the set-associative cache, which one has a smaller cache miss rate for the address stream in 1.C? Explain why. Is that replacement policy always going to yield better miss rates? If not, give a counter example using an address stream.

The miss rate for the direct-mapped cache is 10/13. The miss rate for the 4-way LRU set associative cache is 8/13. For FIFO is 9/13.

The average memory access latency is (hit time) + (miss rate) × (miss penalty).

For the direct-mapped cache, the average memory access latency would be:

(5 cycles) + (10/13) × (15 cycles) = 28.1 cycles.

For the LRU set-associative cache, the average memory access latency would be:

$(6 \text{ cycles}) + (8/13) \times (15 \text{ cycles}) = 24.5 \text{ cycles.}$

For the FIFO set-associative cache, the average memory access latency would be:

$(6 \text{ cycles}) + (9/13) \times (15 \text{ cycles}) = 26.8 \text{ cycles.}$

The set-associative cache with LRU replacement is better than the direct-mapped cache in terms of average memory access latency.

For the above example, LRU has a slightly smaller miss rate than FIFO. This is because the FIFO policy replaced tag {4} block instead of tag {D} during the 10th access, because the {4} block has been in the cache longer, even though the {D} was least recently used. In this case, **the LRU policy took better advantage of temporal locality.**

LRU does not always outperform FIFO. Assume we have a set-associative cache with the same parameters as in 1.C and an access sequence shown below. There is a miss with LRU for the last access while there is a hit with FIFO.

0x11B
0x134
0x20D
0x1A2
0x105
0x360
0x27D
0x121
0x1A3
0x17A
0x307
0x273
0x361

Problem 2: Loop Ordering

This problem requires knowledge of Lecture 7. Please, read it before answering the following questions.

This problem evaluates the cache performances for different loop orderings. You are asked to consider the following two loops, written in C, which calculate the sum of the entries in a 128 by 32 matrix of 32-bit integers:

<i>Loop A</i>	<i>Loop B</i>
<pre>sum = 0; for (i = 0; i < 128; i++) for (j = 0; j < 32; j++) sum += A[i][j];</pre>	<pre>sum = 0; for (j = 0; j < 32; j++) for (i = 0; i < 128; i++) sum += A[i][j];</pre>

The matrix A is stored contiguously in memory in row-major order. Row major order means that elements in the same row of the matrix are adjacent in memory as shown in the following memory layout:

$A[i][j]$ resides in memory location $[4 * (32 * i + j)]$

Memory Location:

0	4	124	128	4*(32*127+31)	
A[0][0]	A[0][1]	...	A[0][31]	A[1][0]	... A[127][31]

For *Problem 2.A* to *Problem 2.C*, assume that the caches are initially empty. Also, assume that only accesses to matrix A cause memory references and all other necessary variables are stored in registers. Instructions are in a separate instruction cache.

Problem 2.A

Consider a 4KB direct-mapped data cache with 8-word (32-byte) cache lines.
Calculate the number of cache misses that will occur when running Loop A.
Calculate the number of cache misses that will occur when running Loop B.

Each element of the 128x32 matrix A can only be mapped to one particular cache location in this direct-mapped data cache. Since each row has 32 32-bit integers, and since each cache line can hold 8 32-bit ints, a row of the matrix occupies the lines in 4 consecutive sets of the cache.

Loop A—where each iteration of the inner loop sums a row of A—accesses memory addresses in a linear sequence. Given this access pattern, the access to the first word in each cache line will miss, but the next seven accesses will hit. After sequentially moving through this line, it will not be accessed again, so its later eviction will not cause any future misses. Therefore, Loop A will only have compulsory misses for the 512 (128 rows x 4 lines per row) that matrix A spans.

The consecutive accesses in Loop B will move in a stride of 32 words. Therefore, the inner loop will touch the first element in 128 cache lines before the next iteration of the outer loop. While intuition might suggest that the 128 lines could all fit in the cache with 128 sets, there is a complicating factor: each row is four cache lines past the previous row, meaning that the lines accessed when traversing the first column go in indices 0, 4, 8, 16, and so on. Since the lines containing the column are competing for only one quarter of the sets, the lines loaded when starting a column are evicted by the time the column is complete, preventing any reuse. Therefore, all 4096 (128 x 32) accesses miss.

The number of cache misses for Loop A: _____512_____

The number of cache misses for Loop B: _____4096_____

Problem 2.B

Consider a direct-mapped data cache with 8-word (32-byte) cache lines. Calculate the minimum number of cache lines required for the data cache if Loop A is to run without any cache misses other than compulsory misses. Calculate the minimum number of cache lines required for the data cache if Loop B is to run without any cache misses other than compulsory misses.

Since Loop A accesses memory sequentially, we can sum all the elements in a cache line and then never touch it again. Therefore, we only need to hold 1 active line at any given time to avoid all but compulsory misses.

For Loop B to run without any cache misses other than compulsory misses, the data cache needs to have the ability to hold one column of matrix A in the cache. Since the consecutive accesses in the inner loop of Loop B will use one out of every four cache lines, and since we have 128 rows, Loop B requires 512 (128×4) lines to avoid all but compulsory misses.

Data-cache size required for Loop A: _____1_____ cache line(s)

Data-cache size required for Loop B: _____512_____ cache line(s)

Problem 2.C

Consider a 4KB set-associative data cache with 4 ways, and 8-word (32-byte) cache lines. This data cache uses a first-in/first-out (FIFO) replacement policy. Calculate the number of cache misses that will occur when running Loop A. Calculate the number of cache misses that will occur when running Loop B.

Loop A still only has 512 (128 rows x 4 lines per row) compulsory misses.

Loop B still cannot fully utilize the cache. The first 8 accesses will allocate into way 1 in sets 0, 4, 8, 16, etc.; the next 8 accesses will allocate into way 2 of those same sets; and so on. After 32 accesses, all four ways will be filled, and the next 8 accesses along the column will evict the previous lines in way 1, preventing any reuse. Therefore, all 4096 (128×32) accesses miss.

The number of cache misses for Loop A: _____512_____

The number of cache misses for Loop B: _____4096_____

Problem 3: Microtagged Cache

In this problem, we explore *microtagging*, a technique to reduce the access time of set-associative caches. Recall that for associative caches, the tag check must be completed before load results are returned to the CPU, because the result of the tag check determines which cache way is selected. Consequently, the tag check is often on the critical path.

The time to perform the tag check (and, thus, way selection) is determined in large part by the size of the tag. We can speed up way selection by checking only a subset of the tag—called a microtag—and using the results of this comparison to select the appropriate cache way. Of course, the full tag check must also occur to determine if the cache access is a hit or a miss, but this comparison proceeds in parallel with way selection. We store the full tags separately from the microtag array.

We will consider the impact of microtagging on a 4-way set-associative 16KB data cache with 32-byte lines. Addresses are 32 bits long. Microtags are 8 bits long. The baseline cache (i.e. without microtagging) is depicted in Figure H2-B in the handout 2. Figure 1, below, shows the modified tag comparison and driver hardware in the microtagged cache.

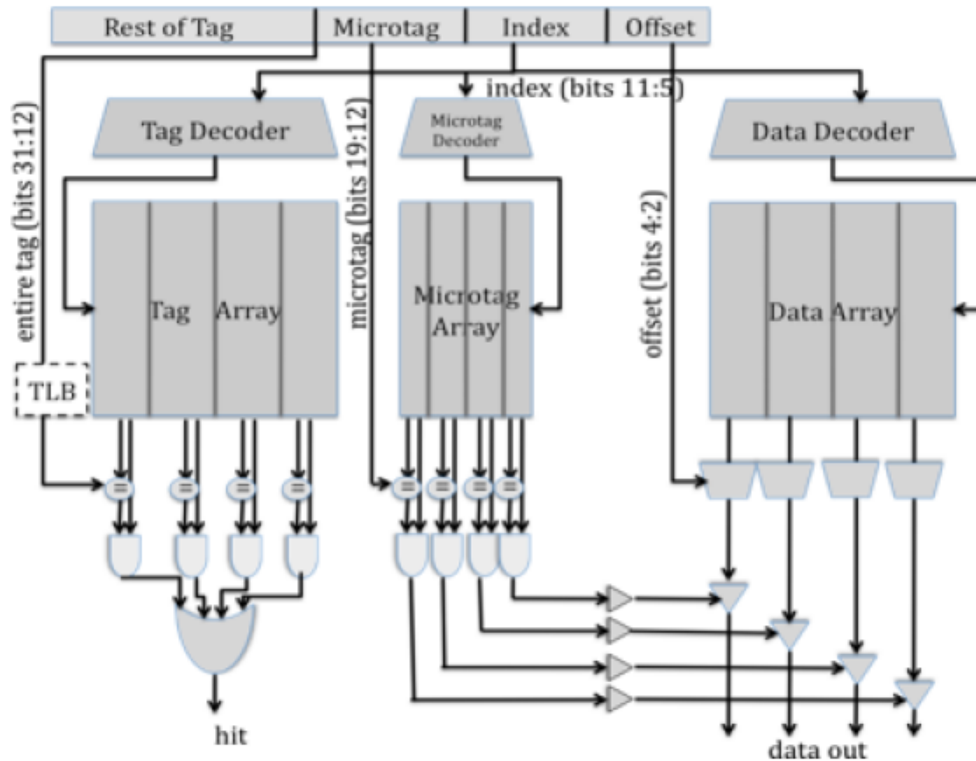


Figure 2.4-1: Microtagged cache datapath

Problem 3.A**Cache Cycle Time**

Table 2.4-1, below, contains the delays of the components within the 4-way set-associative cache, for both the baseline and the microtagged cache. For both configurations, determine the critical path and the cache access time (i.e., the delay through the critical path).

Assume that the 2-input AND gates have a 50ps delay and the 4-input OR gate has a 100ps delay.

Component	Delay equation (ps)		Baseline	Microtagged
Decoder	$20 \times (\# \text{ of index bits}) + 100$	Tag	240	240
		Data	240	240
Memory array	$20 \times \log_2 (\# \text{ of rows}) + 20 \times \log_2 (\# \text{ of bits in a row}) + 100$	Tag	330	330
		Data	440	440
		Microtag		300
Comparator	$20 \times (\# \text{ of tag bits}) + 100$	Tag	500	500
		Microtag		260
N-to-1 MUX	$50 \times \log_2 N + 100$		250	250
Buffer driver	200		200	200
Data output driver	$50 \times (\text{associativity}) + 100$		300	300
Valid output driver	100		100	100

Table 2.4-1: Delay of each Cache Component

What is the old critical path? The old cycle time (in ps)?

Candidate 1: Full tag check

tag decoder → tag read → comparator → 2-in AND → 4-in OR → valid output driver
 $240 \text{ ps} + 330 \text{ ps} + 500 \text{ ps} + 50 \text{ ps} + 100 \text{ ps} + 100 \text{ ps} = 1320 \text{ ps}$

Candidate 2: Data select based on full tag check

tag decoder → tag read → comparator → 2-in AND → buffer driver → data output driver
 $240 \text{ ps} + 330 \text{ ps} + 500 \text{ ps} + 50 \text{ ps} + 200 \text{ ps} + 300 \text{ ps} = 1620 \text{ ps}$

Candidate 3: Data readout

data decoder → data read → 4-to-1 MUX → data output driver
 $240 \text{ ps} + 440 \text{ ps} + 250 \text{ ps} + 300 \text{ ps} = 1230 \text{ ps}$

The critical path is the data select based on the full tag match. The cycle time is 1620 ps.

What is the new critical path? The new cycle time (in ps)?

Candidate 1: Full tag check

same as baseline full tag check => 1320 ps

Candidate 2: Data select based on microtag check

μtag decoder → μtag read → comparator → 2-in AND → buffer driver → data out driver

240 ps + 300 ps + 260 ps + 50 ps + 200 ps + 300 ps = 1350 ps

Candidate 3: Data readout

same as baseline data read => 1230 ps

The critical path is the data select based on the microtag check. The cycle time is 1350 ps.

Problem 3.B

AMAT

Assume temporarily that both the baseline cache and the microtagged cache have the same hit rate, 90%, and the same average miss penalty, 10 ns. Using the cycle times computed in 3.A as the hit times, compute the average memory access time for both caches. **What was the old AMAT (in ns)? What is the new AMAT (in ns)?**

$AMAT = (hit_time) + (miss_rate) \times (miss_penalty)$
 $= X + (0.10)(10ns) = X + 1ns$, where X is the hit time calculated from 3.A

Old AMAT = 1.620 + 1 = 2.620 ns

New AMAT with microtags = 1.350 + 1 = 2.350 ns

Problem 3.C

Constraints

Microtags add an additional constraint to the cache: in a given cache set, all microtags must be unique. This constraint is necessary to avoid multiple microtag matches in the same set, which would prevent the cache from selecting the correct way.

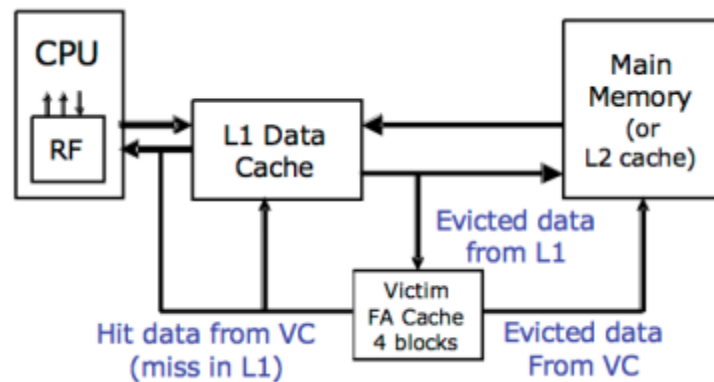
State which of the 3C's of cache misses this constraint affects. **How will the cache miss rate compare to an ordinary 4-way set-associative cache? How will it compare to that of a direct-mapped cache of the same size?**

Because the uniqueness property of microtags restricts the replacement policy, the cache isn't free to make as optimal replacement decisions as it could in the baseline. This will lead to some increase in conflict misses. The magnitude of this effect depends on which 8 bits are selected to form the microtag. In principle, using the bottom 8 bits would result in more potential for microtag collisions and would add the biggest restriction to the ability of the cache to hold spatially local data. However, it will still be better than a directmapped cache of the same size and line size.

Problem 4: Victim Cache Evaluation

Although direct-mapped caches have an advantage of smaller access time than set-associative caches, they have more conflict misses due to their lack of associativity. In order to reduce these conflict misses, Norm Jouppi proposed victim caching, where a small fully-associative back up cache, called a victim cache, is added to a direct-mapped L1 cache to hold recently evicted cache lines.

The following diagram shows how a victim cache can be added to a direct-mapped L1 data cache. Upon a data access, the following chain of events takes place:



1. The L1 data cache is checked. If it holds the data requested, the data is returned.
2. If the data is not in the L1 cache, the victim cache is checked. If it holds the data requested, the data is moved into the L1 cache and sent back to the processor. The data evicted from the L1 cache is put in the victim cache, and put at the end of the FIFO replacement queue.
3. If neither of the caches holds the data, it is retrieved from memory, and put in the L1 cache. If the L1 cache needs to evict old data to make space for the new data, the old data is put in the victim cache and placed at the end of the FIFO replacement queue. Any data that needs to be evicted from the victim cache to make space is written back to memory or discarded, if unmodified.

Note that the two caches are *exclusive*. That means that the same data cannot be stored in both L1 and victim caches at the same time.

Problem 4.A

Baseline Cache Design

The diagram below shows our victim cache, a 32-Byte fully associative cache with four 8-Byte cache lines. Each line contains of two 4-Byte words and has an associated tag and two status bits (valid and dirty). The Input Address is 32-bits and the two least significant bits are assumed to be zero. The output of the cache is a 32-bit word.

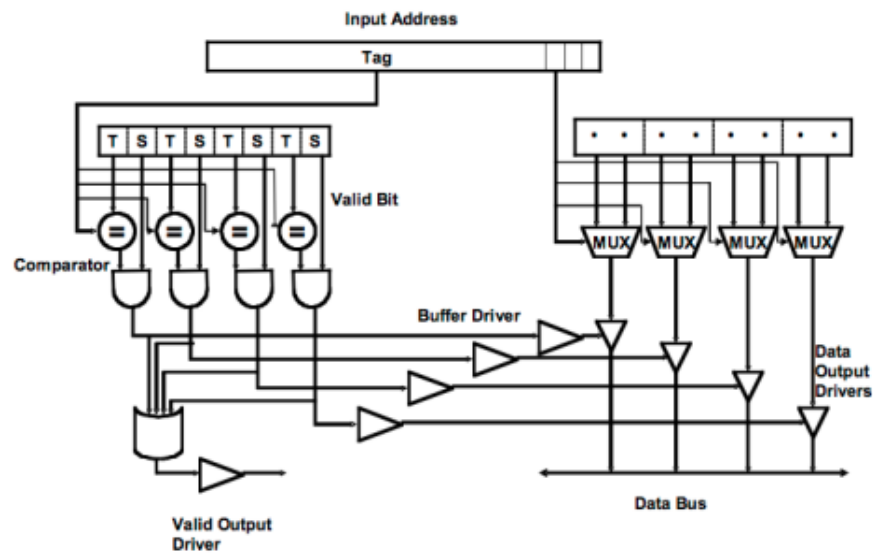


Figure 2.5-1: Victim cache datapath

Please complete Table 2.5-1 with delays across each element of the cache. Using the data you compute in Table 2.5-1, calculate the critical path delay through this cache (from when the Input Address is set to when both Valid Output Driver and the appropriate Data Output Driver are outputting valid data).

Component	Delay equation (ps)	FA(ps)
Comparator	$30 \times (\# \text{ of tag bits}) + 100$	970
N-to-1 MUX	$50 \times \log_2 N + 100$	150
Buffer driver	200	200
AND gate	100	100
OR gate	$50 \times \log_2 N + 100$	200
Data output driver	$50 \times (\text{associativity}) + 100$	300
Valid output driver	100	100

Table 2.5-1: Delay of each cache component

Critical Path Cache Delay:

Below, we evaluate the three major paths through the victim cache to find the critical path and cycle time. Note that the victim cache is fully-associative and uses 29-bit tags.

Candidate 1: Tag check

comparator → 2-in AND → 4-in OR → valid output driver

$$970 \text{ ps} + 100 \text{ ps} + 200 \text{ ps} + 100 \text{ ps} = 1370 \text{ ps}$$

Candidate 2: Data select based on tag check

comparator → 2-in AND → buffer driver → data output driver

$$970 \text{ ps} + 100 \text{ ps} + 200 \text{ ps} + 300 \text{ ps} = 1570 \text{ ps}$$

Candidate 3: Data readout

2-to-1 MUX → data output driver

$$200 \text{ ps} + 300 \text{ ps} = 500 \text{ ps}$$

The critical path is the data select based on the tag match. The cycle time is 1570 ps.

Problem 4.B

Victim Cache Behavior

Now we will study the impact of a victim cache on cache hit rate. Our main L1 cache is a 128 byte, direct-mapped cache with 16 bytes per cache line. The cache is word (4-bytes) addressable. The victim cache in Figure 2.5-1 is a 32-byte fully associative cache with 16 bytes per cache line, and is also word addressable. The victim cache uses the first in first out (FIFO) replacement policy.

Please complete Table 2.5-2 showing a trace of memory accesses. In the table, each entry contains the tag of that line, or “inv”, if no data is present. You should only fill in elements in the table when a value changes. For simplicity, the addresses are only 8 bits. The first 3 lines of the table have been filled in for you. For your convenience, the address breakdown for access to the main cache is depicted below.

7	6		4	3		2	1		0	
TAG	INDEX				WORD SELECT			BYTE SELECT		

Input Address	Main Cache (tag)									Victim Cache (tag)		
	L0	L1	L2	L3	L4	L5	L6	L7	Hit?	Way0	Way1	Hit?
	inv	inv	inv	inv	inv	inv	inv	inv	-	inv	inv	-
0	0								N			N
80	1								N	0		N
4	0								N	8		Y
A0			1						N			N

10		0							N			N
C0					1				N			N
18		0							Y			
20			0						N		A	N
8C	1								N	0		Y
28			0						Y			
AC			1						N		2	Y
38				0					N			N
C4					1				Y			
3C				0					Y			
48					0				N	C		N
0C	0								N		8	N
24			0						N	A		N

Table 2.5-2: Memory access trace

Problem 4.C

Average Memory Access Time

Assume **15%** of memory accesses are resolved in the victim cache. If retrieving data from the victim cache takes **5 cycles** and retrieving data from main memory takes **55 cycles**, by how many cycles does the victim cache improve the average memory access time?

$$AMAT = HitTime + L1MissRate * L1MissPenalty$$

$$AMAT^* = HitTime + L1MissRate * (VictimHitRate * VictimHitTime + (1 - VictimHitRate) * VictimMissPenalty)$$

VictimMissPenalty = L1MissPenalty = DRAMTime, since this is just time to get data from main memory

$$AMAT - AMAT^* = L1MissRate * (DRAMTime - VictimHitRate * VictimHitTime - (1 - VictimHitRate) * DRAMTime)$$

$$= L1MissRate * (55 - 0.15 * 5 - 0.85 * 55)$$

$$= L1MissRate * 7.5$$

Since L1 Miss rate was not provided, answer could have also just been 7.5 cycles, which is average improvement to L1 miss penalty

Problem 5: Three C's of Cache Misses

Mark whether the following modifications will cause each of the categories to **increase**, **decrease**, or whether the modification will have **no effect**. You can assume the baseline cache is set associative. **Explain your reasoning.**

	Compulsory Misses	Conflict Misses	Capacity Misses
<p>Halving the line size (associativity and # sets constant)</p> <p>Halves capacity!</p>	<p>Increase</p> <p>Shorter lines mean fewer adjacent elements are brought in with the first access to a given line.</p>	<p>Increase</p> <p>The program will access more cache lines in total, creating more opportunity for conflict misses.</p>	<p>Increase</p> <p>Capacity has been cut in half.</p>
<p>Doubling the number of sets (capacity and line size constant)</p> <p>Halves associativity!</p>	<p>No effect</p> <p>Halving associativity doesn't change when lines are first brought into the cache</p>	<p>Increase</p> <p>Typically, lower associativity increases conflict misses, since there are fewer places to put the same element.</p>	<p>No effect</p> <p>Capacity does not change.</p>

	Compulsory Misses	Conflict Misses	Capacity Misses
Adding prefetching	<p>Decrease</p> <p>Ideally, a good prefetcher can bring data in before we use it, avoiding compulsory misses.</p>	<p>Best answer: Decrease With good prefetching, conflict misses should decrease, as the prefetcher will often bring lines that have been evicted back into the cache.</p> <p>Okay answer: increase With mediocre prefetching, conflict misses could increase, as the prefetcher could evict useful lines to bring in useless.</p> <p>Other okay answer: no effect With mediocre prefetching that uses a stream buffer or other ancillary structure, there will be little to no effect on conflict misses.</p>	<p>Best answer: Decrease With good prefetching, capacity misses should decrease. In a situation where the working set simply won't fit, the prefetcher can dynamically bring lines in, "Just-In-Time," avoiding what would have been capacity misses.</p> <p>Okay answer: no effect With a mediocre prefetcher that would increase conflict misses, it is unlikely that capacity misses would increase. If prefetcher traffic evicts useful data, newly created misses will almost</p>
Combine ICache and DCache into a single L1 cache with the combined capacity (associativity and line size constant)	No effect	May increase: New opportunities for conflicts between cache lines for data and cache lines for instructions are introduced	Decrease: Greater capacity

Problem 6: Memory Hierarchy Performance

Mark whether the following modifications will cause each of the categories to **increase**, **decrease**, or whether the modification will have **no effect**. You can assume the baseline cache is set associative. **Explain your reasoning**.

	Hit Time	Miss Rate	Miss Penalty
Halving the line size (associativity and # sets constant) Halves # of capacity	Decreases The cache is now physically smaller, which overshadows the slightly increased tag check time (tag grows by 1 bit).	Increases Smaller capacity, less ability to take advantage of spatial locality within a single cache line (more compulsory misses).	Decreases Smaller lines can be brought in more quickly. OR No effect, because cache already brings in critical word first.

<p>Doubling the number of sets (capacity and line size constant)</p> <p>Halves associativity</p>	<p>Decreases</p> <p># of sets increases, so tags get smaller. Fewer tags must be checked, and fewer ways have to be muxed outs.</p>	<p>Increases</p> <p>More conflict misses because associativity gets halved.</p>	<p>No effect</p> <p>This is dominated by the outer memory hierarchy</p>
<p>Adding prefetching</p>	<p>No effect</p> <p>The prefetcher isn't on the hit path.</p>	<p>Decreases</p> <p>The whole purpose of a prefetcher is to reduce the miss rate by bringing in data ahead of time.</p>	<p>Good answer: no effect.</p> <p>May increase due to bandwidth pollution but we can(should) give a priority on cache misses over prefetch requests.</p> <p>May decrease because a prefetch can be inflight when a miss occurs (but this is unlikely).</p>
<p>Combine L1ICache and L1DCache into a single L1 cache with the combined capacity (associativity and line size constant)</p>	<p>Increase:</p> <p>If the cache is dual-ported, it will be slower than a single-ported cache</p> <p>If there is a single port, then frequently data accesses may stall for instruction accesses, or vice-versa</p>	<p>May Decrease</p> <p>Cache can more flexibly allocate space towards either data or instructions, depending on dynamic program behavior.</p> <p>May increase:</p> <p>Edge cases may cause more conflict misses between instruction and data accesses</p>	<p>No effect:</p> <p>This is dominated by outer memory hierarchy</p>

