

The principle of locality justifies the use of:

- A. Interrupts
  - B. DMA
  - C. Polling
  - D. Cache Memory**

A computer system has a  $4K$  word cache organized in block-set-associative manner with 4 blocks per set, 64 words per block. The number of bits in the SET and WORD fields of the main memory address format is:      $n$  way  $\in$   $n$  lines in

- A. 15, 40
  - B. 6, 4
  - C. 7, 2
  - D. 4, 6

$$\# \text{ words} = 64$$

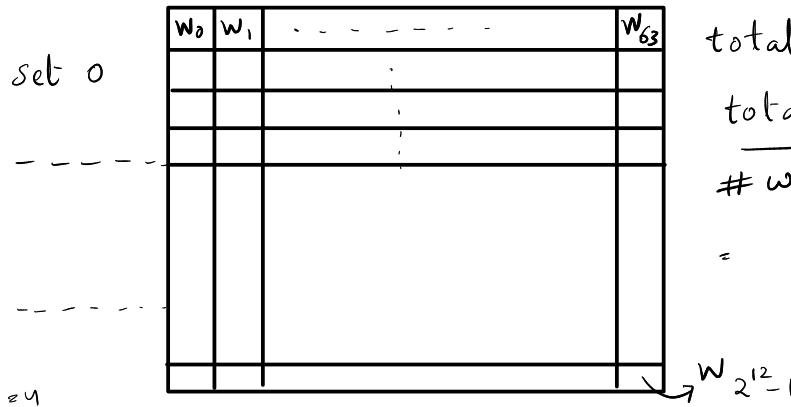
$$\# \text{ word bits} = \log_2(\#\text{words}) = \log_2(64) = 6 \quad (\text{enough})$$

$$\text{LR cache size} = 4 \times 2^{10} \text{ words} = 2^{12} \text{ words}$$

$$\# \text{ Sets} = \frac{\# \text{ lines}}{\# \text{ lines per set}}$$

$$= \frac{2^6}{4}$$

$$= 2^4$$



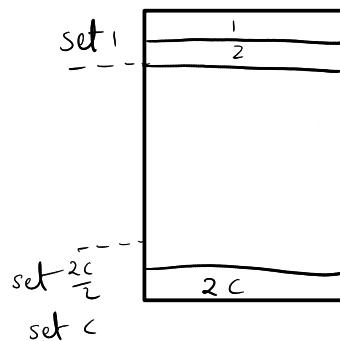
total # lines =

total # words

$$\frac{\# \text{ words per line}}{12} = \frac{?}{\frac{12}{6}} = ?$$

The main memory of a computer has  $2 \text{ cm}$  blocks while the cache has  $2 \text{ c}$  blocks. If the cache uses the set associative mapping scheme with 2 blocks per set, then block  $k$  of the main memory maps to the set:

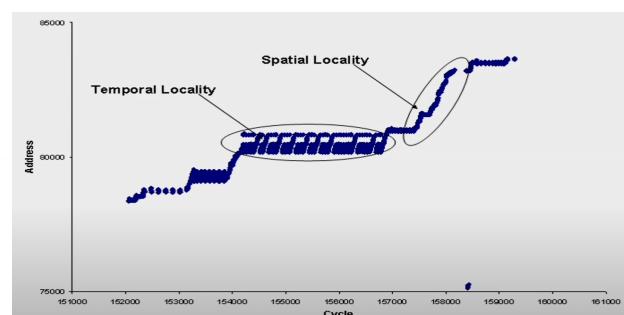
- A.  $(k \bmod m)$  of the cache
  - ~~B.  $(k \bmod c)$  of the cache~~
  - C.  $(k \bmod 2c)$  of the cache
  - D.  $(k \bmod 2cm)$  of the cache



so  $k$  must go into  
one of  $c$  sets

More than one word are put in one cache block to:

- A. exploit the temporal locality of reference in a program
  - B. exploit the spatial locality of reference in a program
  - C. reduce the miss penalty
  - D. none of the above



A CPU has 32-bit memory address and a 256 KB cache memory. The cache is organized as a 4-way set associative cache with cache block size of 16 bytes. Each line can hold # words  $\Rightarrow 2^9 B \rightarrow 2^8 B = 2^{18} B \equiv$  cache size is  $2^{18}$  words byte addressable

- A. What is the number of sets in the cache?  $2^{12}$
  - B. What is the size (in bits) of the tag field per cache block? 16 bits
  - C. What is the number and size of comparators required for tag matching? 16
  - D. How many address bits are required to find the byte offset within a cache block? 4
  - E. What is the total amount of extra memory (in bytes) required for the tag bits?  $2^{15}$

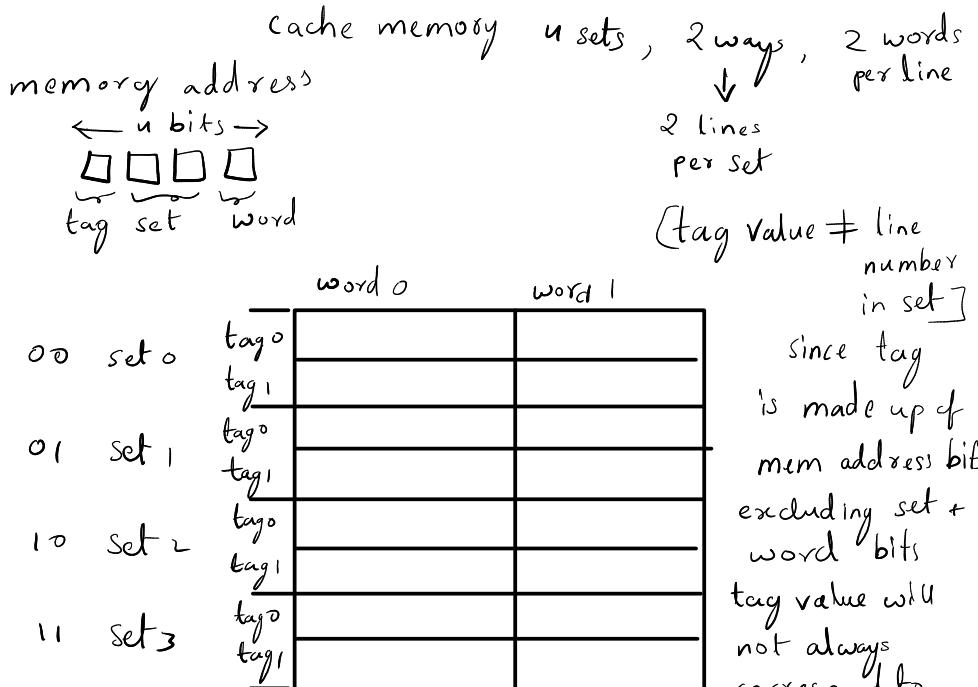
$$\left( \frac{\# \text{ words}}{\# \text{ words per line}} \right) = \left( \frac{2^{18}}{2^4} \right) = \frac{2^{14}}{2^2} = 2^{12} = \# \text{ sets}$$

$$\begin{array}{lll}
 B & \# \text{tag bits} & \# \text{set bits} & \# \text{word bits} \\
 32 - \# \text{set bit} & \log_2(2^{10}) = 10 & \log_2(2^4) = 4 \\
 - \# \text{word bit} & & \\
 32 - 12 - 4 = 16 & & \\
 & \hookrightarrow \text{tag bits} &
 \end{array}$$

- c. purpose of tag/line comparators : to identify if cpu requested memory  
is present in cache or not

Some other example for understanding purpose

main memory	
address	word
0000	a
0001	b
0010	c
0011	d
0100	e
0101	f
0110	g
0111	h
1000	i
1001	j
1010	k
1011	l
1100	m
1101	n
1110	o
1111	p

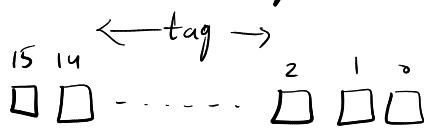


when we receive a memory address from CPU we will know in which set we need to look based on set bits. Inside each set we need to compare every tag, so # comparators needed = # lines per set

Coming to original question, # comparators needed =  $n$  each with 16 bits

Note as # lines per set  $\uparrow$  # comparators needed also  $\uparrow$   
 in turn hardware complexity, latency, cost also  $\uparrow$

E.) # tag bits = 16



extra memory here is tag directory

size of tag directory = # tag bits  $\times$  # total lines in cache

$$= 16 \times \# \text{ sets} \times \# \text{ lines per set}$$

$$= 16 \times 2^{12} \times 4$$

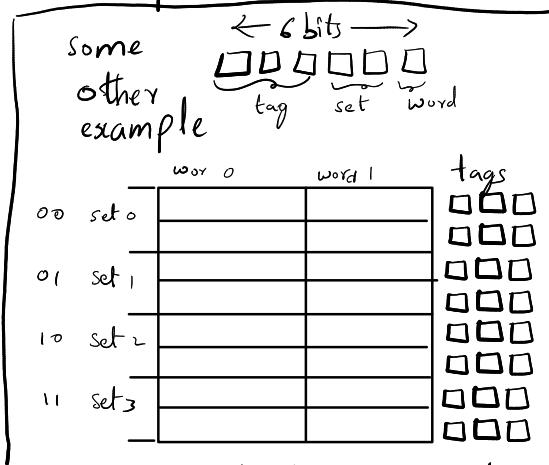
$$= 2^4 \times 2^{12} \times 2^2$$

$$= 2^{18} \text{ bits}$$

$$= 2^{15} \times 2^3 \text{ bits}$$

$$= 2^{15} \times 8 \text{ bits}$$

$$= 2^{15} \text{ bytes}$$



every cache line store a tag

#### 1.2.15 Cache Memory: GATE CSE 2005 | Question: 67

<https://gateoverflow.in/1390>

Consider a direct mapped cache of size 32 KB with block size 32 bytes. The CPU generates 32 bit addresses. The number of bits needed for cache indexing and the number of tag bits are respectively,

- A. 10, 17
- B. 10, 22
- C. 15, 17
- D. 5, 17

line number

$$\text{cache size } 32 \text{ KB} = 2^5 \times 2^{10} \text{ B} = 2^{15} \text{ Bytes}$$

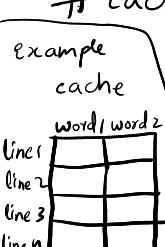
block size = 32 bytes [memory block = cache line]

block size = line capacity

assuming byte addressable memory  $\Rightarrow$  1 byte = 1 word

so one cache line can hold 32 words

$$\begin{aligned} \# \text{ bits required to locate word in a line} &= \log_2(\# \text{ words}) \\ &= \log_2(32) = 5 \end{aligned}$$



1 byte = 1 byte

cache size = total no. of bytes

$$= 8 \text{ bytes}$$

line capacity = 2 bytes

$$\# \text{ cache line} = \frac{8 \text{ bytes}}{2 \text{ bytes}} = 4$$

# bits required to locate a line in cache

= # bits required for cache indexing

$$= \log_2(\# \text{ cache lines}) = \log_2(2^4) = 4$$

# tag bits = total no. of memory address bits - # line bits

$$= 32 - 4 - 5 = 17$$

- # word bits

## Scenario Setup

Imagine we have a very small **direct-mapped cache** that has only **4 blocks** (or lines). In a direct-mapped cache, each block from main memory can only go to one specific location in the cache. We can determine this location using the formula:

```
(Memory Block Address) % (Number of Cache Blocks)
```

Let's say we have the following memory blocks and their corresponding cache block destinations:

Memory Block Address	Calculation	Cache Block Destination
0	$0 \% 4$	0
1	$1 \% 4$	1
2	$2 \% 4$	2
3	$3 \% 4$	3
4	$4 \% 4$	0
5	$5 \% 4$	1
6	$6 \% 4$	2
7	$7 \% 4$	3
8	$8 \% 4$	0

### 1. Compulsory Miss (Cold Miss)

This happens the very first time we access a block of data.

Let's access **Memory Block 2**.

**Initial State:** The cache is empty.

Code snippet

```
Cache Block 0: [     ]
Cache Block 1: [     ]
Cache Block 2: [     ] <--- We want to access data from Memory Block 2
Cache Block 3: [     ]
```

**Result:** It's a **miss** because the cache is empty. The data from Memory Block 2 is fetched from main memory and placed into Cache Block 2.

**Final State:**

Code snippet

```
Cache Block 0: [     ]
Cache Block 1: [     ]
Cache Block 2: [Data from Mem Block 2] <--- Data is now cached
Cache Block 3: [     ]
```

This first miss was unavoidable because the cache had to be populated.

### 2. Conflict Miss

This happens when the cache has space, but the specific block we need is already occupied by data that maps to the same location.

Let's continue from the previous state and now access **Memory Block 0**, and then **Memory Block 4**.

**Step 1: Access Memory Block 0** (This is also a compulsory miss)

Code snippet



```
Cache Block 0: [     ] <--- Accessing Mem Block 0
Cache Block 1: [     ]
Cache Block 2: [Data from Mem Block 2]
Cache Block 3: [     ]
```

**Result:** The cache places data from Memory Block 0 into Cache Block 0.

Code snippet



```
Cache Block 0: [Data from Mem Block 0]
Cache Block 1: [     ]
Cache Block 2: [Data from Mem Block 2]
Cache Block 3: [     ]
```

### Step 2: Now, access Memory Block 4

Remember, Memory Block 4 also maps to Cache Block 0 ( $4 \% 4 = 0$ ).

**Initial State:**

Code snippet

```
Cache Block 0: [Data from Mem Block 0] <--- Want Mem Block 4, but Block 0 is here!
Cache Block 1: [     ]
Cache Block 2: [Data from Mem Block 2]
Cache Block 3: [     ]
```

**Result:** This is a **conflict miss**. Even though Cache Blocks 1 and 3 are empty, Memory Block 4 must go into Cache Block 0. The cache is forced to **evict** the data from Memory Block 0 to make room.

**Final State:**

Code snippet

```
Cache Block 0: [Data from Mem Block 4] <--- Old data is replaced
Cache Block 1: [     ]
Cache Block 2: [Data from Mem Block 2]
Cache Block 3: [     ]
```

### 3. Capacity Miss

This happens when the cache is full, and we need to fetch data that was previously evicted simply because there wasn't enough space for all the data we are actively using.

Let's fill our entire cache and then request an old item.

**Step 1: Fill the cache**

We access Memory Blocks 0, 1, 2, and 3.

Code snippet



```
Cache Block 0: [Data from Mem Block 0]
Cache Block 1: [Data from Mem Block 1]
Cache Block 2: [Data from Mem Block 2]
Cache Block 3: [Data from Mem Block 3]
```

Our cache is now completely full.

**Step 2: Access a new block, forcing an eviction**

Let's access **Memory Block 5**. This maps to Cache Block 1 ( $5 \% 4 = 1$ ), evicting the data from Memory Block 1.

Code snippet

```
Cache Block 0: [Data from Mem Block 0]
Cache Block 1: [Data from Mem Block 5] <--- Mem Block 1 was evicted
Cache Block 2: [Data from Mem Block 2]
Cache Block 3: [Data from Mem Block 3]
```

**Step 3: Now, try to access Memory Block 1 again**

**Initial State:** The cache is still full.

Code snippet

```
Cache Block 0: [Data from Mem Block 0]
Cache Block 1: [Data from Mem Block 5] <--- We want Mem Block 1, but it's gone!
Cache Block 2: [Data from Mem Block 2]
Cache Block 3: [Data from Mem Block 3]
```

**Result:** This is a **capacity miss**. Memory Block 1 was in the cache, but we had to evict it because the cache wasn't large enough to hold all the data we needed (Blocks 0, 1, 2, 3, and 5). It's a miss caused by the **limited size** of the cache. The cache will now evict the data from Memory Block 1 to load the data from Memory Block 5.

A CPU has a  $32KB$  direct mapped cache with  $128$  byte-block size. Suppose  $A$  is two dimensional array of size  $512 \times 512$  with elements that occupy  $8$ -bytes each. Consider the following two C code segments,  $P1$  and  $P2$ .

P1:

$$\text{cache size} = 2^5 \times 2^{10} \text{ bytes} = 2^{15} \text{ bytes}$$

$$\text{line capacity} = 128 \text{ bytes} = 2^7 \text{ bytes}$$

$$\# \text{lines} = \frac{2^{15}}{2^7} = 2^8$$

P2:

```
for (i=0; i<512; i++)
{
    for (j=0; j<512; j++)
    {
        x += A[i][j];
    }
} column major
```

goclasses.in

tests.gate

$P1$  and  $P2$  are executed independently with the same initial state, namely, the array  $A$  is not in the cache and

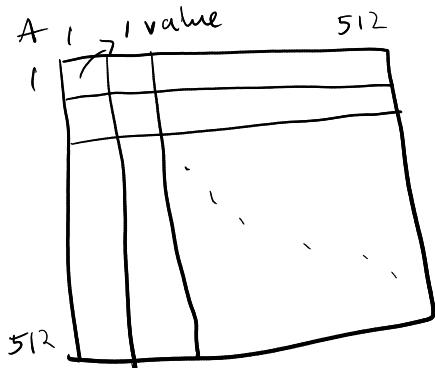
 $i$ , $j$ ,

$x$  are in registers. Let the number of cache misses experienced by  $P1$  be  $M_1$  and that for  $P2$  be  $M_2$ .

The value of  $M_1$  is:

- A. 0
- B. 2048
- C. 16384
- D. 262144

main memory

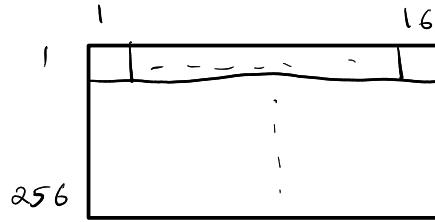


512 rows, 512 columns

one cache line is 128 bytes

one value in array cell occupies 8 bytes

No. of values one cache line can hold =  $\frac{128}{8} = 16$



$i$      $j$      $A[i][j]$     hit/miss

0    0     $A[0][0]$     miss    move block of 16 values

0    1     $A[0][1]$     hit

0    15     $A[0][15]$     hit

0    16     $A[0][16]$     miss

$A[0][0]$  to  $A[0][15]$  into cache line 1

$A[0][16]$  to  $A[0][31]$  into cache line 2

So in every row out of 512 values for every 16 values there is a miss

$$\# \text{misses} = M_1 = 512 \times \frac{512}{16}$$

$$= 512 \times 32 = 16384$$

for same setting above

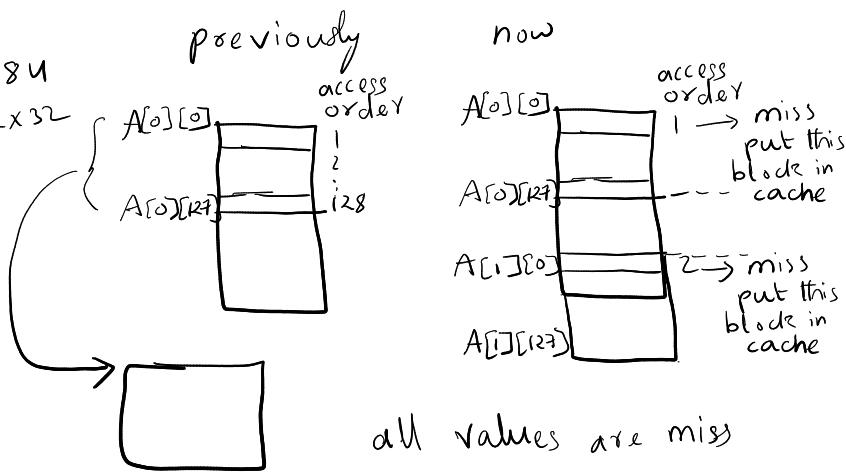
The value of the ratio  $\frac{M_1}{M_2}$ :

- A. 0 →  $M_1 \neq 0$
- B.  $\frac{1}{16}$
- C.  $\frac{1}{8}$
- D.  $16 \rightarrow M_1 < M_2$

$$\text{so } \frac{M_1}{M_2} < 1$$

$$\frac{M_1}{M_2} = \frac{512 \times 32}{512 \times 512} = \frac{2^5}{2^9} = \frac{1}{2^4} = \frac{1}{16}$$

$$M_2 = \# \text{ misses} = 512 \times 512$$



all values are miss

#### 1.2.20 Cache Memory: GATE CSE 2007 | Question: 10

<https://gateoverflow.in/1208>

Consider a 4-way set associative cache consisting of 128 lines with a line size of 64 words. The CPU generates a 20-bit address of a word in main memory. The number of bits in the TAG, LINE and WORD fields are respectively:

- A. 9, 6, 5
- B. 7, 7, 6
- C. 7, 5, 8
- D. 9, 5, 6

given,

$n\text{-way} \equiv \text{each set has } n \text{ lines}$

$$\begin{aligned} \# \text{ sets} &= \frac{\# \text{ lines}}{\# \text{ lines per set}} \\ &= \frac{128}{n} = 32 \end{aligned}$$

$$\text{so } \# \text{ sets} = 32$$

$$\text{then } \# \text{ bits required to locate set} = \log_2(\# \text{ sets}) = \log_2(32) = 5$$

$$\# \text{ words in a line} = 64$$

$$\text{then } \# \text{ bits required to locate word} = \log_2(\# \text{ words}) = \log_2(64) = 6$$

$$\# \text{ bits required to locate tag} = \text{total no. of bits} - \# \text{ set bits} - \# \text{ word bits}$$

$$= 20 - 5 - 6$$

$$= 20 - 11 = 9$$

#### 1.2.21 Cache Memory: GATE CSE 2007 | Question: 80

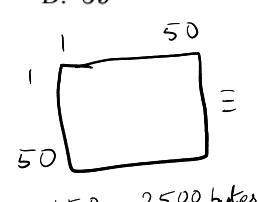
<https://gateoverflow.in/1273>

Consider a machine with a byte addressable main memory of  $2^{16}$  bytes. Assume that a direct mapped data cache consisting of 32 lines of 64 bytes each is used in the system. A  $50 \times 50$  two-dimensional array of bytes is stored in the main memory starting from memory location  $1100H$ . Assume that the data cache is initially empty. The complete array is accessed twice. Assume that the contents of the data cache do not change in between the two accesses.

How many data misses will occur in total?

- A. 48
- B. 50
- C. 56
- D. 59

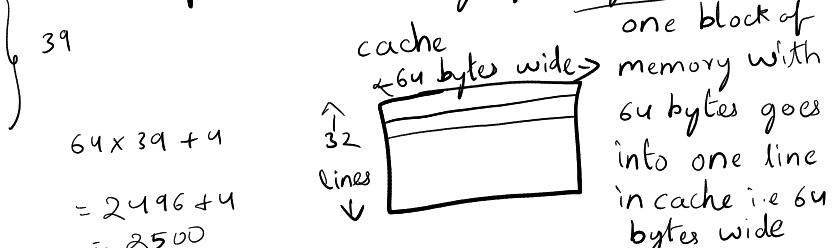
block of memory

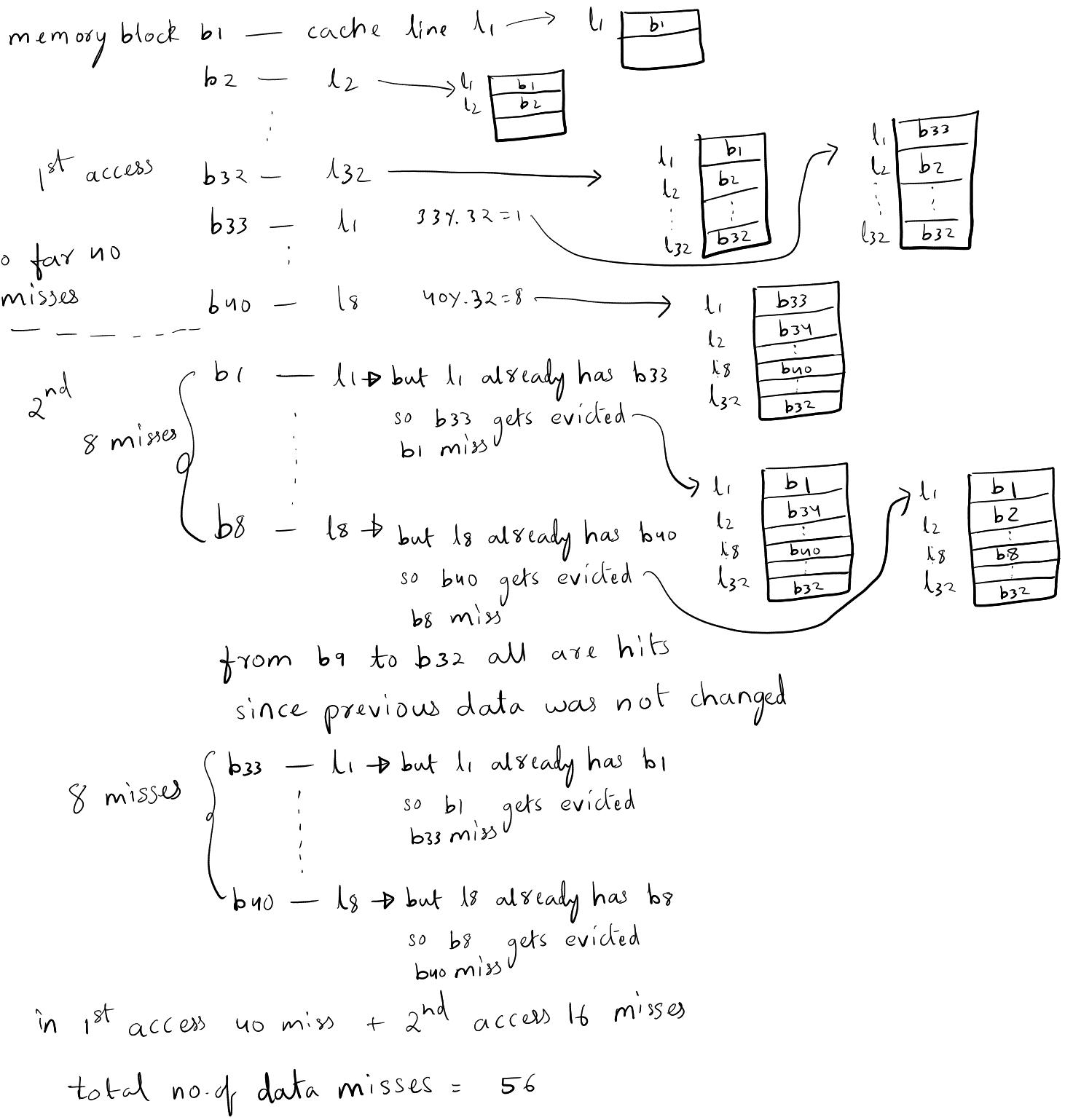


1	50	b1	1—64	64 bytes
1	50	b2	65—128	64 bytes
			⋮	
		b39	2433—2496	64 bytes
		b40	2497—2500	4 bytes

$$50 \times 50 = 2500 \text{ bytes}$$

Assuming access order to be row-major order, why if it was column-major then # misses = # elements in array but options are significantly less than that





total 128 lines

n line sets = n way associative

A block-set associative cache memory consists of 128 blocks divided into four block sets. The main memory consists of 16,384 blocks and each block contains 256 eight bit words.

- How many bits are required for addressing the main memory? 22
- How many bits are needed to represent the TAG, SET and WORD fields?

9 5 8

1 main memory blocks goes into 1 cache line

main memory    1 word = 8 bits }  $\Rightarrow$  memory is byte addressable  
 $= 1 \text{ byte}$

block 1    word 1    word 2    :    word 256    # bits required to locate a word in mem block / cache line =  $\log_2(\# \text{words})$   
 $= \log_2(256) = 8$

block 16384

$$\begin{aligned}
 \text{main memory size} &= \text{total no. of words that it can hold} \\
 &= \text{no. of blocks} \times \text{no. of words per block} \\
 &= 16384 \times 256 \text{ words} \\
 &= 2^{14} \times 2^8 \text{ bytes} \\
 &= 2^{22} \text{ bytes}
 \end{aligned}$$

# bits required to locate a word or byte in memory

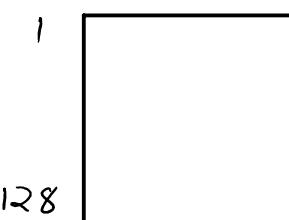
$$\begin{aligned}
 &= \# \text{ bits required to address main memory} \\
 &= \log_2(\text{main memory size}) = \log_2(2^{22}) = 22
 \end{aligned}$$

from given,

total lines in cache = 128

lines per set = 4

$$\Rightarrow \# \text{ sets} = \frac{\text{total lines}}{\text{line per set}} = \frac{128}{4} = 32$$



128

# bits required to locate a set =  $\log_2(\# \text{sets}) = \log_2(32) = 5$ 

# tag bits = total address bits - # set bits - # word bit = 22 - 5 - 8 = 9



Consider a machine with a 2-way set associative data cache of size 64 Kbytes and block size 16 bytes. The cache is managed using 32 bit virtual addresses and the page size is 4 Kbytes. A program to be run on this machine begins as follows:

```
double ARR[1024][1024];
int i, j;
/*Initialize array ARR to 0.0 */
for(i = 0; i < 1024; i++)
    for(j = 0; j < 1024; j++)
        ARR[i][j] = 0.0;
```

The size of double is 8 bytes. Array  $ARR$  is located in memory starting at the beginning of virtual page 0xFF000 and stored in row major order. The cache is initially empty and no pre-fetching is done. The only data memory references made by the program are those to array  $ARR$ .

The total size of the tags in the cache directory is:

- A. 32 Kbits
- B. 34 Kbits
- C. 64 Kbits
- D. 68 Kbits

to know tag directory size we need to know  
# tag bits

to know # tag bits need to know

# set bits & # word bits      total # bytes

$$\text{total \# lines in cache} = \frac{\text{cache size}}{\text{line size}} = \frac{\text{cache can hold}}{\# \text{bytes per line}}$$

$$= \frac{64 \text{ KB}}{16 \text{ B}} = \frac{2^6 \times 2^{10} \text{ B}}{2^4 \text{ B}} = 2^{12}$$

$$\# \text{sets} = \frac{\text{total \# lines in cache}}{\# \text{lines per set}} = \frac{2^{12}}{2} = 2^{11}$$

→ 2 since 2 way set associative

# word bits = no. of bits required to locate word or byte in one memory block or one cache line  
 $= \log_2(\text{line size}) = \log_2(16) = 4$

# tag bits = no. of bits in memory address - # set bits  
- # word bits

$$= 32 - 11 - 4 = 17$$

tag directory size  $\neq 2^{\# \text{tag bits}} = 2^{17}$

# all possible tags =  $2^{17}$

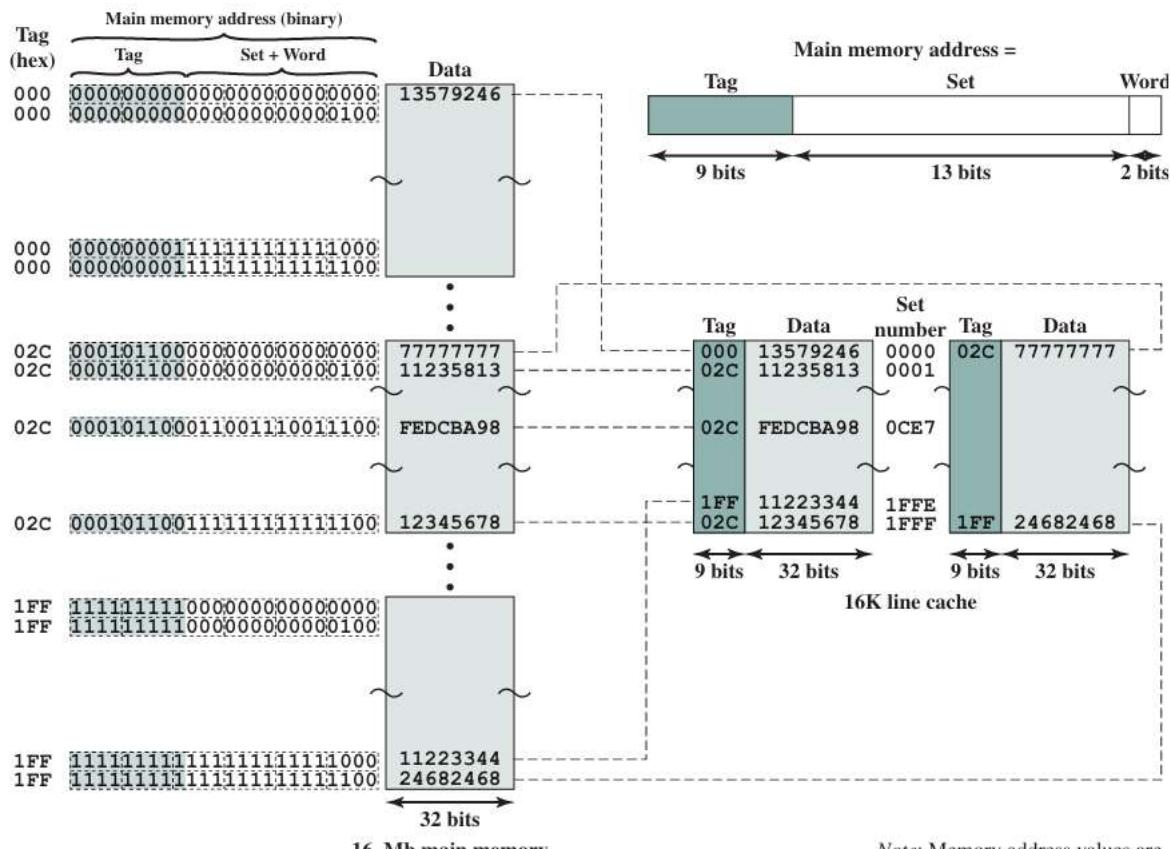


Figure 4.15 Two-Way Set-Associative Mapping Example

Note: Memory address values are  
in binary representation;  
other values are in hexadecimal.

total no. of tags = total no. of lines in cache i.e. for every line there is a tag

so tag directory size =

$$(\# \text{tags or } \# \text{cache lines}) \times \# \text{tag bits}$$

$$= 2^{12} \times 17 \text{ bits}$$

$$= 4096 \times 17 \text{ bits} = 69632 \text{ bits}$$

convert bits to kilo bits

$$1 \text{ Kilo bits} = 1024 \text{ bits (exact)}$$

$$? = 69632 \text{ bits}$$

$$? = \frac{69632}{1024} = \frac{2^{12} \times 17}{2^{10}} = 2^2 \times 17 = 68$$

$\therefore$  tag directory size = 68 kilo bits

Consider a machine with a 2-way set associative data cache of size 64 Kbytes and block size 16 bytes. The cache is managed using 32 bit virtual addresses and the page size is 4 Kbytes. A program to be run on this machine begins as follows:

```
double ARR[1024][1024];
int i, j;
/*Initialize array ARR to 0.0 */
for(i = 0; i < 1024; i++)
    for(j = 0; j < 1024; j++)
        ARR[i][j] = 0.0;
```

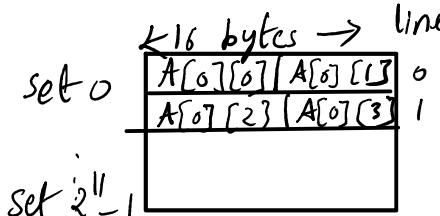
Good explanation including Virtual page address  
This conversation has lot more info do check it  
<https://gemini.google.com/app/91c25fb7ed2008a3>

The size of double is 8 bytes. Array *ARR* is located in memory starting at the beginning of virtual page 0xFF000 and stored in row major order. The cache is initially empty and no pre-fetching is done. The only data memory references made by the program are those to array *ARR*.

Which of the following array elements have the same cache index as *ARR*[0][0]?

- A. *ARR*[0][4]
- B. *ARR*[4][0]
- C. *ARR*[0][5]
- D. *ARR*[5][0]

↳ Set number



single value in A occupies 8 bytes, so one cache line can hold 2 values

The problem tells us:

Array *ARR* is located at the beginning of virtual page 0xFF000.

- This 0xFF000 is the **Virtual Page Number (VPN)**. Let's check if it's 20 bits. In hexadecimal, F is 4 bits (1111). So, FF000 is 5 hex digits, which is  $5 \times 4 = 20$  bits. Perfect!
- "At the beginning" means the **offset is 0**. We need to represent this with 12 bits, which is simply 0000 0000 0000 in binary, or 0x000 in hex.

Now, we just combine the two parts to form the full 32-bit starting address:

- **VPN (20 bits):** 0xFF000
- **Offset (12 bits):** 0x000

Combine them: 0xFF000 000 = 0xFF000000

## 2. Finding the Cache Index for *ARR*[0][0]

The array *ARR* starts at the beginning of virtual page 0xFF000. With a page size of 4 KBytes (2<sup>12</sup> bytes), the starting address is 0xFF000 followed by 12 zero bits.

- **Address of *ARR*[0][0]:** 0xFF000000

Let's look at the lower 15 bits of this address in binary:

...000 0000 0000 0000

- **Index (bits 14-4):** 00000000000 (which is 0 in decimal).
- **Offset (bits 3-0):** 0000

We need to find which option also has an index of 0 then it falls in same set. The address of an element *ARR*[i][j] is calculated as: **Address = Base Address + (i \* 1024 + j) \* sizeof(double)** where the Base Address is 0xFF000000 and sizeof(double) is 8 bytes

**A. *ARR*[0][4]** Mul 1024 Since arr is 1024 x 1024 we need to skip 1024 elements whenever we move from one row to another

- Byte Offset =  $(0 \times 1024 + 4) \times 8 = 32$ .
- Address = 0xFF000000 + 32 = 0xFF000020.
- Binary (lower bits): ...0010 0000.
- Index bits (bits 14-4): ...0000010. This is 2. ✗

**C. *ARR*[0][5]**

- Byte Offset =  $(0 \times 1024 + 5) \times 8 = 40$ .
- Address = 0xFF000000 + 40 = 0xFF000028.
- Binary (lower bits): ...0010 1000.
- Index bits (bits 14-4): ...0000010. This is 2. ✗

"cache is managed using 32 bit virtual addresses and the page size is 4 KB". This indicates that word or byte in cache is located using virtual address only, memory management unit does address translation when needed and we don't have to worry about it. Remember in

Book [Page Number | Line Number] is like a Main Memory [Physical Frame Number | Word Offset], i.e like Virtual Memory [Virtual Page Number | Page Offset] after translation. Virtual address = VPN + page offset

**B. *ARR*[4][0]**

- Byte Offset =  $(4 \times 1024 + 0) \times 8 = 4 \times 1024 \times 8 = 32768$ .
- Note that 32768 =  $2^{15}$ .
- Address = 0xFF000000 + 32768 = 0xFF008000.
- Binary (lower bits): ...1000 0000 0000 0000.
- Index bits (bits 14-4): ...00000000000. This is 0. ✓ and ends at 15, so to locate next value we multiply by 8

**D. *ARR*[5][0]**

- Byte Offset =  $(5 \times 1024 + 0) \times 8 = 40 \times 1024 = 40960$ .
- Address = 0xFF000000 + 40960 = 0xFF00A000.
- Binary (lower bits): ...1010 0000 0000 0000.
- Index bits (bits 14-4): ...01000000000. This is 512. ✗

Mul sizeof(double) i.e 8 since in memory one value occupies 8 bytes so if one value starts at address 0 it spans till 7 next value starts at 8 and ends at 15, so to locate next value we multiply by 8



Consider a machine with a 2-way set associative data cache of size 64 Kbytes and block size 16 bytes. The cache is managed using 32 bit virtual addresses and the page size is 4 Kbytes. A program to be run on this machine begins as follows:

```
double ARR[1024][1024];
int i, j;
/*Initialize array ARR to 0.0 */
for(i = 0; i < 1024; i++)
    for(j = 0; j < 1024; j++)
        ARR[i][j] = 0.0;
```

goclasses.in

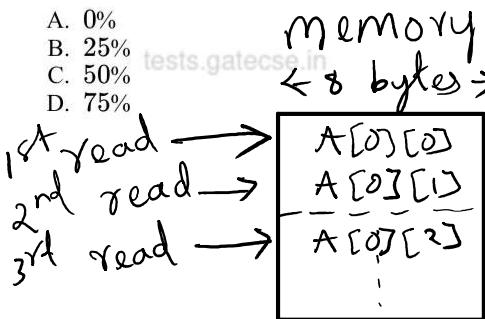
tests.gatecse.in

The size of double is 8 bytes. Array  $ARR$  is located in memory starting at the beginning of virtual page  $0xFF000$  and stored in row major order. The cache is initially empty and no pre-fetching is done. The only data memory references made by the program are those to array  $ARR$ .

The cache hit ratio for this initialization loop is:

Memory block  $\rightarrow$  cache line/block

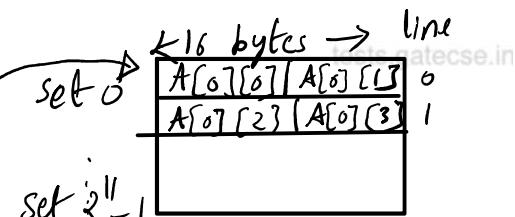
- A. 0%  
B. 25%  
C. 50%  
D. 75%



1st is a miss so  
 this memory  
 block is written  
 to appropriate  
 cache line

2nd read is Hit

3rd read is miss



so out of every  
 2 reads 1 is a hit  
 & 1 is a miss  
 since one cache line  
 can hold only 2 values  
 at a time

so hit ratio =  $\frac{1}{2} = 50\%$

hypothetical scenario:

if cache could hold  $u$  values in a line  
 then hit ratio =  $\frac{3}{u}$  + miss ratio =  $\frac{1}{u}$

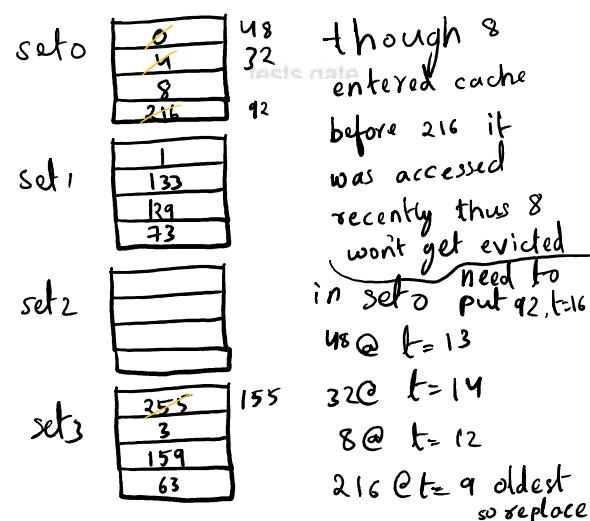
Consider a 4-way set associative cache (initially empty) with total 16 cache blocks. The main memory consists of 256 blocks and the request for memory blocks are in the following order:

0, 255, 1, 4, 3, 8, 133, 159, 216, 129, 63, 8, 48, 32, 73, 92, 155.

Which one of the following memory block will NOT be in cache if LRU replacement policy is used?

- A. 3      t=1       $t=2$        $t=3$   
 B. 8       $0 \% .u = 0, 255 \% .u = 3, 1 \% .u = 1$   
 C. 129      $t=u$       $t=5$       $t=6$   
 D. 216      $9 \% .u = 0, 3 \% .u = 3, 8 \% .u = 0$   
  
 $t=7$       $t=8$       $t=9$   
 133 \% .u = 1, 159 \% .u = 3, 216 \% .u = 0  
  
 $t=10$      $t=11$      $t=12$   
 129 \% .u = 1, 63 \% .u = 3, 8 \% .u = 0 (hit)  
  
 $t=13$      $t=14$      $t=15$   
 8 \% .u = 0, 32 \% .u = 0, 73 \% .u = 1  
  
 $t=16$      $t=17$   
 92 \% .u = 0, 155 \% .u = 3

$u\text{-way} \equiv u \text{ lines per set} \quad \# \text{sets} = \frac{\# \text{cache lines}}{\text{lines per set}} = \frac{16}{u} = 4$

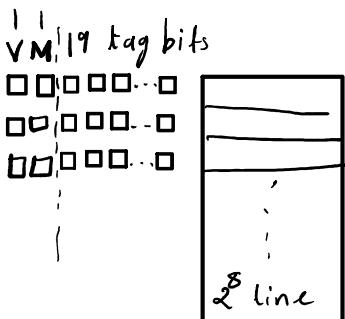


An 8KB direct-mapped write-back cache is organized as multiple blocks, each size of 32-bytes. The processor generates 32-bit addresses. The cache controller contains the tag information for each cache block comprising of the following.

- 1 valid bit
- 1 modified bit
- As many bits as the minimum needed to identify the memory block mapped in the cache.

What is the total size of memory needed at the cache controller to store meta-data (tags) for the cache?

- A. 4864 bits  
B. 6144 bits  
C. 6656 bits  
D. 5376 bits



each cache line has a tag of 19 bits

then only tag part from address requested can be compared vs tags stored for already cached

data in tag directory (part of cache memory system)

and whether data is present or not can be determined.

$$\text{given cache size} = 8 \text{ kB} = 2^3 \times 2^{10} \text{ B} = 2^{13} \text{ B}$$

$$\text{cache line can hold } 32 \text{ B} = 2^5 \text{ B}$$

$$\# \text{ lines} = \frac{\text{cache size}}{\text{size of line}} = \frac{2^{13}}{2^5} = 2^8$$

$$\begin{array}{c} \text{memory address 32 bits} \\ \hline \text{32-8-5} \\ = 19 \end{array} \quad \begin{array}{c} \xleftarrow{8} \text{bits for} \\ \text{cache line} \end{array} \quad \begin{array}{c} \xleftarrow{5} \text{bits for} \\ \text{word inline} \end{array}$$

$$\begin{aligned} \text{tag directory size} &= \text{no. of bits for tag} \\ &\quad \times \\ &\quad \text{no. of cache lines} \end{aligned}$$

$$\begin{aligned} &= (1+1+19) \times 2^8 \\ &= 21 \times 256 \\ &= 5376 \end{aligned}$$

Note: if write-through policy was used then no need of modified/dirty bit

Feature	Write-Through	Write-Back
<b>Write Operation</b>	Data is written to both cache and main memory simultaneously.	Data is written only to the cache.
<b>Performance</b>	Slower write operations due to main memory latency.	Faster write operations.
<b>Data Consistency</b>	High. Main memory is always up-to-date.	Lower. Main memory can be temporarily out of sync.
<b>Complexity</b>	Simpler to implement.	More complex, requires a dirty bit.
<b>Data Loss Risk</b>	Low. Data is immediately saved to main memory.	Higher. Data can be lost in case of power failure.
<b>Bandwidth Usage</b>	High. Every write goes to main memory.	Low. Reduces writes to main memory.



A computer has a 256-KByte, 4-way set associative, write back data cache with block size of 32-Bytes. The processor sends 32-bit addresses to the cache controller. Each cache tag directory entry contains, in addition to address tag, 2 valid bits, 1 modified bit and 1 replacement bit.

The number of bits in the tag field of an address is

- A. 11
- B. 14
- C. 16
- D. 27

$$\text{cache size} = 256 \text{ KB} = 2^8 \times 2^{10} \text{ B} = 2^{18} \text{ B}$$

$$\text{cache line capacity} = 32 \text{ byte} = 2^5 \text{ B} \text{ or } 2^5 \text{ words}$$

n-way set = n lines per set

$$\# \text{ cache lines} = \frac{\text{cache size}}{\text{line capacity}} = \frac{\text{total no. of words}}{\text{cache can hold}} = \frac{2^{18} \text{ B}}{2^5 \text{ B}} = 2^{13}$$

$$\# \text{ sets} = \frac{\# \text{ cache lines}}{\# \text{ lines per set}} = \frac{2^{13}}{4} = \frac{2^{13}}{2^2} = 2^{11}$$

$$\begin{aligned} \# \text{ tag bits} &= \text{total address bits} - \# \text{ set bits} - \# \text{ word bits} \\ &= 32 - 11 - 5 = 16 \end{aligned}$$

# word bits  
=  $\log_2(2^5) = 5$



A computer has a 256-KByte, 4-way set associative, write back data cache with block size of 32 Bytes. The processor sends 32 bit addresses to the cache controller. Each cache tag directory entry contains, in addition to address tag, 2 valid bits, 1 modified bit and 1 replacement bit.

The size of the cache tag directory is:

$$\text{no. of tag bits} \rightarrow \text{actual tag bits} + \text{additional bits like these} = (2+1+1+16) \times 2^{11} \text{ bits} = 20 \times 2^{11} = 20 \times 2 \times 2^{10} \text{ bits} = 40 \text{ kilobits}$$

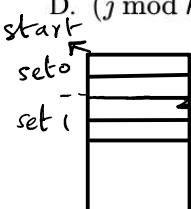
- A. 160 Kbits
- B. 136 Kbits
- C. 40 Kbits
- D. 32 Kbits

Mapping Technique	Number of Tags	Tag Size	Number of Tag Comparisons per Access	Reason for Comparisons
Direct-Mapped	Number of Cache Lines	Smallest	1	The memory address's index points to a <b>single, specific cache line</b> . The system only needs to check the tag of that one line to determine a hit or miss.
N-Way Set-Associative	Number of Cache Lines	Medium	N (ways)	The address's index points to a <b>specific set</b> . The desired block could be in any of the 'N' lines within that set, so all 'N' tags in the set must be checked simultaneously.
Fully Associative	Number of Cache Lines	Largest	All lines in the cache	There is <b>no index</b> , so a memory block could be in any line in the entire cache. To find it, the system must compare the address tag against the tags of every single line in the cache at once.



In a  $k$ -way set associative cache, the cache is divided into  $v$  sets, each of which consists of  $k$  lines. The lines of a set are placed in sequence one after another. The lines in set  $s$  are sequenced before the lines in set  $(s+1)$ . The main memory blocks are numbered 0 onwards. The main memory block numbered  $j$  must be mapped to any one of the cache lines from

- A.  $(j \bmod v) * k$  to  $(j \bmod v) * k + (k-1)$
- B.  $(j \bmod v)$  to  $(j \bmod v) + (k-1)$
- C.  $(j \bmod k)$  to  $(j \bmod k) + (v-1)$
- D.  $(j \bmod k) * v$  to  $(j \bmod k) * v + (v-1)$



2-way  
set num x 2  
tells us how many lines we must skip from start address to reach intended line

mem block j must go into one of v sets  
so its  $j \bmod v = i \bmod v \rightarrow$  set number  
 $(j \bmod v) \times k \rightarrow$  tells where the starting line in each set is  
 $(j \bmod v) \times k + (k-1) \rightarrow$  last line in set  
 $(j \bmod v) \times k + 0, (j \bmod v) \times k + 1, \dots, (j \bmod v) \times k + (k-1)$  i.e.  $k$  lines



An access sequence of cache block addresses is of length  $N$  and contains  $n$  unique block addresses. The number of unique block addresses between two consecutive accesses to the same block address is bounded above by  $k$ . What is the miss ratio if the access sequence is passed through a cache of associativity  $A \geq k$  exercising least-recently-used replacement policy?

- A.  $\left(\frac{n}{N}\right)$
- B.  $\left(\frac{1}{N}\right)$
- C.  $\left(\frac{1}{A}\right)$
- D.  $\left(\frac{k}{n}\right)$

lets take an example <https://gemini.google.com/app/ec86fc67acc14d88>

unique blocks =  $\{B_0, B_1, B_2, B_3, B_4\} n=5$

access sequence  $B_0, B_1, B_2, \underbrace{B_3}_{1 < K}, B_2, \underbrace{B_1}_{2 < K}, B_3 \quad N=7$   
 $K=3$  i.e.  $B_i, \dots, B_i$

associativity  $A=1$ , no info on set associative or fully associative. we can assume two cases:

1) set associative with all blocks mapping to same set, since no mechanism was given to calculate into which set a block goes into.

2) fully associative, a memory block can go into any line.

access      hit or miss      cache [ ] initially empty

$B_0$       miss       $[B_0]$

$B_1$       miss       $[B_1, B_0]$

$B_2$       miss       $[B_2, B_1, B_0]$

$B_3$       miss       $[B_3, B_2, B_1, B_0]$

$B_2$       hit       $[B_2, B_3, B_1, B_0]$

$B_4$       miss       $[B_4, B_2, B_3, B_1]$   $B_0$  gets evicted

$B_3$       hit       $[B_3, B_4, B_2, B_1]$

# misses = # unique block =  $n$  i.e. all are compulsory misses  
 out of  $N=7$  accesses = miss ratio =  $\frac{n}{N}$

In designing a computer's cache system, the cache block (or cache line) size is an important parameter. Which one of the following statements is correct in this context?

- A. A smaller block size implies better spatial locality
- B. A smaller block size implies a smaller cache tag and hence lower cache tag overhead
- C. A smaller block size implies a larger cache tag and hence lower cache hit time
- D. A smaller block size incurs a lower cache miss penalty

option A, B and C are not correct and D is correct.

A: for higher spacial locality, if we access a word x from cache then as many words as possible around x must be in same cache line.

B: smaller block size implies more cache lines/sets each of which has a tag, then more tag comparisons are needed so high cache tag overhead.

D: miss penalty is extra time needed to fetch and write group of words from main memory, if cache block size is small then less words needs to be fetched and written, this lower miss penalty.

#### Claim 1: "A smaller block size implies a larger cache tag"

C: This is not true under standard assumptions. A physical memory address is typically divided into three parts to find data in the cache:

[ Tag | Index | Block Offset ]

Let's see how changing the block size affects these parts, assuming the total data capacity of the cache is fixed:

1. **Block Offset:** This part identifies a specific word within a block. The number of offset bits is determined by the block size ( $\text{Offset bits} = \log_2(\text{Block Size})$ ). If you make the **block size smaller**, you need **fewer offset bits**.
2. **Index:** This part identifies which cache line (or set) the data belongs to. The number of lines is  $\text{Cache Size} / \text{Block Size}$ . If you make the **block size smaller**, the number of lines in the cache **increases**, which means you need **more index bits**.
3. **Tag:** This part is what's left over. The total number of bits in a physical address is constant.

As it turns out, the **decrease in offset bits is exactly canceled out by the increase in index bits**. This means the **size of the tag remains the same**.

#### Claim 2: "...and hence lower cache hit time"

This part of the statement is also flawed. Cache hit time is the time taken to retrieve data that is already in the cache. It depends on factors like the time to compare tags and the time to select the data from the block.

- A smaller block size means less data to choose from within a cache line, which could slightly speed up the data selection process (the final step of a hit).
- However, a smaller block size also means more cache lines, which requires a larger and potentially slower decoder to select the correct line based on the index.

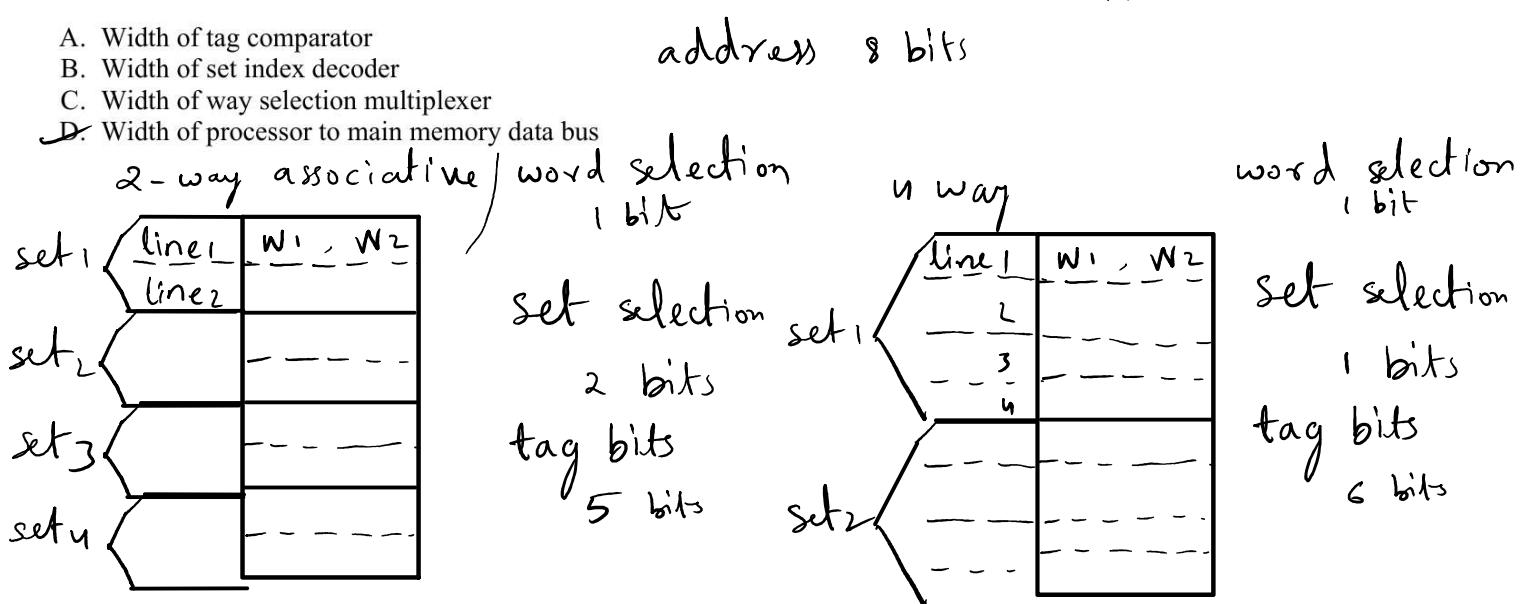
The overall effect on hit time is complex and not guaranteed to be lower. More importantly, the reasoning provided ("hence") is based on a false premise (a larger tag). A larger tag would require a larger comparator, which would likely **increase**, not decrease, the hit time.

From <<https://www.perplexity.ai/search/i-know-why-option-a-and-b-are-h6NTjb79SquZaBheWHbMrA>>

Since the **first part of the statement is false**, the entire proposition is **incorrect**.

If the associativity of a processor cache is doubled while keeping the capacity and block size unchanged, which one of the following is guaranteed to be NOT affected?

- A. Width of tag comparator
- B. Width of set index decoder
- C. Width of way selection multiplexer
- D. Width of processor to main memory data bus



way/line selection multiplexer

previously needed to route 2 ways now has to route n ways

A 4-way set-associative cache memory unit with a capacity of 16 KB is built using a block size of 8 words. The word length is 32 bits. The size of the physical address space is 4 GB. The number of bits for the TAG field is \_\_\_\_\_

$$\text{cache size} = 16 \text{ KB} = 2^4 \times 2^{10} \text{ B}$$

$$\text{line/block size} = 8 \text{ words} = 8 \times 32 \text{ bits} = 8 \times 4 \times 8 \text{ bits} \\ = 32 \text{ bytes}$$

$$\# \text{ cache lines} = \frac{\text{cache size}}{\text{line size}} = \frac{2^4 \text{ B}}{2^5 \text{ B}} = 2^4 = 16 \quad \begin{matrix} \text{B} \\ \# \text{ word offset bits} = \log_2(32) \end{matrix}$$

$$\# \text{ sets} = \frac{\# \text{ cache lines}}{\# \text{ lines per set}} = \frac{2^4}{4} = 2^2 = 4 \quad \begin{matrix} \text{B} \\ \# \text{ set index bits} = \log_2(2^2) = 2 \end{matrix}$$

$$\text{physical address space} = 4 \text{ GB} = 2^{32} \text{ B} \quad \begin{matrix} \text{B} \\ = \log_2(2^{32}) = 32 \end{matrix}$$

$$\text{i.e. main memory size} = 2^2 \times 2^{30} \text{ B} = 2^{32} \text{ B}$$

# bits needed to locate word in main memory

$$= \log_2(\text{main memory size}) = \log_2(2^{32}) = 32$$

$$\# \text{ tag bits} = 32 - \# \text{ set bits} - \# \text{ word bits} = 32 - 2 - 5 = 25$$

Consider a machine with a byte addressable main memory of  $2^{20}$  bytes, block size of 16 bytes and a direct mapped cache having  $2^{12}$  cache lines. Let the addresses of two consecutive bytes in main memory be  $(E201F)_{16}$  and  $(E2020)_{16}$ . What are the tag and cache line addresses (in hex) for main memory address  $(E201F)_{16}$ ?

~~A. E, 201~~

~~B. F, 201~~

~~C. E, E20~~

~~D. 2, 01F~~

Tag is left most / most significant bits

option B has right most / least significant bits  
D has some middle bits as tag

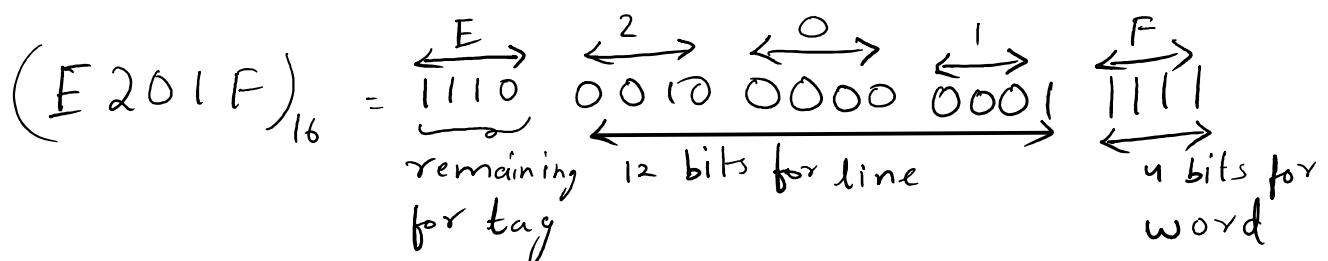
Tag bits & line bits are clearly separate with no overlap but in option C has hex value E in both tag & line parts

Systematic ways:

$$\begin{aligned}\# \text{ bits in physical memory address} &= \log_2(\text{size of memory}) \\ &= \log_2(2^{20}) = 20 \text{ bits}\end{aligned}$$

$$\begin{aligned}\text{block size} = 16 \text{ B} \Rightarrow \# \text{ word offset bits} &= \log_2(\text{block size}) \\ &= \log_2(16) = 4\end{aligned}$$

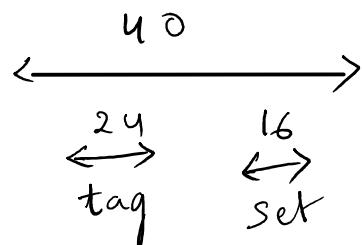
$$\begin{aligned}\# \text{ cache lines} = 2^{12} \Rightarrow \# \text{ line bits} &= \log_2(\# \text{ lines}) \\ &= \log_2(2^{12}) = 12\end{aligned}$$



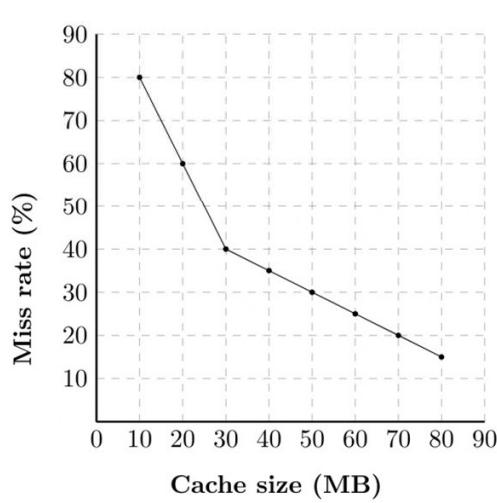
$$\therefore \text{Tag} = (E)_{16} \text{ & line} = (201)_{16}$$

The width of the physical address on a machine is 40 bits. The width of the tag field in a 512 KB 8-way set associative cache is 24 bits.

- since no info given about block/line size we will assume 1 byte - 1 word, 8 way = 1 set has 8 lines/block  
 1 block/line can hold 1 word / byte cause that is minimum # sets =  $\frac{2^{19} B}{8 B} = \frac{2^{19} B}{2^3 B} = 2^{16}$
- \* since there is only 1 word per line no need of word offset bits, just locating set is enough
  - \* if tag matches we will retrieve word from line corresponding to matched tag.



A file system uses an in-memory cache to cache disk blocks. The miss rate of the cache is shown in the figure. The latency to read a block from the cache is 1 ms and to read a block from the disk is 10 ms. Assume that the cost of checking whether a block exists in the cache is negligible. Available cache sizes are in multiples of 10 MB.



avg latency = hitratio \* cache latency + miss ratio \* memory latency  
 always strive for lowest miss rate

miss rate cache avg latency

20%	$70 \text{ MB} (0.8) 1 \text{ ms} + (0.2) 10 \text{ ms} = (0.8 + 2) \text{ ms} = 2.8$
30%	$50 \text{ MB} (0.7) 1 + (0.3) 10 = (0.7 + 3) \text{ ms} = 3.7$
40%	$30 \text{ MB} (0.6) 1 + (0.4) 10 = (0.6 + 4) \text{ ms} = 4.6$
50%	not available
60%	$20 \text{ MB} (0.4) 1 + (0.6) 10 = (0.4 + 6) \text{ ms} = 6.4$

The smallest cache size required to ensure an average read latency of less than 6 ms is 30 MB.



A cache memory unit with capacity of  $N$  words and block size of  $B$  words is to be designed. If it is designed as a direct mapped cache, the length of the TAG field is 10 bits. If the cache unit is now designed as a 16-way set-associative cache, the length of the TAG field is \_\_\_\_\_ bits.

$$\# \text{word bits} = \log_2(B) \quad \# \text{tag bits} = 10$$

$$\# \text{line bits} = \log_2(\# \text{lines}) = \log_2\left(\frac{N}{B}\right)$$

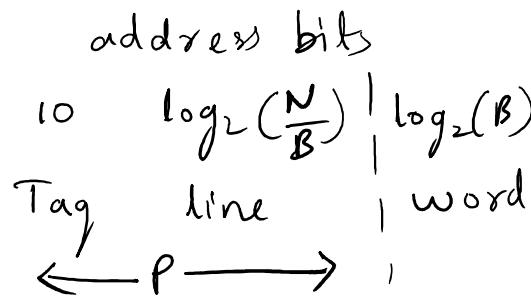
$$\# \text{sets} = \frac{\# \text{lines}}{\# \text{lines per set}} = \frac{(N/B)}{16}$$

$$\# \text{set bits} = \log_2\left(\frac{N}{16B}\right)$$

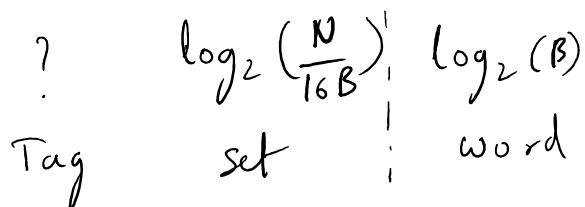
block size / line size doesn't change

So # word bits also don't change

old



new



$$P = 10 + \log_2\left(\frac{N}{B}\right) = ? + \log_2\left(\frac{N}{16B}\right) \quad \text{let } \frac{N}{B} = X$$

$$10 + \log_2(X) = ? + \log_2\left(\frac{X}{16}\right) \quad \text{we know } \log \frac{a}{b} = \log a - \log b$$

$$10 + \log_2(X) = ? + \cancel{\log_2(X)} - \log_2(16)$$

$$10 = ? - \log_2(2^4)$$

$$? = 10 + 4 = 14 \text{ i.e. no. of new Tag bits}$$

#### 1.2.47 Cache Memory: GATE CSE 2017 Set 2 | Question: 53

<https://gateoverflow.in/118613>

Consider a machine with a byte addressable main memory of  $2^{32}$  bytes divided into blocks of size 32 bytes. Assume that a direct mapped cache having 512 cache lines is used with this machine. The size of the tag field in bits is 18

$$\# \text{ word bits} = \log_2(32) = 5 \quad \# \text{ line bits} = \log_2(512) = 9$$

$$\# \text{ tag bits} = 32 - 9 - 5 = 18$$

#### 1.2.48 Cache Memory: GATE CSE 2018 | Question: 34

<https://gateoverflow.in/204108>

The size of the physical address space of a processor is  $2^P$  bytes. The word length is  $2^W$  bytes. The capacity of cache memory is  $2^N$  bytes. The size of each cache block is  $2^M$  words. For a K-way set-associative cache memory, the length (in number of bits) of the tag field is

$$\text{cache size} = 2^N \text{ Bytes}$$

$$\text{block/line size} = 2^M \text{ words} = 2^M \times 2^W \text{ Bytes} = 2^{M+W} \text{ Bytes}$$

$$\# \text{ bits to locate word in line} = \log_2(2^{M+W}) = M+W$$

$$\begin{aligned} \# \text{ lines} &= \frac{\text{cache size}}{\text{line size}} = \frac{2^N}{2^{M+W}} \\ \# \text{ sets} &= \frac{\# \text{ lines}}{\# \text{ lines per set}} = \frac{2^{N-M-W}}{K} = 2^{N-M-W} \end{aligned}$$

$$\begin{aligned} \# \text{ set bits} &= \log_2\left(\frac{2^{N-M-W}}{K}\right) = N-M-W - \log_2 K \\ \# \text{ tag bits} &= P - (N-M-W - \log_2 K) - (M+W) \\ &= P - N + M + W + \log_2 K = M - W \end{aligned}$$

### 1.2.49 Cache Memory: GATE CSE 2019 | Question: 1 top

A certain processor uses a fully associative cache of size 16 kB. The cache block size is 16 bytes. Assume that the main memory is byte addressable and uses a 32-bit address. How many bits are required for the *Tag* and the *Index* fields respectively in the addresses generated by the processor?

- A. 24 bits and 0 bits
- B. 28 bits and 4 bits
- C. 24 bits and 4 bits
- D. 28 bits and 0 bits

$$\text{cache size} = 16 \text{ kB} = 2^4 \times 2^{10} \text{ B} = 2^{14} \text{ B}$$

$$\text{block/line size} = 16 \text{ B} = 2^4 \text{ B} \Rightarrow \# \text{word/index} \\ \text{bits} = \log_2(2^4) \\ = 4$$

$$\# \text{tag bits} = 32 - 4 = 28$$

in fully associative cache no need of set bits or line bits

tag bits are enough to locate line/block in cache.

### 1.2.52 Cache Memory: GATE CSE 2020 | Question: 30 top

<https://gateoverflow.in/333201>



A computer system with a word length of 32 bits has a 16 MB byte-addressable main memory and a 64 KB, 4-way set associative cache memory with a block size of 256 bytes. Consider the following four physical addresses represented in hexadecimal notation.

- A1 = 0x42C8A4,
- A2 = 0x546888,
- A3 = 0x6A289C,
- A4 = 0x5E4880

$$\text{main memory } 16 \text{ MB} = 2^4 \times 2^{20} \text{ B} = 2^{24} \text{ B}$$

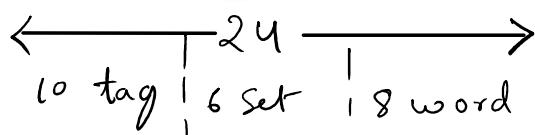
$$\text{main mem address bits} = \log_2(\text{mem size}) = \log_2(2^{24}) = 24$$

$$\text{total cache size} = 64 \text{ KB} = 2^6 \times 2^{10} \text{ B} = 2^{16} \text{ B}$$

$$\text{block/line size} = 256 \text{ bytes} = 2^8 \text{ B} \rightarrow 1 \text{ line has } 2^8 \text{ bytes or words} \\ \text{so 8 bits will locate word in line}$$

Which one of the following is TRUE?

- A. A1 and A4 are mapped to different cache sets.
- B. A2 and A3 are mapped to the same cache set.
- C. A3 and A4 are mapped to the same cache set.
- D. A1 and A3 are mapped to the same cache set.



A1    4 2 C 8    | A4

11 00 1000

A2 & A3 same set

A2    5 6 8    | 88

0110 1000

A1 & A4 same set

A3    6 A 2 8    | 9C

0010 1000

A4    5 E 4 8    | 80

01100 1000

0000	0	1000	8
0001	1	1001	9
0010	2	1010	A
0011	3	1011	B
0100	4	1100	C
0101	5	1101	D
0110	6	1110	E
0111	7	1111	F

### 1.2.53 Cache Memory: GATE CSE 2021 Set 1 | Question: 22 top

<https://gateoverflow.in/357429>

Consider a computer system with a byte-addressable primary memory of size  $2^{32}$  bytes. Assume the computer system has a direct-mapped cache of size 32 KB (1 KB =  $2^{10}$  bytes), and each cache block is of size 64 bytes.

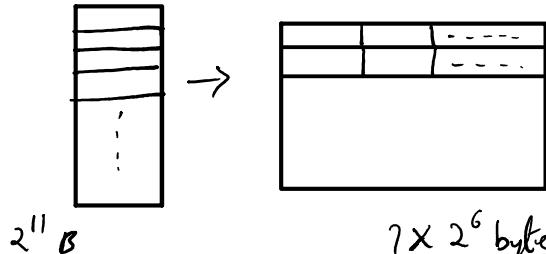
The size of the tag field is \_\_\_\_\_ bits. Tag bits =  $\text{PAS}_{\text{bits}} - \log_2(\text{Cache Size}) + \log_2(K)$  (where  $K$  is associativity),  
 $= 32 - 15 + 0 = 17$  bits



Consider a set-associative cache of size 2KB ( $1\text{KB} = 2^{10}$  bytes) with cache block size of 64 bytes. Assume that the cache is byte-addressable and a 32-bit address is used for accessing the cache. If the width of the tag field is 22 bits, the associativity of the cache is 2.  $\text{tagT} = 23 \text{ bits} - 1 = 6$        $22 = T + S + W$

$$- \text{bits for tag } T = 22 \text{ & word } W = 6 \quad 32 = T + s + w$$

$$32 = 22 + 5 + 6$$



$2^6$  bytes # set bits  $s = 4 \rightarrow \# \text{sets} = 2^4$

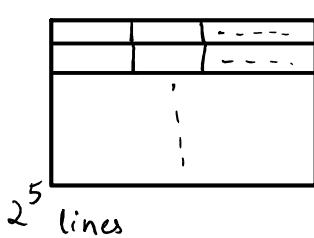
$$\text{cache size} = 2 \text{KB} = 2 \times 2^10 \text{B} = 2^{11} \text{B}$$

$$? \times 2^6 \text{ bytes} = 2^{11} \text{ bytes}$$

$$\text{no.of lines} \times \text{no.of words/bytes per line} = \begin{array}{l} \text{total no.of words} \\ \text{or byte a cache can} \\ \text{store} \end{array}$$

? =  $\frac{2^{11}}{2^6} = 2^5$  lines

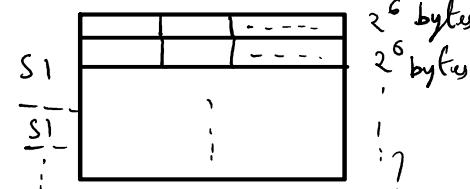
$$? = \frac{2^{11}}{3^6} = 2^5 \text{ lines}$$



$2^6$  bytes  
 $2^6$  bytes

there are  $2^u$  sets

how many lines per set



No. of sets  $\times$  No. of lines per set

= no. of total lines

$$2^4 \times ? = 2^5 \Rightarrow ? = 2 \text{ ie 2 way associative}$$

Question: 48 [top](#) <https://gateoverflow.in/3691>

Consider a fully associative cache with 8 cache blocks (numbered 0 – 7) and the following sequence of memory block requests:

4, 3, 25, 8, 19, 6, 25, 8, 16, 35, 45, 22, 8, 3, 16, 25, 7

If LRU replacement policy is used, which cache block will have memory block 7?



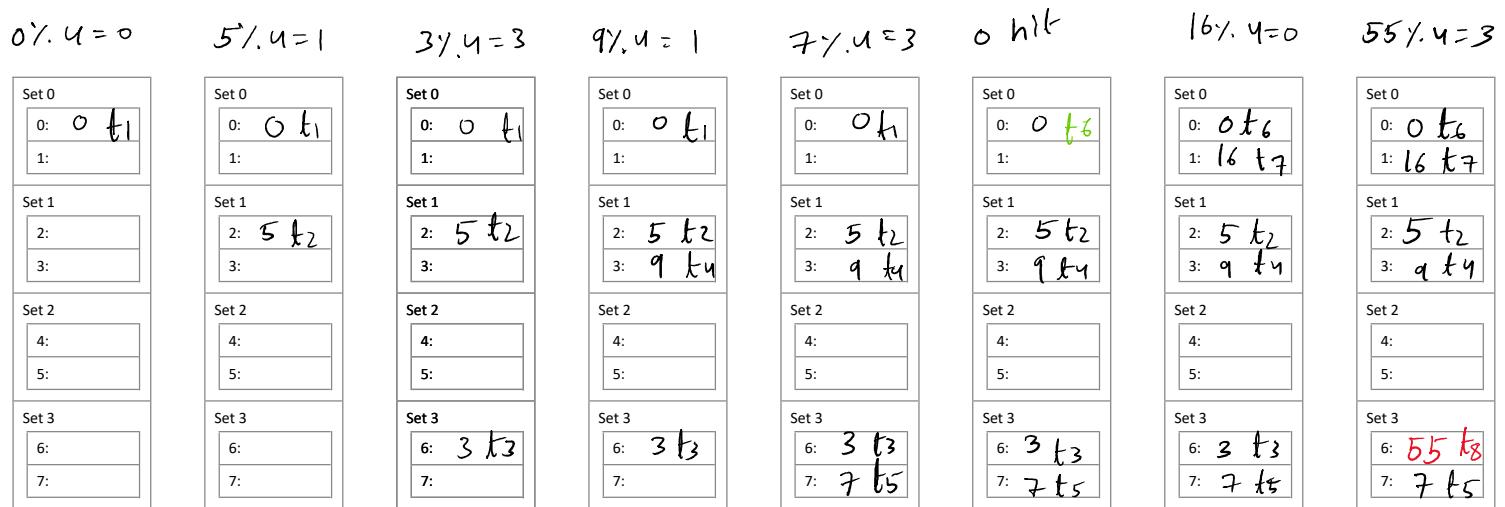
35	miss	45
miss	4	X
0: 4 t1	0: 45 t11	
1: 3 t2	1: 3 t2	
2: 25 t7	2: 25 t7	
3: 9 t8	3: 9 t8	
4: 19 t5	4: 19 t5	
5: 6 t6	5: 6 t6	
6: 16 t9	6: 16 t9	
7: 35 t10	7: 35 t10	

### 1.2.58 Cache Memory: GATE IT 2005 | Question: 61

Consider a 2-way set associative cache memory with 4 sets and total 8 cache blocks (0 – 7) and a main memory with 128 blocks (0 – 127). What memory blocks will be present in the cache after the following sequence of memory block references if LRU policy is used for cache block replacement. Assuming that initially the cache did not have any memory block from the current job?

0 5 3 9 7 0 16 55

- A. 0 3 5 7 16 55
- B. 0 3 5 7 9 16 55
- C. 0 5 7 9 16 55
- D. 3 5 7 9 16 55



### 1.2.61 Cache Memory: GATE IT 2007 | Question: 37

<https://gateoverflow.in/3470>

Consider a Direct Mapped Cache with 8 cache blocks (numbered 0 – 7). If the memory block requests are in the following order

3, 5, 2, 8, 0, 63, 9, 16, 20, 17, 25, 18, 30, 24, 2, 63, 5, 82, 17, 24.

Which of the following memory blocks will not be in the cache at the end of the sequence ?

- A. 3
- B. 18
- C. 20
- D. 30

time →

3, 5, 2, 8, 0, 63, 9, 16, 20, 17, 25, 18, 30, 24, 2, 63, 5, 82, 17, 24.

line: 3 5 2 0 0 7 1 0 4 1 1 2 6 0 2 7 5 2 1 0  
1/8

### 1.2.62 Cache Memory: GATE IT 2008 | Question: 80

<https://gateoverflow.in/3403>

Consider a computer with a 4-ways set-associative mapped cache of the following characteristics: a total of 1 MB of main memory, a word size of 1 byte, a block size of 128 words and a cache size of 8 KB.  $\Rightarrow 2^3 \times 2^{10} B = 2^{13}$

The number of bits in the TAG, SET and WORD fields, respectively are:

- A. 7, 6, 7
- B. 8, 5, 7
- C. 8, 6, 6
- D. 9, 4, 7

$$\# \text{Sets} = \frac{\# \text{lines}}{\# \text{lines per set}} = \frac{\frac{\text{cache size}}{\text{block/line size}}}{\# \text{lines per set}} = \frac{2^{13} B}{2^7 B} = \frac{2^6}{2^2} = 2^4$$

$$\# \text{set bits} = \log_2 (\# \text{sets}) = \log_2 (2^4) = 4$$

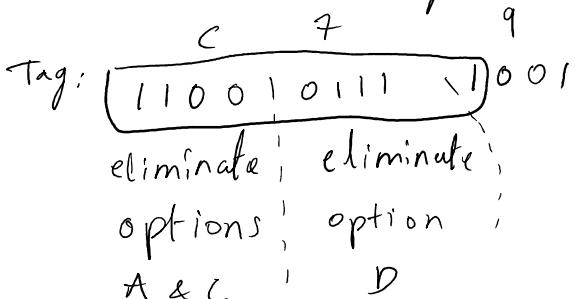
Consider a computer with a 4-way set-associative mapped cache of the following characteristics: a total of 1 MB of main memory, a word size of 1 byte, a block size of 128 words and a cache size of 8 KB. While accessing the memory location  $0C795H$  by the CPU, the contents of the TAG field of the corresponding cache line is:

- X 000011000
- B. 110001111
- X 000110000
- R 110010101

$$\# \text{sets} = \frac{\# \text{lines}}{\# \text{lines per set}} = \frac{\frac{\text{cache size}}{\text{block/line size}}}{\# \text{lines per set}} = \frac{2^{13}}{2^7} = \frac{2^6}{2^2} = 2^4$$

$$\# \text{set bits} = \log_2 (\# \text{sets}) = \log_2 (2^4) = 4$$

$$\# \text{tag bits} = 20 - 4 - 7 = 9$$



0000	0	1000	8
0001	1	1001	9
0010	2	1010	A
0011	3	1011	B
0100	4	1100	C
0101	5	1101	D
0110	6	1110	E
0111	7	1111	F

access time basics: lets say we made 10 memory access cache-line time taken for access i =  $t_{a_i}$   $\xrightarrow{\# \text{M}} \text{hit or miss}$  main mem - 1 ms

time taken for access 1 if hit =  $t_{a_1}$ , time taken for access 2 if miss =  $t_{a_2}$

total time for all memory access 7 hits 3 Miss

$$= t_{a_1} + t_{a_2} + t_{a_3} + t_{a_4} + t_{a_5} + t_{a_6} + t_{a_7} + t_{a_8} + t_{a_9} + t_{a_{10}}$$

$$= 1 \text{ ns} + 1 \text{ ns}$$

$$+ 1 \text{ ms} + 1 \text{ ms} + 1 \text{ ms} + 1 \text{ ms} + 1 \text{ ms}$$

$$+ 1 \text{ ns} + 1 \text{ ns} + 1 \text{ ns} + 1 \text{ ns}$$

$$= 7 \times 1 \text{ ns} + 3 \times (1 \text{ ns} + 1 \text{ ms} + 1 \text{ ns})$$

↑ 1st access cache and (1 ns)  
no. of hits cache no. of miss  
access time on hit

update cache happens implicitly so  
its time is largely ignored

Miss penalty { then read memory (1 ms)  
and update cache (1 ns)

average time taken for all memory access : total No. of access

$$\frac{7 \times 1 \text{ ns} + 3 \times (1 \text{ ns} + 1 \text{ ms} + 1 \text{ ns})}{10} = \frac{7 \times 1 \text{ ns} + 3 \times (1 \text{ ns} + 1 \text{ ms} + 1 \text{ ns})}{10}$$

hit ratio miss ratio



In the three-level memory hierarchy shown in the following table,  $p_i$  denotes the probability that an access request will refer to  $M_i$ .

Hierarchy Level ( $M_i$ )	Access Time ( $t_i$ )	Probability of Access ( $p_i$ )	Page Transfer Time ( $T_i$ )
$M_1$	$10^{-6}$	0.99000	0.001 sec
$M_2$	$10^{-5}$	0.00998	0.1 sec
$M_3$	$10^{-4}$	0.00002	---

If a miss occurs at level  $M_i$ , a page transfer occurs from  $M_{i+1}$  to  $M_i$  and the average time required for such a page swap is  $T_i$ . Calculate the average time  $t_A$  required for a processor to read one word from this memory system.

$$\begin{aligned}
 t_a = & \underbrace{p_1}_{\text{Chance the word is in M1}} \times \underbrace{(t_1)}_{\text{Time to read from M1}} \\
 + & \underbrace{p_2}_{\text{Chance the word is in M2}} \times \left( \underbrace{\begin{array}{c} 1. \text{ Check M1} \\ (\text{it's a miss}) \end{array}}_{t_1} + \underbrace{\begin{array}{c} 2. \text{ Read from M2} \\ (\text{it's a hit}) \end{array}}_{t_2} + \underbrace{\begin{array}{c} 3. \text{ Update M1 with} \\ \text{data from M2} \end{array}}_{T_1} \right) \\
 & \qquad \qquad \qquad \text{Total time cost when the word is found in M2} \\
 \times & \left( \underbrace{\begin{array}{c} 1. \text{ Check M1} \\ (\text{miss}) \end{array}}_{t_1} + \underbrace{\begin{array}{c} 2. \text{ Check M2} \\ (\text{miss}) \end{array}}_{t_2} + \underbrace{\begin{array}{c} 3. \text{ Read from M3} \\ (\text{hit}) \end{array}}_{t_3} + \underbrace{\begin{array}{c} 4. \text{ Update M2 with} \\ \text{data from M3} \end{array}}_{T_2} + \underbrace{\begin{array}{c} 5. \text{ Update M1 with} \\ \text{data from M2} \end{array}}_{T_1} \right) \\
 & \qquad \qquad \qquad \text{Total time cost when the word is found in M3}
 \end{aligned}$$

## Step 1: Identify the Three Possible Scenarios

#### Scenario 1: Hit at M<sub>1</sub>

- Probability:  $p_1 = 0.99000$
  - Time required:  $t_1 = 10^{-6}$  sec

#### Scenario 2: Miss at M<sub>1</sub>, Hit at M<sub>2</sub>

- Probability:  $p_2 = 0.00998$
  - Time required:  $t_1 + T_1 + t_2$
  - Calculation:  $10^{-6} + 0.001 + 10^{-5} = 0.001011 \text{ sec}$

### Scenario 3: Miss at $M_1$ and $M_2$ , Hit at $M_3$

- Probability:  $p_3 = 0.00002$
  - Time required:  $t_1 + T_1 + t_2 + T_2 + t_3$
  - Calculation:  $10^{-6} + 0.001 + 10^{-5} + 0.1 + 10^{-4} = 0.101111$  seconds

#### Step 2: Apply the Average Access Time Formula

The **average access time**  $t_a$  is calculated using:

$$t_a = p_1 \times t_1 + p_2 \times (t_1 + T_1 + t_2) + p_3 \times (t_1 + T_1 + t_2 + T_2 + t_3)$$

### Step 3: Calculate Each Term

$$\text{Term 1: } 0.99000 \times 10^{-6} = 9.9 \times 10^{-7} \text{ sec}$$

$$\text{Term 2: } 0.00998 \times 0.001011 = 1.009 \times 10^{-5} \text{ sec}$$

**Term 3:**  $0.00002 \times 0.10111 = 2.022 \times 10^{-6}$  sec

#### Step 4: Sum All Terms

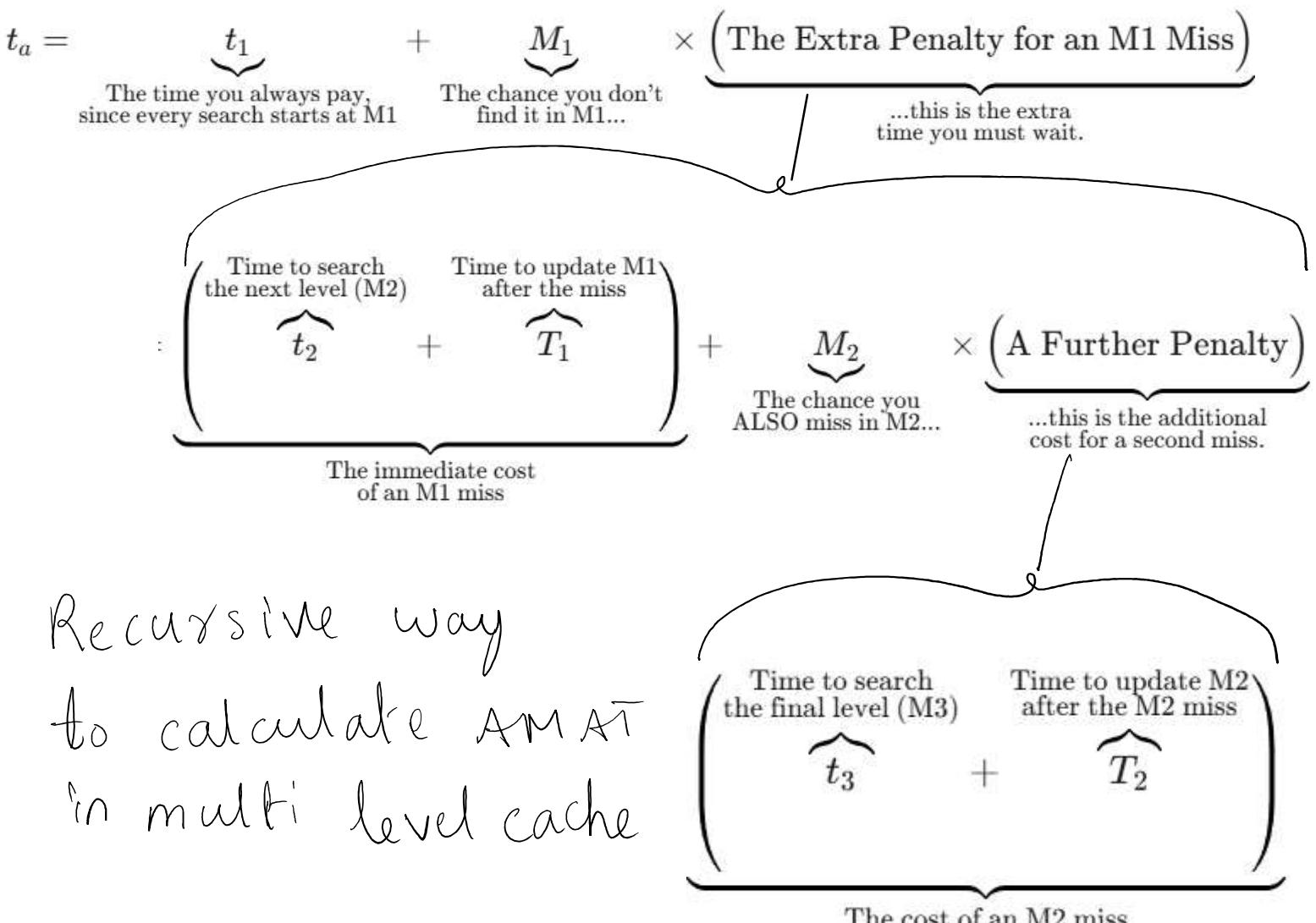
$$t_a = 9.9 \times 10^{-7} + 1.009 \times 10^{-5} + 2.022 \times 10^{-6}$$

Converting to the same power of 10:

$$t_a = 0.99 \times 10^{-6} + 10.09 \times 10^{-6} + 2.022 \times 10^{-6}$$

$$= (2.22 + 12.22 - 2.22) = 12.22 = 12.22 - 12.22$$

$$t_a = (0.99 + 10.09 + 2.022) \times 10^{-6} = 13.102 \times 10^{-6} \text{ sec}$$



Recursive way  
to calculate AMAT  
in multi level cache

#### 1.2.8 Cache Memory: GATE CSE 1996 | Question: 26

<https://gateoverflow.in/2778>

A computer system has a three-level memory hierarchy, with access time and hit ratios as shown below:

Level 1 (Cache memory) Access time = 50nsec/byte

Level 2 (Main memory) Access time = 200nsec/byte

Size	Hit ratio
8M bytes	0.80
16M bytes	0.90
64M bytes	0.95

Size	Hit ratio
4M bytes	0.98
16M bytes	0.99
64M bytes	0.995

Size	Hit ratio
260M bytes	1.0

- A. What should be the minimum sizes of level 1 and 2 memories to achieve an average access time of less than 100nsec?
- B. What is the average access time achieved using the chosen sizes of level 1 and level 2 memories?



For a set-associative Cache organization, the parameters are as follows:

$t_c$	Cache Access Time
$t_m$	Main memory access time
$l$	Number of sets
$b$	Block size
$k \times b$	Set size

Calculate the hit ratio for a loop executed 100 times where the size of the loop is  $n \times b$ , and  $n = k \times m$  is a non-zero integer and  $1 \leq m \leq l$ .

Give the value of the hit ratio for  $l = 1$ .

#### Illustrative walk-through (concrete numbers)

In the relation  $n = k \times m$

- $k$  = lines (blocks) per set
- $n$  = total blocks touched by the loop
- $l$  = total sets in the cache
- $m$  is chosen so that  $1 \leq m \leq l$

Hence  $m$  counts how many cache sets the loop's data occupy.

Why this is the right interpretation:

- One set can hold exactly  $k$  blocks.
- If the loop needs  $n$  blocks and  $n = k \times m$ , it needs  $m$  whole sets' worth of space.
- The bound  $m \leq l$  ensures the loop's footprint never references more sets than the cache actually has.
- When  $l = 1$ ,  $m$  must be 1, meaning the entire loop fits in that single set—this is the special case you computed earlier.

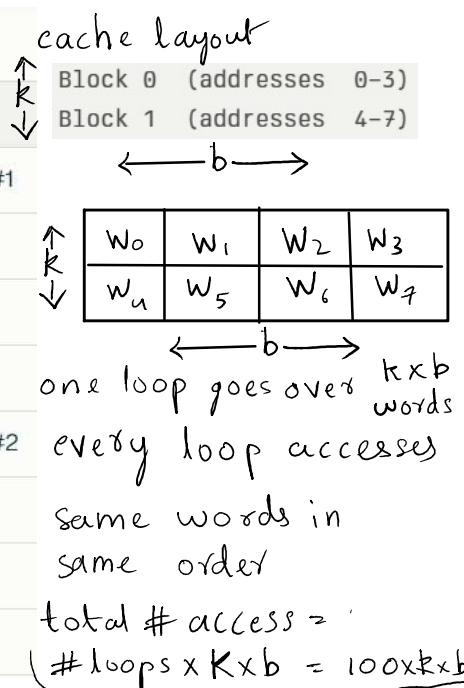
So, practically,  $m$  = number of distinct sets accessed by the loop.

To keep everything small and visible, pick the following values that respect the given constraints  $l = 1$  (one set) and  $n = k$ :

Parameter	Chosen value	Reason
Block size $b$	4 words	Small enough to list every access
Lines per set $k$	2	Gives one 2-line set (the whole cache)
Number of sets $l$	1	Matches the case we solved earlier
"Loop size" $n$	$k = 2$ blocks	So the loop touches exactly 2 blocks = 8 words
# iterations	100	Required by the problem

#### First pass inside the loop (iteration 1)

Access	Hit / miss	Action in cache	Comments
Word 0	Miss	Load Block 0 into line 0.	Compulsory miss #1
Word 1	Hit	—	Still in Block 0
Word 2	Hit	—	
Word 3	Hit	—	
Word 4	Miss	Load Block 1 into line 1.	Compulsory miss #2
Word 5	Hit	—	Still in Block 1
Word 6	Hit	—	
Word 7	Hit	—	



only in 1st iteration there misses, since we access iterating over same words and no eviction happened thus from next loop onwards all are hits only

$$\text{Hit ratio} = \frac{\text{\# times accessed word was found in cache } (N_H)}{\text{total no. of accesses } (N_A)} = \frac{N_A - N_M}{N_A} = \frac{800 - 2}{800} = \frac{798}{800} = \frac{399}{400}$$

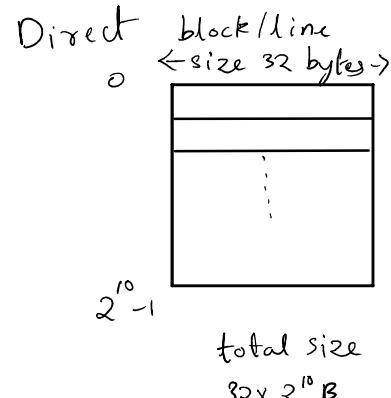
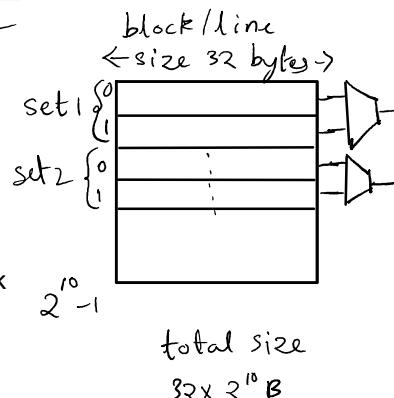
so we got  $1 - \frac{2}{800}$  i.e.  $1 - \frac{\text{\# misses}}{\text{\# accesses}} = 1 - \frac{N_M}{100 \times k \times b}$  here  $N_M$  was equal to 2 since all compulsory miss =  $1 - \frac{1}{100b}$



$2^5 \times 2^{10} B$   
 Consider two cache organizations. First one is 32 KB 2-way set associative with 32 byte block size, the second is of same size but direct mapped. The size of an address is 32 bits in both cases. A 2-to-1 multiplexer has latency of 0.6 ns while a  $k$ -bit comparator has latency of  $\frac{k}{10}$  ns. The hit latency of the set associative organization is  $h_1$  while that of direct mapped is  $h_2$ .

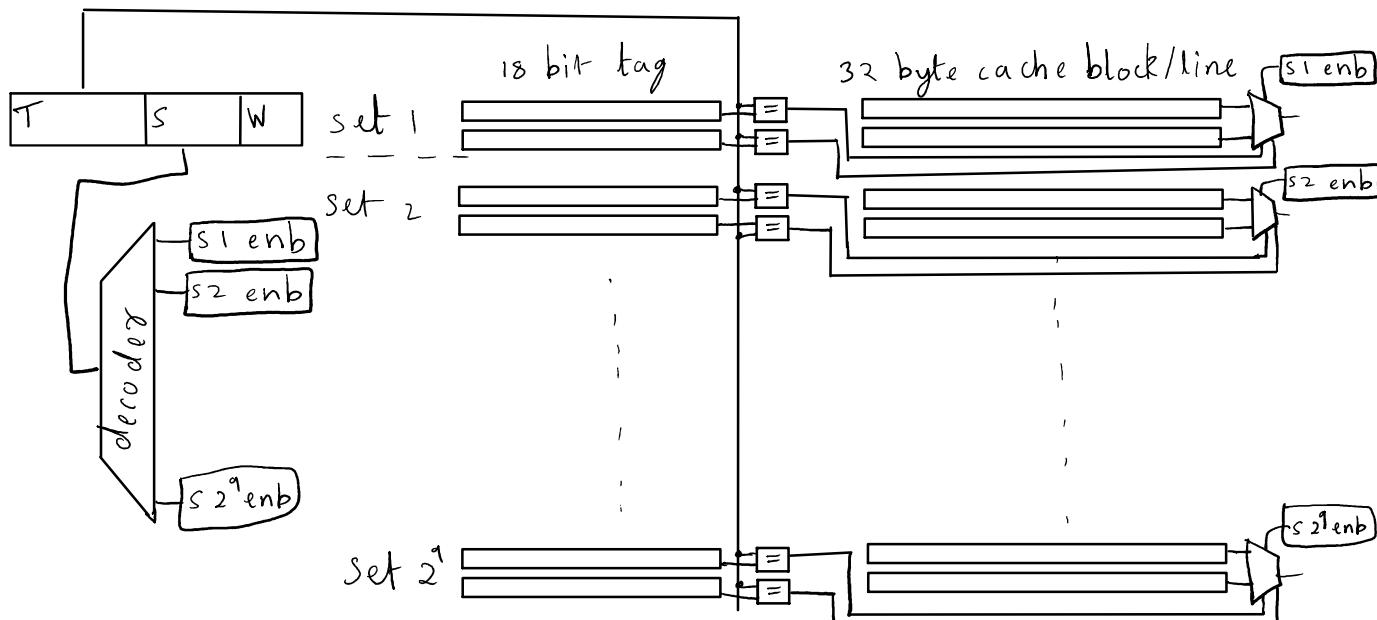
The value of  $h_1$  is: set

- A. 2.4 ns
- B. 2.3 ns
- C. 1.8 ns
- D. 1.7 ns



$$\begin{aligned} \# \text{sets} &= \frac{\# \text{lines}}{\text{lines per set}} & W_s &= \log_2(32) = 5 \\ &= \frac{2^{10}}{2} = 2^9 & T_s &= 32 - 9 - 5 \\ S &= \log_2(2^9) = 9 & &= 18 \end{aligned}$$

$$\begin{aligned} W_d &= 5 \\ I &= \log_2(\# \text{lines}) = \log_2(2^{10}) = 10 \\ T_d &= 32 - 10 - 5 = 17 \end{aligned}$$



set associative  
comparator latency =  $\frac{18}{10} \text{ ns} = 1.8 \text{ ns}$

way selection mux latency = 0.6 ns

hit latency =  $1.8 + 0.6 = 2.4 \text{ ns}$

The value of  $h_2$  is: in direct map

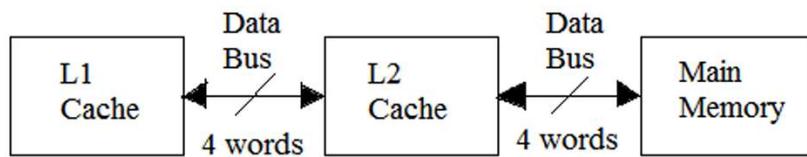
- A. 2.4 ns
  - B. 2.3 ns
  - C. 1.8 ns
  - D. 1.7 ns
- there is no way selection mux  
direct map comparator latency

within a set combined result from all comparators acts as select bits for that set's mux  
another mux will select wanted word from matched block

$$> \frac{17}{10} = 1.7 \text{ ns}$$



A computer system has an  $L_1$  cache, an  $L_2$  cache, and a main memory unit connected as shown below. The block size in  $L_1$  cache is 4 words. The block size in  $L_2$  cache is 16 words. The memory access times are 2 nanoseconds, 20 nanoseconds and 200 nanoseconds for  $L_1$  cache,  $L_2$  cache and the main memory unit respectively.

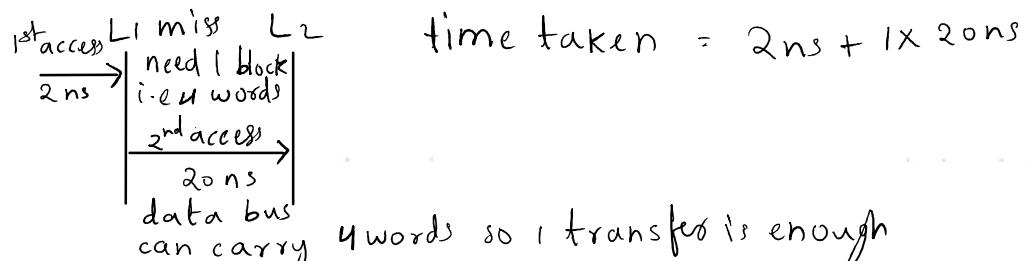


tests.gatecse.in

tests.gatecse.in

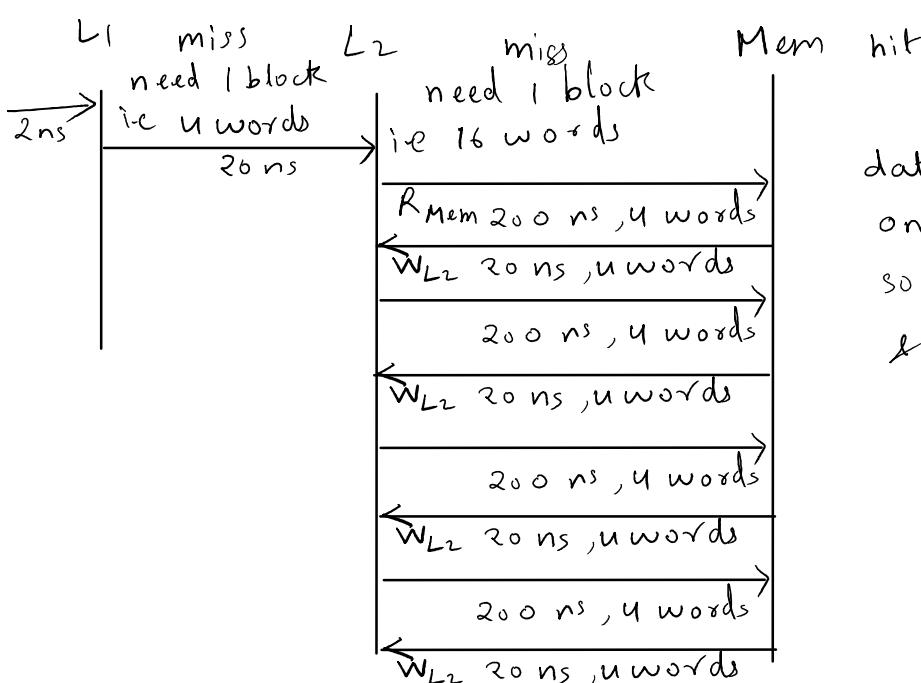
When there is a miss in  $L_1$  cache and a hit in  $L_2$  cache, a block is transferred from  $L_2$  cache to  $L_1$  cache. What is the time taken for this transfer?

- A. 2 nanoseconds
- B. 20 nanoseconds
- C. 22 nanoseconds
- D. 88 nanoseconds



When there is a miss in both  $L_1$  cache and  $L_2$  cache, first a block is transferred from main memory to  $L_2$  cache, and then a block is transferred from  $L_2$  cache to  $L_1$  cache. What is the total time taken for these transfers?

- A. 222 nanoseconds
- B. 888 nanoseconds
- C. 902 nanoseconds
- D. 968 nanoseconds



data bus can transfer  
only 4 words at a time  
so need to read memory u times  
& write to  $L_2$  u times

$$\begin{aligned}
 & 4 \times (200 + 20) \\
 & = 4 \times 220 \\
 & = 880
 \end{aligned}$$

total access time : CPU to  $L_1$  +  $L_1$  to  $L_2$  +  $L_2$  to Mem

$$= 2 + 20 + 880$$

$$= 902$$





In a two-level cache system, the access times of  $L_1$  and  $L_2$  caches are 1 and 8 clock cycles, respectively. The miss penalty from the  $L_2$  cache to main memory is 18 clock cycles. The miss rate of  $L_1$  cache is twice that of  $L_2$ . The average memory access time (AMAT) of this cache system is 2 cycles. The miss rates of  $L_1$  and  $L_2$  respectively are

- L<sub>1</sub> 2 M tests.gatecse.in goclasses.in tests.gatecse.in
- A. 0.111 and 0.056  
 B. 0.056 and 0.111  
 C. 0.0892 and 0.1784  
 D. 0.1784 and 0.0892

clock cycles is a measure of time

$L_1 \quad 1 \text{ cc}$

let  $\Rightarrow$   
 miss rate of  $L_2$

be  $M$

$L_2 \quad 8 \text{ cc} \quad L_2 \rightarrow \text{mem } 18 \text{ cc}$

then miss rate of  $L_1$   
 will be  $2M$

$$2 \text{ cc} = \underbrace{1 \text{ cc}}_{L_1 \text{ access time}} + \underbrace{2M \times \left( \underbrace{8 \text{ cc}}_{L_2 \text{ access time}} + \underbrace{M \times 18 \text{ cc}}_{L_2 \rightarrow \text{mem miss penalty}} \right)}_{L_1 \text{ miss rate}}$$

in options B

$$2 \text{ cc} = 1 \text{ cc} + 2M \times 8 \text{ cc} + 2 \times M^2 \times 18 \text{ cc}$$

if  $M = 0.111$  then

$$2M = 0.222 \text{ not } 0.056$$

$$2 \text{ cc} = 1 \text{ cc} + 16M \text{ cc} + 36M^2 \text{ cc}$$

∴ option C so

$$\sigma = (36M^2 + 16M - 1) \text{ cc}$$

eliminate both.

option A substitute  $M = 0.056$

option D

$$\frac{36 * (\text{sqr}(0.056)) + 16 * 0.056 - 1}{0.008896000}$$

$$\frac{36 * (\text{sqr}(0.0892)) + 16 * 0.0892 - 1}{0.7136390}$$

↳ very close to  $\sigma$

↳ not close to  $\sigma$

$$\therefore M = 0.056$$

$$2M = 0.111$$

option A ✓

Two ways to calculate AMAT

①

hit ratio  $\times$  cache access time + miss ratio  $\times$   $(\underbrace{\text{cache access time} + \text{main memory time}}_{\text{memory time}})$

=) cache access time  $\times$   $(\underbrace{\text{hit ratio} + \text{miss ratio}}_{1}) + \text{miss ratio} \times \text{memory time}$  i.e extra time

∴ cache access time + miss ratio  $\times$  Miss penalty  $\xrightarrow[③ \text{ upon miss}]{}$   
 used this form since miss penalty info was given directly in question



The read access times and the hit ratios for different caches in a memory hierarchy are as given below:

Cache	Read access time (in nanoseconds)	Hit ratio
I-cache	2	0.8
D-cache	2	0.9
L2-cache	8	0.9

since hit ratio is given lets use Weighted form of AMAT instead of penalty AMAT

tests.gatecse.in

goclasses.in

tests.gatecse.in

The read access time of main memory is 90 nanoseconds. Assume that the caches use the referred-word-first read policy and the write-back policy. Assume that all the caches are direct mapped caches. Assume that the dirty bit is always 0 for all the blocks in the caches. In execution of a program, 60% of memory reads are for instruction fetch and 40% are for memory operand fetch. The average read access time in nanoseconds (up to 2 decimal places) is \_\_\_\_\_

update only  
cache & not  
main memory

avg read access time

$$= \text{Hit ratio} \times \text{read access time} + \text{Miss ratio} \times (\text{read access time} + \text{larger memory access time})$$

or

$$= \text{read access time} + \text{Miss ratio} \times \text{larger memory}$$

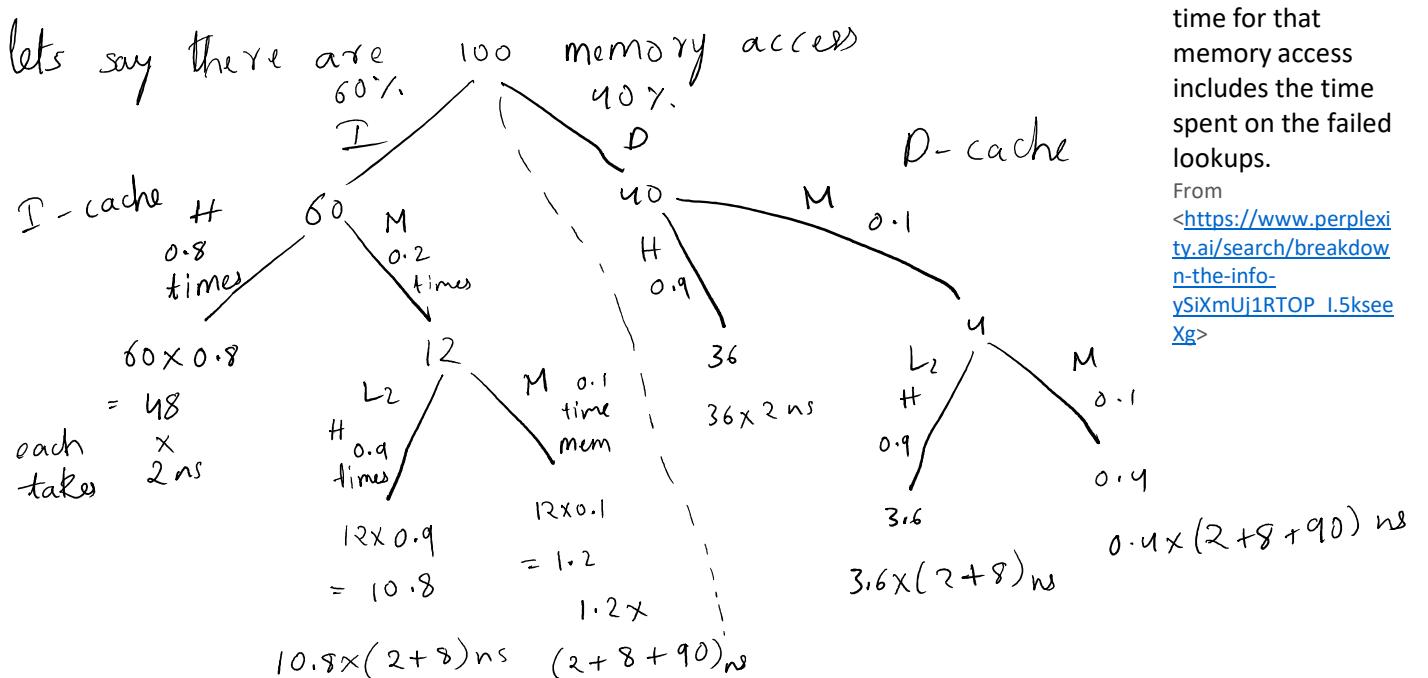
**referred word first** policy is a technique where,

without "referred word first" policy, on a miss data word from memory to cache and then to processor, two sequential transfers with "referred word first" policy, on a miss 1 data word from memory to processor 2 data word from memory to cache 1 and 2 happens simultaneously

From <<https://www.perplexity.ai/search/in-caches-what-is-referred-word-7afutkMrSYyaAxe0s0gv8Q>>

When a cache miss occurs, the total time for that memory access includes the time spent on the failed lookups.

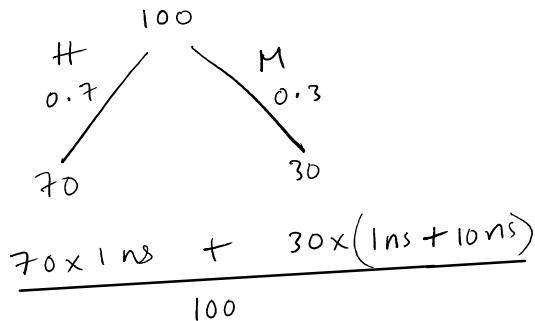
From <[https://www.perplexity.ai/search/breakdown-the-info-ySiXmUi1RTOP\\_1.5kseeXg](https://www.perplexity.ai/search/breakdown-the-info-ySiXmUi1RTOP_1.5kseeXg)>



$$\begin{aligned}
 &= 48 \times 2 \text{ ns} + 10.8 \times 10.8 \text{ ns} + 1.2 \times 100 \text{ ns} + 36 \times 2 \text{ ns} + 3.6 \times 10 \text{ ns} + 0.4 \times 100 \text{ ns} \\
 &= \frac{96 + 108 + 120 + 72 + 36 + 40}{100} = \frac{300 + 172}{100} = \frac{472}{100} = 4.72 \text{ ns}
 \end{aligned}$$

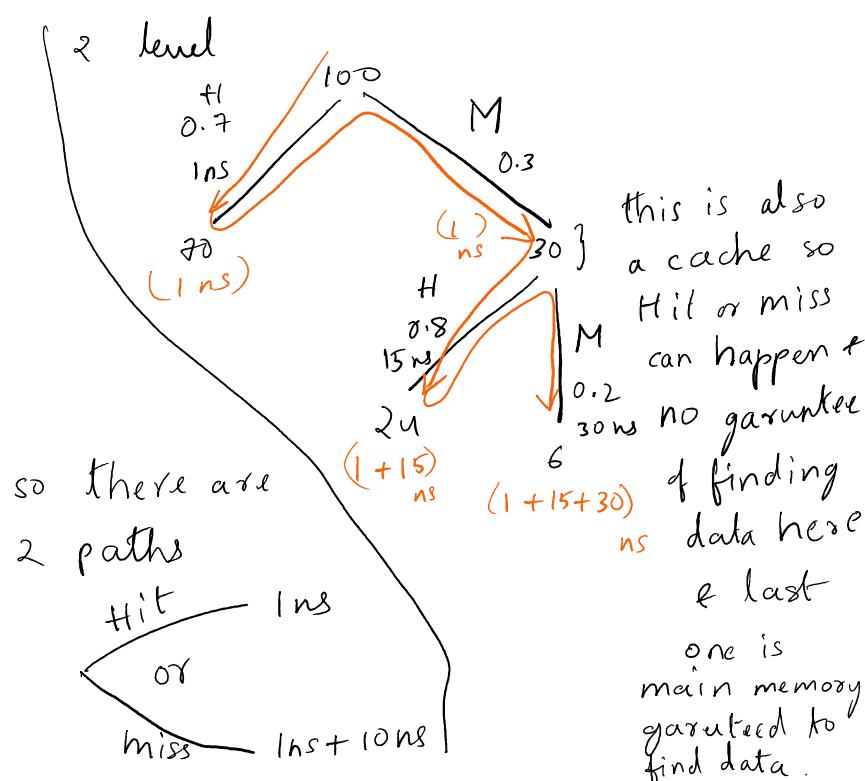
# Recap of AMAT

1 level



access cache 1ns

if hit return  
to CPU  
else if miss  
go to next  
memory level  
i.e. 10ns



## 1.2.50 Cache Memory: GATE CSE 2019 | Question: 45

<https://gateoverflow.in/302803>



A certain processor deploys a single-level cache. The cache block size is 8 words and the word size is 4 bytes. The memory system uses a 60-MHz clock. To service a cache miss, the memory controller first takes 1 cycle to accept the starting address of the block, it then takes 3 cycles to fetch all the eight words of the block, and finally transmits the words of the requested block at the rate of 1 word per cycle. The maximum bandwidth for the memory system when the program running on the processor issues a series of read operations is  $\underline{\quad} \times 10^6$  bytes/sec

$$60 \text{ MHz} \quad 60 \times 10^6 \text{ Hz} \Rightarrow 60 \times 10^6 \text{ clock cycles per 1 second}$$

block 8 words, 1 word - 4 bytes

1 cycle to accept starting address

3 cycles to fetch a block i.e. 8 words

& each word takes 1 cycles to transmit

no. of cycles to get & send a block = 1 cycle + 3 cycle + 8 words each 1 cycle = 12 cycles

so in 12 cycle 1 block i.e. 8 words ie  $8 \times 4$  bytes 32 bytes are fetched & transmitted.

$$12 \text{ cycle} - 32 \text{ bytes}$$

$$\frac{5}{60 \times 10^6 + 2} = \frac{?}{32} \Rightarrow ? = 160 \times 10^6 \text{ bytes}$$

$60 \times 10^6$  cycles - ? bytes in 1 second  $60 \times 10^6$  cycles happens

so in 1 second  $160 \times 10^6$  bytes gets transferred



A direct mapped cache memory of 1 MB has a block size of 256 bytes. The cache has an access time of 3 ns and a hit rate of 94%. During a cache miss, it takes 20 ns to bring the first word of a block from the main memory, while each subsequent word takes 5 ns. The word size is 64 bits. The average memory access time in ns (round off to 1 decimal place) is \_\_\_\_\_.

word - 8 bytes

block - 256 bytes

block -  $\frac{256}{8}$  words

line - block - 32 words

given, in a block

first word takes 20 ns

& remaining each word takes 5 ns

so in a block of 32 words first word takes 20 ns

remaining 31 takes 5 ns each

$$\text{i.e } 1 \times 20 \text{ ns} + 31 \times 5 \text{ ns} = 20 + 155 = 175 \text{ ns}$$

to bring 1 block from memory it takes 175 ns

$$\text{AMAT} \rightarrow 3 \text{ ns} + 0.06 \times 175 \text{ ns} = 3 \text{ ns} + 10.5 \text{ ns} = 13.5 \text{ ns}$$



Assume a two-level inclusive cache hierarchy, L1 and L2, where L2 is the larger of the two. Consider the following statements.

- $S_1$ : Read misses in a write through L1 cache do not result in writebacks of dirty lines to the L2
- $S_2$ : Write allocate policy *must* be used in conjunction with write through caches and no-write allocate policy is used with writeback caches.

Which of the following statements is correct?

- A.  $S_1$  is true and  $S_2$  is false
- B.  $S_1$  is false and  $S_2$  is true
- C.  $S_1$  is true and  $S_2$  is true
- D.  $S_1$  is false and  $S_2$  is false

A dirty line in cache is a block of cached data that has been modified by the CPU but not yet written back to main memory, marked by a dirty bit to ensure its synchronization later for data consistency.

From <<https://www.perplexity.ai/search/what-is-a-dirty-line-in-cache-EiwwztX6RbCAKXXc.n0otQ>>

tests.gatecse.in

#### 1. Cache Write Policies

These govern how data writes are handled between the cache and main memory 1 2 .

- **Write-Through:** Data is written to both cache and main memory at the same time for immediate consistency 1 2 3 .
- **Write-Back** (also called Write-Behind): Data is written only to the cache initially, with main memory updated later (e.g., on eviction) for better performance 1 2 3 .
- **Write-Around:** Writes bypass the cache and go directly to main memory, reducing cache pollution from rarely read data 1 2 .
- **Write-Allocate** (or Fetch on Write): On a write miss, the block is loaded into the cache before writing 2 4 .
- **No-Write-Allocate:** On a write miss, data is written directly to memory without loading it into the cache 2 4 .

**Combinations:**

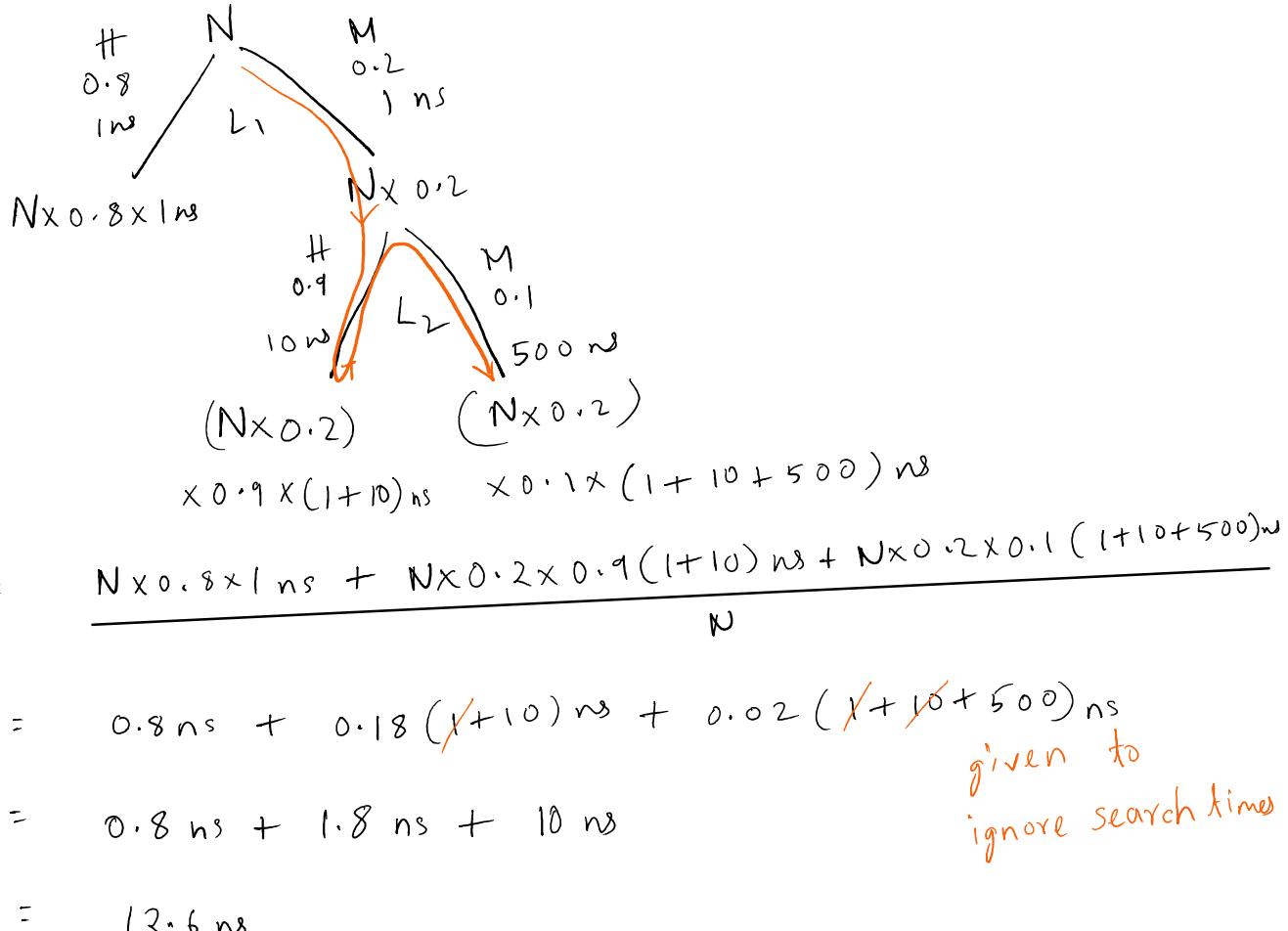
- Write-Through with Write-Allocate: Writes update both cache and memory and cache block is loaded on write miss.
- Write-Back with Write-Allocate: Writes update cache and defer memory update; cache block loaded on write miss.
- Write-Through with No-Write-Allocate: Writes update both cache and memory but block not loaded on write miss.
- Write-Back with No-Write-Allocate: Writes update cache, memory updated on eviction; write miss writes directly to memory without loading cache block.

Revisit



Consider a system with 2 level cache. Access times of Level 1 cache, Level 2 cache and main memory are  $1\text{ ns}$ ,  $10\text{ ns}$ , and  $500\text{ ns}$  respectively. The hit rates of Level 1 and Level 2 caches are 0.8 and 0.9, respectively. What is the average access time of the system ignoring the search time within the cache?

- A. 13.0
- B. 12.8
- C. 12.6
- D. 12.4



*possibly due to Hardware components like wires, etc*

A cache line is 64 bytes. The main memory has latency  $32\text{ ns}$  and bandwidth  $1\text{ GBytes/s}$ . The time required to fetch the entire cache line from the main memory is:

- A.  $32\text{ ns}$
- B.  $64\text{ ns}$
- C.  $96\text{ ns}$
- D.  $128\text{ ns}$

so it takes  $64\text{ ns}$  to transfer 64 bytes

$$64\text{ ns} + 32\text{ ns} = 96\text{ ns}$$

$1 \times 10^9 \text{ bytes per second}$

$10^9 \text{ bytes per } 10^9 \text{ ns}$

$\Rightarrow 1 \text{ byte} - 1\text{ ns}$

we need to transfer 64 bytes  
64 bytes -  $64\text{ ns}$



A computer system has a level-1 instruction cache (I-cache), a level-1 data cache (D-cache) and a level-2 cache (L2-cache) with the following specifications:

# set	# lines		Capacity	Mapping Method	Block Size
-	1K	I-Cache	4K words	Direct mapping	4 words
1K/2	1K	D-Cache	4K words	2-way set associative mapping	4 words
4K/4 = 1K	4K	L2-Cache	64K words	4-way set associative mapping	16 words

↳ every line will have a tag, after knowing  
# lines eliminate option C & D

The length of the physical address of a word in the main memory is 30 bits. The capacity of the tag memory in the I-cache, D-cache and L2-cache is, respectively,

- $\frac{1}{I}$        $\frac{1}{D}$        $\frac{1}{L_2}$
- A. 1 K x 18-bit, 1 K x 19-bit, 4 K x 16-bit
  - B. 1 K x 16-bit, 1 K x 19-bit, 4 K x 18-bit
  - C. 1 K x 16-bit, 512 x 18-bit, 1 K x 16-bit
  - D. 1 K x 18-bit, 512 x 18-bit, 1 K x 18-bit

goclasses.in

Word bits	T + I or T + S
2 <sup>10</sup>	2 <sup>8</sup>
2 <sup>10</sup>	2 <sup>8</sup>
4	2 <sup>6</sup>

goclasses.in

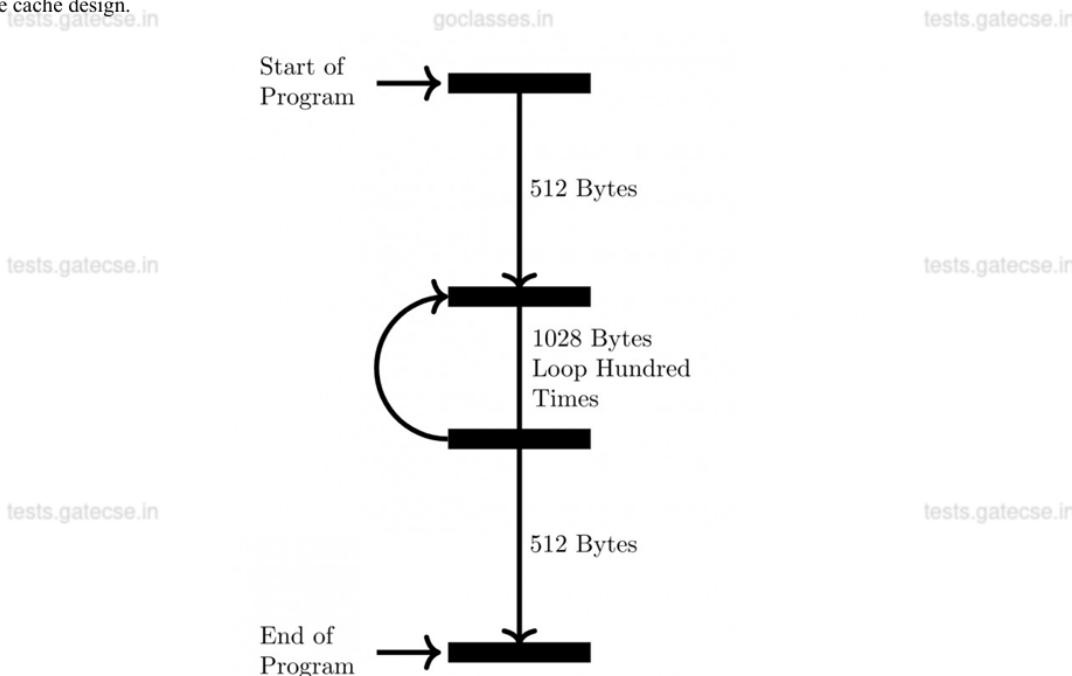
tests.gatecse.in

goclasses.in

tests



A certain computer system was designed with cache memory of size 1 Kbytes and main memory of 256 Kbytes. The cache implementation was fully associative cache with 4 bytes per block. The CPU memory data path was 16 bits and the memory was 2-way interleaved. Each memory read request presents two 16-bit words. A program with the model shown below was run to evaluate the cache design.



Answer the following questions:

- What is the hit ratio?
- Suggest a change in the program size of model to improve the hit ratio significantly.

$$\text{in cache } \# \text{blocks} = \# \text{lines}$$

<https://gateoverflow.in/85403>

A block-set associative cache memory consists of 128 blocks divided into four block sets. The main memory consists of 16,384 blocks and each block contains 256 eight bit words.

$$= 256 \text{ bytes}$$

- How many bits are required for addressing the main memory? 22
- How many bits are needed to represent the TAG, SET and WORD fields?

$$9 \quad 5 \quad 8 \quad \text{TAG} \quad \text{SET} \quad \text{WORD}$$

← 22 bits →

main memory  
no. of blocks × no. of line per blocks  
 $16384 \times 256$   
4194304 → Total no. of lines

$$\log_2(4194304)$$

22

↓  
no. of bits required to address main memory

$$\text{no. of words} = 256 \Rightarrow \text{no. of bits required for word} = \log_2(256)$$

$$= 8$$

← 22 bits →

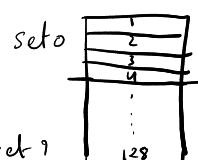
$$\text{TAG} \quad \text{SET} \quad \text{WORD}$$

$$\leftarrow 11 \rightarrow \leftarrow 8 \rightarrow$$

$$\leftarrow 9 \rightarrow \leftarrow 5 \rightarrow$$

four block sets = n line sets = n way

i.e. in each set there are n lines



$$? \times 4 = 128$$

$$? = \frac{128}{4} = 32$$

$$\# \text{sets} = 32 \Rightarrow \# \text{set bits} = \log_2(32) = 5$$

