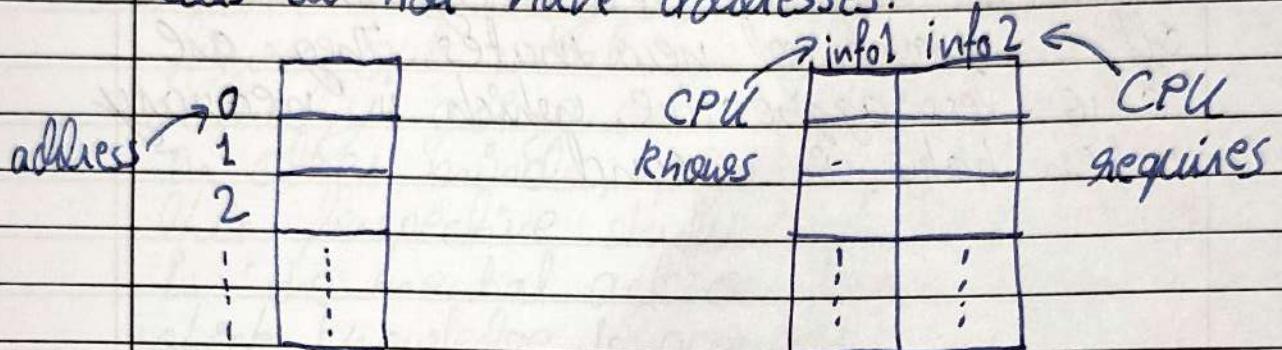


## Lecture 25

### Associative Memory

- Known as Content Addressable Memory
- Cells do not have addresses.



Addressable Memory      Associative Memory

To access the data, CPU will give some info to the memory and by return will get required data.

The required data is searched for in associative memory, unlike Addressable memory where address is directly accessed.

- ⇒ Each cell in associative memory has a matching logic, this allows each cell to be checked parallelly and this memory hence is very fast. (faster than SRAM)

But this also makes associative memory very fast expensive.

Practically, we try to make all info 1 unique.

But if there is a conflict, the one which occurs ~~first sequentially~~ is returned. all matches are sent sequentially (one after another).

### Application of associative memory

- i) can be used for cache implementation.  
(for fully associative mapping)
- ii) for TLB (Translation Lookaside Buffer)  
(TLB is studied in OS)

### Locality of Reference

If CPU requests one address in memory, then that particular address or nearby addresses will be accessed soon.

It was noticed that during operation of a computer, this statement is very probable to be true.

The quality of computer accessing the memory is called locality of reference. This is of two types.

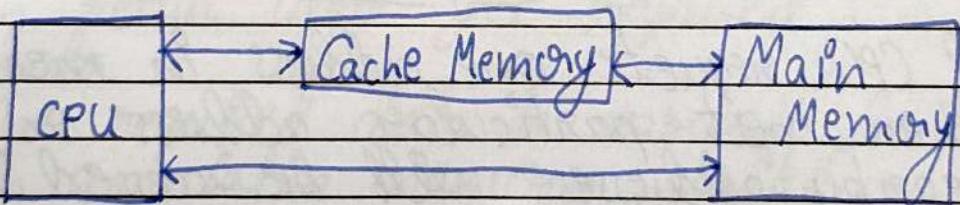
- i) Temporal locality of reference  
(Based on time, same address accessed)
- ii) Spatial locality of ref.  
(Based on space, nearby address accessed)

Example of temporal is loop variable  
Example of spatial is program instructions.

## Cache Memory

Going by the concept of locality of reference, we keep the most accessed parts of main memory in a faster memory called cache memory.

The parts of memory which are most requested are called current demanded locality (or blocks).



## Lecture 26

### Working of Cache Memory

Cache Hit - If CPU demanded content is present in cache

Cache Miss - Content not present in cache

\* Performance is given by Cache Hit Ratio

$$\text{Hit Ratio} = \frac{\text{no of cache hits}}{\text{Total references}} = \frac{\text{no of hits}}{\text{no of hits} + \text{no of misses}}$$

$$\text{Miss Ratio} = 1 - \text{Hit Ratio}$$

⇒ If there is a cache miss, then when we are fetching that data from main memory, we also fetch a whole block of it into cache memory. (block which contains demanded byte is copied).

### Average Memory Access Time

$$\text{Average Memory Access Time (T}_{avg}\text{)} = (\text{Hit ratio} \times \text{Hit time}) + (\text{Miss ratio} \times \text{Miss time})$$

$$T_{avg} = \text{Hit ratio} \times \text{Hit time} + \text{Miss ratio} \times \text{Miss time}$$

Hit Time = Time ~~wast~~ taken for hit

Miss time = Time taken on miss

## Types of Cache Accesses

### 1. Simultaneous Access $\Rightarrow$

Request for cache & main memory are generated simultaneously.  
So CPU accesses them simultaneously and if cache has data then it can access it and request to main memory is discarded.

$$T_{avg} = \text{Hit ratio} \times t_{cm} + \text{Miss ratio} \times t_{mm}$$

$t_{cm}$  = Time to access cache memory

$t_{mm}$  = Time to access main memory

### 2. Hierarchical Access $\Rightarrow$

The faster memory (cache memory) is accessed first and only on miss main memory is accessed.

$$T_{avg} = \text{Hit ratio} \times t_{cm} + \text{Miss ratio} \times (t_{cm} + t_{mm})$$

$t_{cm}$  = Time to access cache memory

$t_{mm}$  = Time to access main memory

~~Tip for which formula to use when,~~  
~~If  $t_{mm}$  /  $t_{cm}$  not given use general formula.~~

~~By default assume hierarchical simultaneous.~~

~~If level/hierarchy in ques  $\Rightarrow$  hierarchical access~~

~~Else  $\Rightarrow$  simultaneous access~~

If ques says to ignore initial checking time  $\Rightarrow$  Simultaneous

## Tug When locality of reference included

### Lecture 27

#### Tug When locality of reference included

⇒ Simultaneous access ⇒

$$T_{avg} = \text{Hit ratio} \times t_{cm} + \text{Miss ratio} \times t_{block}$$

where,  $t_{block} = \text{block size} \times t_{mem}$

This is because when accessed simultaneously, on a miss,  $t_{block} \gg t_{mem}$ , hence the formulae

⇒ Hierarchical access ⇒

$$T_{avg} = \text{Hit ratio} \times t_{cm} + \text{Miss ratio} \times (t_{cm} + t_{block})$$

Similar logic to previous applies.

Note how in both formulae, only  $t_{mem}$  is replaced with  $t_{block}$

#### Tug when block transfer time is included

⇒ Simultaneous access ⇒

$$T_{avg} = \text{Hit ratio} \times t_{cm} + \text{Miss ratio} \times (t_{block} + t_{cm})$$

⇒ Hierarchical access ⇒

$$T_{avg} = \text{Hit ratio} \times t_{cm} + \text{Miss ratio} \times (t_{cm} + t_{block} + t_{cm})$$

Note: if ques mentions miss penalty, use block transfer time included formulae.

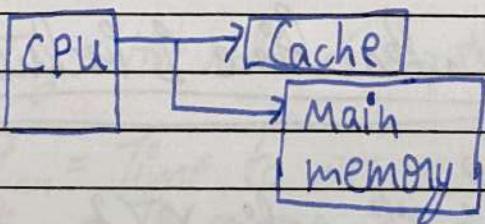
## Cache Write

As cache only keeps a copy of the block in main memory.  
So when CPU writes same data in cache, the block has to be updated in main memory then. This is called write propagation/cache write.

### Write propagation

- └ Write through
- └ Write back

Write through  $\Rightarrow$  When CPU updates Cache content, Main memory also gets updated at the same time



$\Rightarrow$  Regardless of hit or miss, main memory is always accessed in this case on every write by CPU.

## Tavg in write through

$$T_{avg\ read} = \text{Hit ratio} \times t_{on} + \text{Miss ratio} \times t_{mn} \quad [\text{Simultaneous}]$$

$$T_{avg} = \text{Hit ratio} \times t_{on} + \text{Miss ratio} \times (t_{on} + t_{mn}) \quad [\text{Hierarchical}]$$

$$T_{avg\ write} = \max(t_{on}, t_{mn}) = t_{mn} \text{ mostly}$$

$$T_{avg} = (\text{fraction of read} \times T_{avg\ read}) + (\text{fraction of write} \times T_{avg\ write})$$

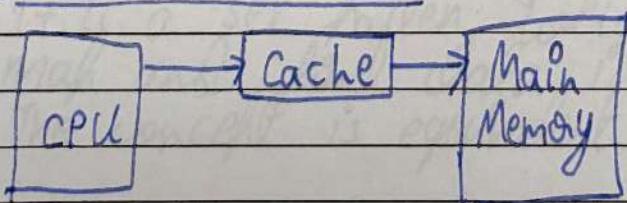
## Effective hit ratio

% of time CPU's request is fulfilled only by cache memory. To find effective hit ratio, use  $T_{avg}$  formulae with both  $T_{writethrough}$  and  $T_{read}$  average.

so basically

$$\text{Effective H.R.} \times \text{Hit cost} + \text{Effective M.R.} \times \text{Miss cost} = \\ \text{fraction of read} \times T_{avg\ read} + \text{fraction of write} \times T_{avg\ write}$$

## Write Back Cache



When block in Cache is being replaced/replaced, only then the old block is copied back to main memory.

But in this type of cache, if I/O device access the main memory, it won't get the correct data which is in cache, this is called memory coherence problem.

### Tag In Write Back

⇒ Simultaneous :

$$t_{avg\ read} / t_{avg\ write} = H.R \times t_{cn} + M.R \times \left( t_{block} + \frac{t_{write\ back}}{2} \right)$$

⇒ Hierarchical :

$$t_{avg\ read} / t_{avg\ write} = \text{Hit ratio} \times (t_{cn}) + \text{Miss ratio} \times (t_{cn} + t_{block} + t_{write\ back})$$

$$t_{write\ back} = x \times t_{block}$$

$x = \text{fraction of dirty block}$

## Lecture 28

### Write Allocate vs No Write Allocate

Write Allocate  $\Rightarrow$  The block is loaded on a write miss (both read and write present) followed by write-hit action.

No Write Allocate  $\Rightarrow$  The block is modified in main memory and not loaded into the cache.

Write Through  
 ↓ uses  
 No write allocate

Write back  
 ↓ uses  
 Write allocate.

## Cache Mapping

Note : CPU only generates address for main memory.

So when CPU generates address we use cache mapping to search in cache using address generated for main memory.

If it is a set partition, which is used to map and find data in cache.  
 The concept is equivalent to hash functions

## Cache Mapping

- Direct Mapping
- Set Associative mapping
- Fully Associative mapping

### Direct Mapping

Note : Mapping is always done on blocks.

Suppose the following,  
Blocks in cache  $\Rightarrow$  10 (0 - 9)

Blocks in main memory  $\Rightarrow$  100 (0 - 99)

In direct mapping, example of mapping will be for example,

$$\text{Cache block location} = \text{Memory block location \% 10}$$

Or generally,

$$\text{Cache block location} = \text{Memory block location \% no. of block in cache}$$

So a certain block in memory is always mapped to same block in cache.

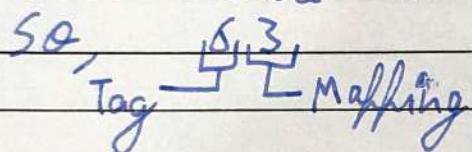
So direct mapping causes conflictions among a group of blocks in main memory.

But the cache itself does not know which block in main memory holds data of.

So each block in cache will hold identification data to store the data about which block from main memory is stored there. This is called a tag.

so for example, if storing 63 in cache,

$63 \% 10 = 3$  is the cache mapping and 6 is stored as a tag

so,  Tag ————— Mapping

This is a common pattern in mapping. The same concept works in binary numbering system.

## Lecture 29

As stated previously, the same idea of tag and cache address from main memory address works in actual systems in binary.

So assume, there are 8 blocks in main memory and 4 blocks in cache memory.

So in the process of storing block no 5 in main to cache,

$$\begin{aligned}\text{Cache block location} &= \text{Memory block location \% No of blocks in cache} \\ &= 5 \% 4 \\ &= 1 \\ 5 &= (101)_2 \quad 1 = (01)_2\end{aligned}$$

Therefore,  $(5)_0 = \underline{\underline{1,01}}$

Tag                  Location in cache memory

This is how the process works in binary

→ Note: Blocks do not refer to the unique addresses assigned to memory. A block is just a continuous chunk of memory, so a single block can have multiple addresses.

For example, if block size is 2 bytes and memory is byte addressable, each block has 2 memory addresses.

### Main memory address to block number

Suppose a block is 2 bytes in a byte addressable memory.

For block number we take the first  $\log_2(\text{no. of blocks})$  bits from the main memory address for the block number. Here, no. of blocks refer to blocks in main memory.

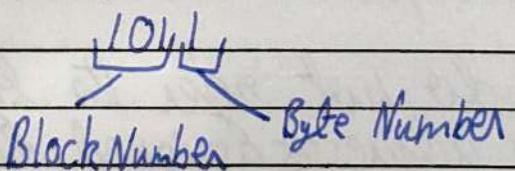
So if there are 8 blocks, each of 2 bytes,

size of memory address = 4 bits

size of block no address = 3 bits

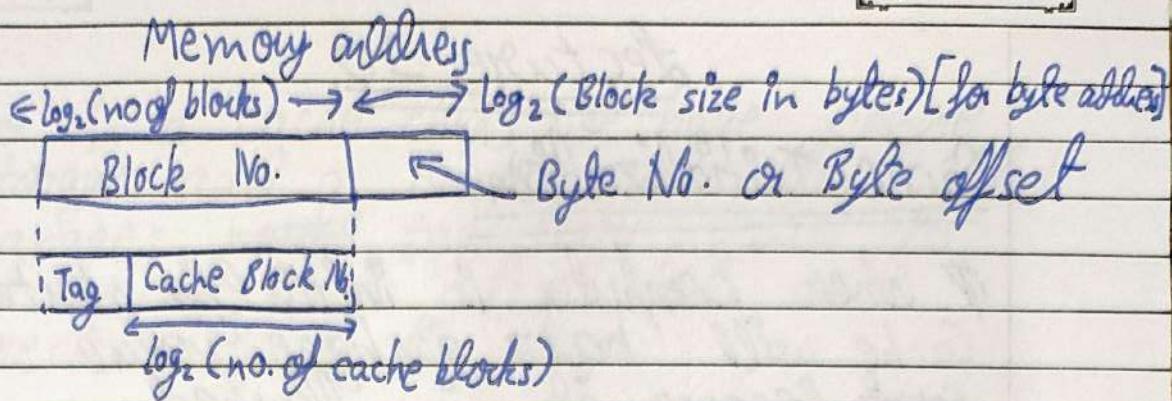
So first 3 bits of memory address show the block number.

Suppose the memory address is 1011,  
Then,



The remaining bits are referred to as byte number.

The block no. is further divided to tag and cache ~~or address~~ block No.



- ⇒ Note: Cache block and cache line refer to same thing. Also called Index field when referring to no. of bits in address form.
- ⇒ Note: Tags are metadata. So total size of all tags is also referred to as metadata of cache memory. Total tag data is also called tag directory.

## Lecture 30

### Cache Initialization

If when computer is initially started, Cache will have garbage value and because it is flushed.

This may cause a hit on CPU request to one of the garbage value.

To resolve this, we keep a valid bit for each block of cache and initialize it to zero, now before checking in cache, CPU will check the valid bit.

If the valid bits are considered a part of Tag directory.

### Write Back Cache Performance Improvement

If CPU only wants to perform read on a cache block, then that block can be ~~replaced~~ directly w/o write back.

This can be done using Modified bit or dirty bit.

This bit is also a part of tag directory (other bits can also be mentioned in ques)

So in a cache chip, we have cache storage and a cache controller which manages most tag, dirty, valid bits of cache and also sends and receives signals to CPU.

So when talking about cache storage, we don't talk about these metadata bits.

Metadata in cache  $\Rightarrow$  Tag, modified bit, Valid bit

In write back cache, simultaneous :-

$$t_{\text{tag read or write}} = H \times t_{\text{on}} + M \times (t_{\text{block}} + t_{\text{write back}})$$

Tag hierarchical :-

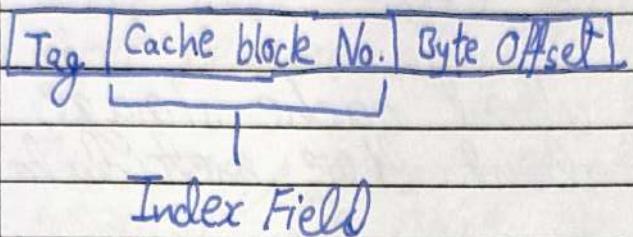
$$t_{\text{tag read or write}} = H \times t_{\text{on}} + M(t_{\text{on}} + t_{\text{block}} + t_{\text{block}} \times t_{\text{write back}})$$

$$t_{\text{write back}} = \text{fraction of dirty block} \times t_{\text{block}}$$

## Lecture - 31

### Set Associative Mapping

Memory address

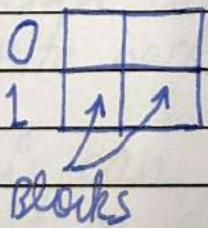


### 2-way set associative

2 blocks are present per index in cache.

block no.  
So, Cache ~~index~~ =  $\frac{\text{Main memory}}{\text{block Number}} \%$   $(\frac{\text{No. of sets}}{\text{in cache}})$

The main reason to do this is so that if CPU is frequently accessing two pages then and they access are mapped to same block in cache the blocks are not frequently replaced.



So two tags are present for each set in cache.

We refer to this collection of block to a single address number as a set.

So main memory address has a change in terminology

[Tag | Set Offset | Byte Offset]

The process to get set offset remains same.

Number of blocks per address in cache = Associativity of cache

Number of sets in cache =  $\frac{\text{Number of blocks in cache}}{\text{Associativity}}$

Note: If block size not given in question remember that the size in bits in direct mapped cache of cache block no. and byte offset is  $\log_2(\text{cache size})$  or  $\log_2(\frac{\text{cache size}}{\text{associativity}})$

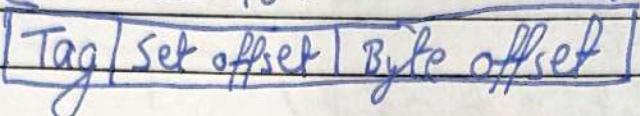
[Tag | Cache Block No. | Byte Offset]

$\log_2(\text{cache size})$  bits or  
 $\log_2(\frac{\text{cache size}}{\text{associativity}})$

Can also assume block size to be  $2^x$ . Assuming it to  $2^x$  will also help solve questions of associative mapping

- Q) For example, width of physical address on a machine is 40 bits. The width of tag field in a 512 KB 8-way associative cache is \_\_\_\_\_ bits?

Sol)



Let block size be  $2^x$

Byte offset =  $\log_2(\text{Block size}) = x$

~~Set offset~~

~~Number of sets in cache~~

Number of blocks in cache = Cache size  
Block size

$$= 2^9 \times 2^{10} \text{ Bytes}$$

$$= 2^{19-x} \text{ Bytes}$$

Number of sets in cache = Number of blocks in cache  
Associativity

$$= 2^{19-x}$$

$$= 2^3$$

$$= 2^{16-x}$$

Set offset =  $\log_2(\text{Number of sets in cache})$   
 $= 16-x$  bits

$$\text{Byte offset + Set offset} = 16 - x + x \\ = 16$$

$$\text{Bits for tag} = 40 - 16$$

## Fully Associative Mapping

In fully associative mapping,

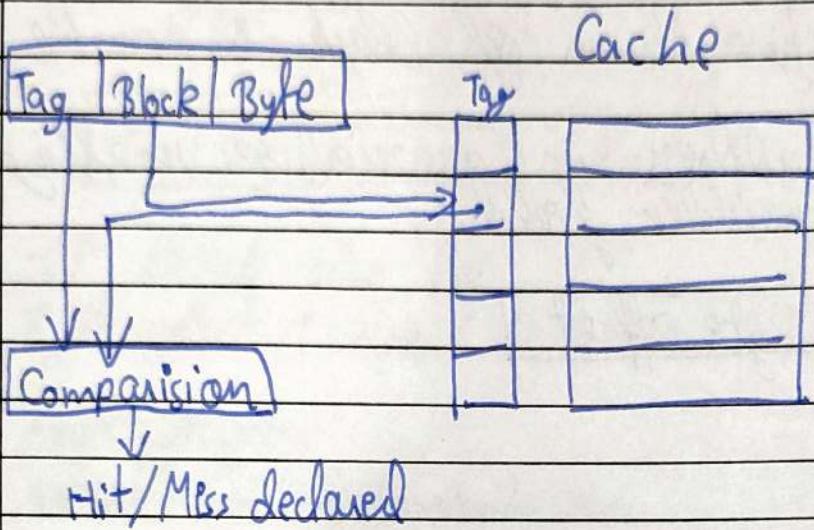
associativity = Number of blocks in cache.

Memory address in associative mapping has no middle field.

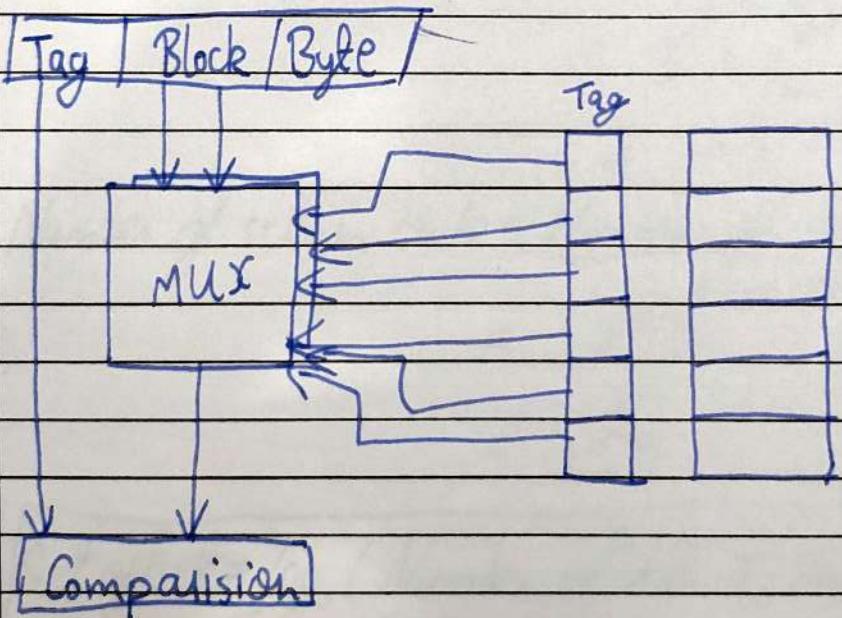
[ Tag | Byte offset ]

## Lecture - 32

### Checking Hit/Miss in Direct Mapping



### Hardware Implementation in Direct Mapping



No. of MUX = Number of bits in tag  
 Size of MUX = Number of blocks in cache  $\times 1$   
 Size of Comparator = Number of bits in tag.

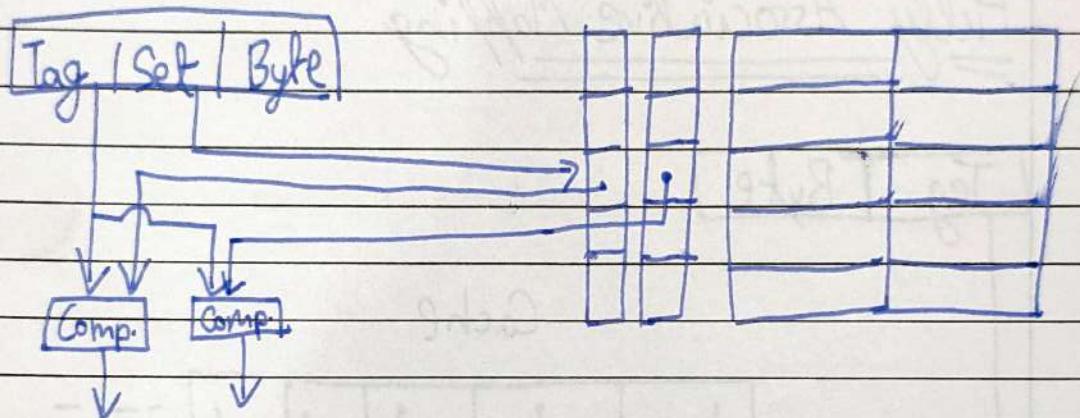
## Hit Latency in Direct Mapping

Latency caused to check for a hit or miss in cache.

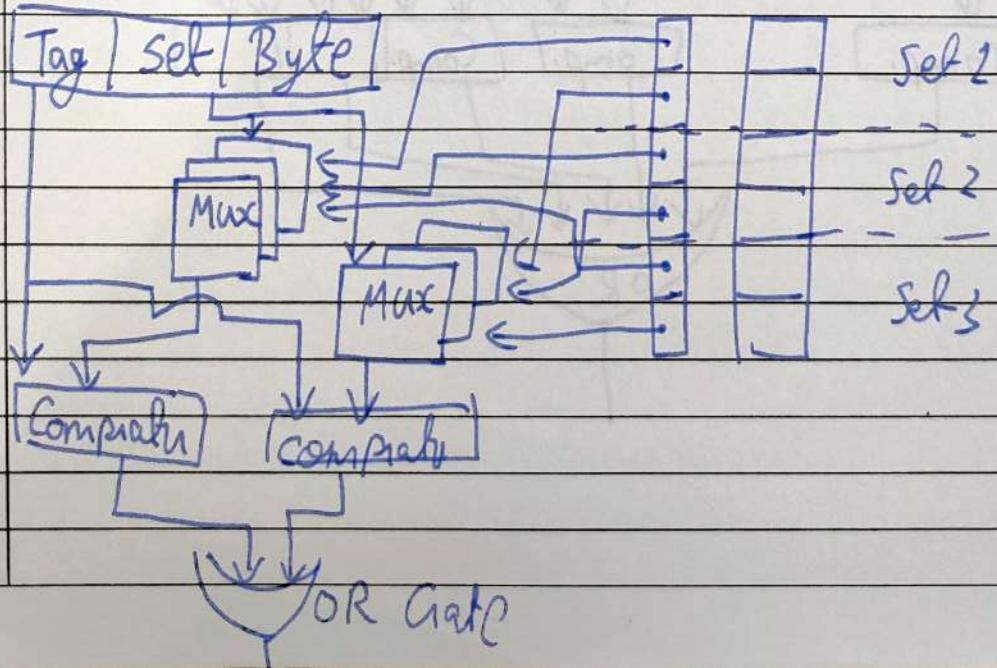
Hit latency = Mux delay + Combinatorial delay for tag selection + Delay.

## Checking Hit/Miss in Set Associative Mapping

Cache



## Hardware Implementation



For k-way associative

Size of MUX = No. of sets in cache  $\times 1$

No. of MUX for tag selection =  $k \times$  No. of tag bits

No. of comparators =  $k$

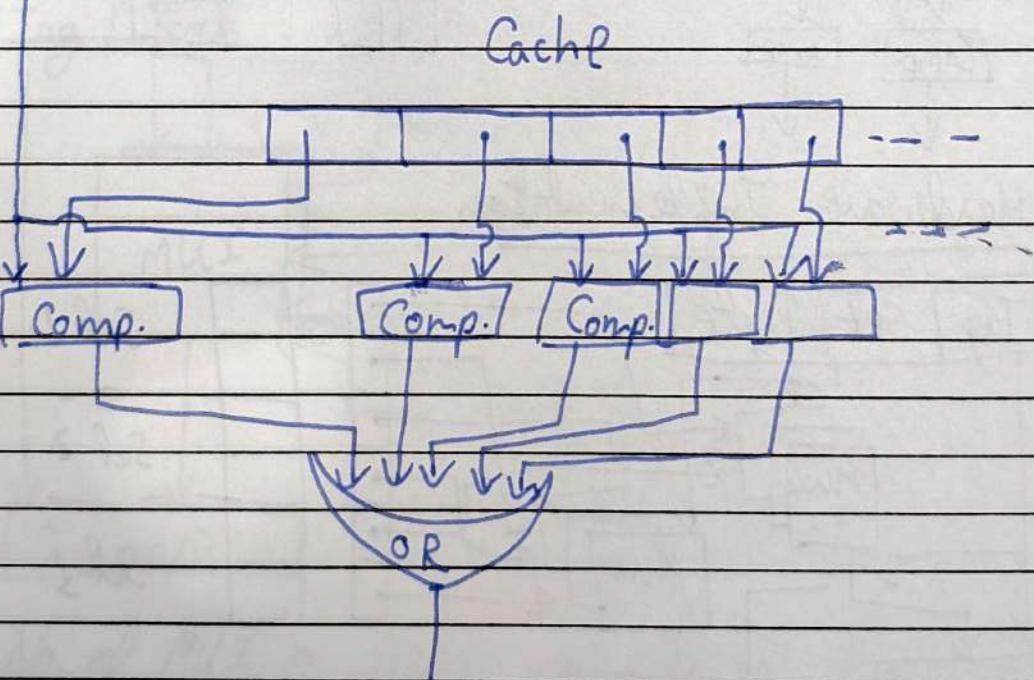
Size of comparator = No. of tag bits.

Hit/miss & latency in k-way set associative

Latency = MUX delay for tag selection + Comparator Delay + OR Gate Delay.

## Fully Associative Mapping

Tag | Byte



No. of comparators = No. of blocks in cache  
Size of comparators = tag bits.

Note: Sometimes question will mention a  $2 \times 1$  Mux rather than an OR Gate. Here  $2 \times 1$  Mux is used to implement the OR Gate.

## Lecture - 33

### Block Replacement : Direct Mapping

In direct mapping, we don't have a special block replacement policy and replace block if a new block is to be brought in cache.

### Block Replacement : Set Associative Mapping

Since multiple blocks are in the same set, we need a policy to know which block is replaced.

#### Replacement policies

- FIFO (First in First out)
- Optimal  $\Rightarrow$  Keep the block which is used more in the future.  
(Not practical)
- LRU (Least Recently Used)  $\Rightarrow$  Replace the block which has not been referred for longest time.

These three policies are also used in fully associative cache.

## Lecture - 34

### Cache Misses

- Cold or Compulsory Miss
- Capacity Miss
- Conflict Miss

↓ check for type of miss in this order

#### Cold or Compulsory Miss

First time access of a block will always cause a miss

Solution: Increase block size to decrease compulsory Miss

#### Capacity Miss

If cache is full and hence miss occurs.

Solution: Increase the cache size.

#### Conflict Miss

If cache set is full and hence miss occurs due to tag mismatch. (Not a cold or capacity miss)

Solution: Increase the associativity.

## Cache Miss Penalty

Time required to bring a missed block from main memory to cache.

Three main info are needed to calculate miss penalty.

- Cycles required to send address to memory
- Cycles required to access 1 main memory cell
- Cycles required to transfer 1 cell data to cache.

Example-

Cache block size	Memory cell size	Miss Penalty
4 bytes	1 byte	$1 + (4 \times 10) + (2 \times 1)$
4 bytes	2 bytes	$1 + (2 \times 10) + (2 \times 1)$
4 bytes	4 bytes	$1 + (1 \times 10) + (1 \times 1)$

The given info for above is

Cycles to send memory address = 1 cycle

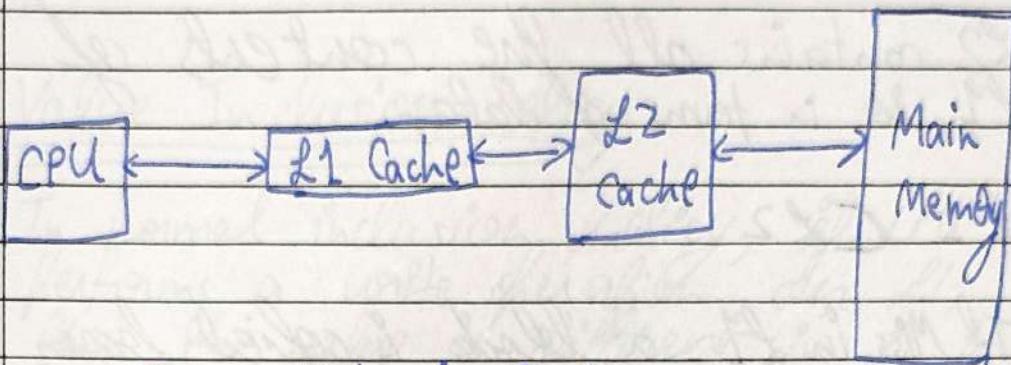
Cycles to access 1 main memory cell = 10 cycles

Cycles to transfer 1 cell data to cache = 1 cycles

## Lecture 35

### Goals of Using cache memory

- Minimize Access Time  $\Rightarrow$  Smaller Cache  $\Rightarrow$  Multi-level Cache.
- Maximize Hit Rate  $\Rightarrow$  Larger Cache  $\Rightarrow$  Cache.
- Minimize Miss Penalty



### Time in Multi-level Cache

#### 1. Simultaneous Access

$$T_{avg} = H_1 \times t_1 + M_1 [H_2 \times t_2 + M_2 \times t_{mm}]$$

<sup>24</sup> Subscript 1  $\Rightarrow$  For Level 1

Subscript 2  $\Rightarrow$  For Level 2

#### 2. Hierarchical Access

$$T_{avg} = H_1 \times t_1 + M_1 [H_2 \times (t_1 + t_2) + M_2 (t_1 + t_2 + t_{mm})]$$

## Cache Inclusion Policy

- Inclusion
- Exclusion

### Inclusion Policy

L<sub>2</sub> contains all the contents of L<sub>1</sub>. So in terms of data,

L<sub>1</sub> ⊂ L<sub>2</sub>

- \* If Miss in L<sub>1</sub>, a block is copied from L<sub>2</sub> if it is hit in L<sub>2</sub>.  
If block is evicted from L<sub>1</sub>, then no changes in L<sub>2</sub>.
- \* If a block is evicted from L<sub>2</sub> which is in L<sub>1</sub>, then send invalidation to L<sub>1</sub>.
- \* If miss on both L<sub>1</sub> & L<sub>2</sub>, copy block to both L<sub>1</sub> & L<sub>2</sub>.

### Exclusion Policy

It is not necessary that content of L<sub>1</sub> cache is also present in L<sub>2</sub>.

- \* If miss in L<sub>1</sub> and hit in L<sub>2</sub>, then move block from L<sub>2</sub> to L<sub>1</sub>, and evicted block from L<sub>1</sub> is

moved to L2.

- \* Misses in both L1 & L2, copy block from main mem to L1 only and move evicted block to L2.

L2 hence is called victim cache as it is only taking data from L1.

### Value Inclusion Policy

In normal inclusion policy, if CPU performs a write operation on L1, then the value is different in L1 & L2.

To have the same value of blocks in L1 & L2, we use value inclusion policy.

Lec 36 and 37 are practice questions on cache

Mind  
Date: / /20  
Page:

## Lecture - 38

### Magnetic Disk

In a magnetic disk, there are multiple dis

In a hard drive, there are multiple magnetic disks arranged in a stack. Each disk is called a platter.

Unlike a optical disk, in a magnetic disk, both surfaces can be used to store data.

so if a drive has 3 platters, we say we have 6 surfaces

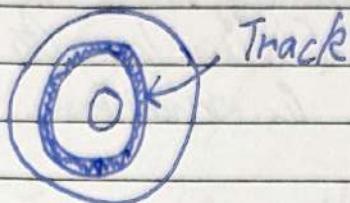
The column that spins all the platters is called spindle

On each surface we have a read/write head, all heads are attached to an arm, all arms are attached to arm assembly.

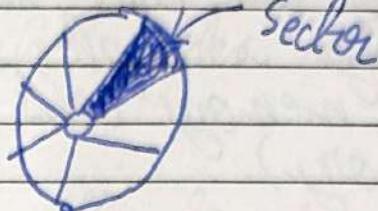
The arm assembly is a perpendicular column, arms rotate with assembly to cover the whole surface of platters.

If the arm was fixed, only a single line of circle would be read from the platter.

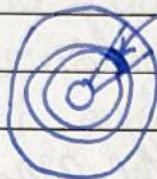
In a magnetic disk,



Track



Sector



sector of a track  
(referred simply as  
sector when  
talking about  
disk)



cluster

## Disk Capacity

- Constant sector capacity (Variable storage density)
- Variable sector capacity (Constant storage density)

The storage density is the amount of data stored per unit surface area.  
The above difference or inverse effect is caused because sectors near inner radius will have smaller areas compared to sectors near outer radius.

In questions consider constant sector capacity as this is more widely used.

Sector in disk is the smallest unit of disk which can be read or written at once.

i.e for one disk read/write op. is accessing a single sector (similar to how memory cells work in main memory).

## Disk Access Time

Time taken to access 1 sector.

Seek time  $\Rightarrow$  Time to move the arm to correct position.

Rotational Latency  $\Rightarrow$  Time taken to rotate disk to bring sector under head to the head.

Transfer Time  $\Rightarrow$  Time required for read/write.

Disk Access Time  $\Rightarrow$  seek time + rotational latency + transfer time + additional latency given in question.

Sometimes, we need to calculate seek time using disk scheduling, this is a concept of OS.

\* The head is not on top of the head sector, it is in between two sectors, so at the start on one sector and at the

end of another.

- \* The rotation of platter happens in a single direction.
- \* The head should be at the start of the sector which is to be under read/write op.

In questions, either we will be given data about head position and rotation time to calculate rotational latency. Or we can use average rotational latency if that is given.

If average rotational latency is also not given, then,

$$\text{rotational latency} = \frac{1}{2} \text{ disk rotation time}$$

For transfer time:

In 1 rotation of disk = 1 track of data can be transferred.

Then,

$$1 \text{ sector transfer time} = \frac{1}{\text{No. of sectors per track}} \text{ rotation of disk time}$$

Constant sector capacity



Variable storage density



Constant angular velocity

Variable sector capacity



Constant storage density



Constant linear velocity

## Lecture - 39

### Multiple Sectors Access Time

#### Sequential Multiple Access

The sectors are on the same track in a continuous manner.

So if we transfer  $n$  sectors :

$$\text{n sector access time} = \text{seek time} + \text{rotational latency} + n \neq 1 \text{ sector access time.}$$

#### Random multiple Access

Since all the sectors are scattered

$$\text{n sector access time} = n * [\text{seek time} + \text{rotational latency} + 1 \text{ sector transfer time}]$$

#### Cylinder in Disk

We want to always reduce seek time.  
(i.e. the time to position read/write head).

So if suppose we wanted store a large file, we will store it on a single track. When this track is completely filled, rather than moving

The head, we store next part in the same track on below surface. When that is also full, we use the same track on the next disk. This reduces seek time.

Hence we store data on disk in form of cylinders.

Number of cylinders = No of tracks.