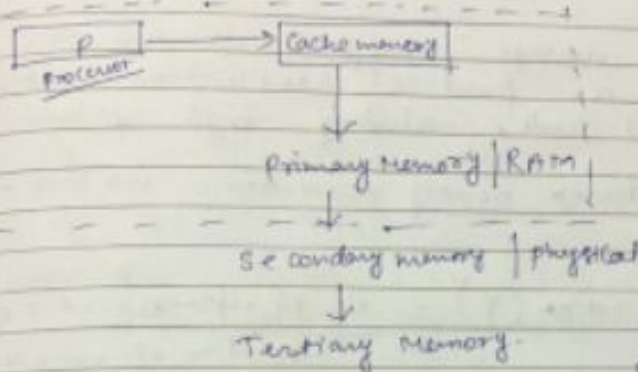


## Thread Level Parallelism Memory Hierarchy design.

6/3/02

### Memory Hierarchy design



Cache hit  
Cache miss

- # Memory mapping is done by virtual memory.
- # RISC and CISC

CU  
└── Hardware    Microprogram control.  
                  (combination of software + hardware)

Prob. locality of reference (b/w P and Cache)

## # Evaluation of cache performance and optimization of cache.

The average memory access time is calculated as a part of instruction execution at all levels, including high & low.

$$\text{average memory access time} = \text{Hit time} + \text{miss rate} \times \text{miss penalty}$$

$$\text{Hit ratio (Y)} = \frac{\text{No. of references found in cache}}{\text{total no. of memory references}}$$

miss rate is the fraction of memory references not found in cache.

$$\text{miss rate} = \frac{\text{No. of misses}}{\text{references}}$$

miss penalty is the additional time required because of a 'miss'.

CPU time in the presence of misses, the CPU time can be considered with the following <sup>operational</sup> entities: —

$$\text{CPU time} = (\text{CPU execution clock cycles} + \text{memory stall clock cycles}) \times \text{clock cycle time}$$



There are 16 Cache Optimizations into 4 categories:

1) Reducing the miss penalty

- a) Multilevel Caches
- b) Critical word first
- c) Read miss before write miss
- d) Merging Right buffers.
- e) ~~M~~ No. of victim Caches

2) Reducing the miss rate

- a) ~~arranging~~ Larger block sizes
- b) Larger Cache size
- c) Higher associativity
- d) Pseudo activity (Pseudo associativity and compiler optimization)

3) Reducing the miss penalty and miss rate via Parallelism: —

- a) Non-blocking Caches
- b) Hardware Pre-fetching
- c) Compiler Pre-fetching

4) Reducing the time to hit in the cache

- a) Arranging Small and simple Caches.
- b) Avoiding Address translation using VM
- c) Pipeline Cache access

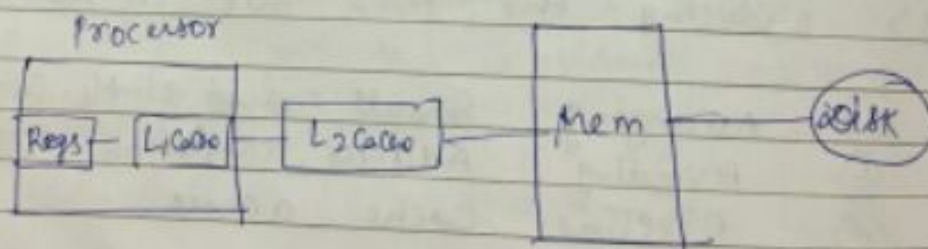
clock  
cycle time

locality of reference

Many of times, a computer executes instructions from consecutive memory location serially until the information happens to be a branch and branch takes a flow in a different locality in the memory.

# With the advancement in the execution of a program, <sup>slowly</sup> changes with time. This is referred as Spatial locality.

# It is possible that same location may be accessed again in short intervals of time or a nearby location may be accessed in another reference. This is referred as Principle of locality.

Various techniques to reduce miss penalty.1) Arranging Multilevel Caches.

In order to rate the miss penalty, the avg. mem. access time for 2 level caches is calculated as

$$\text{Avg mem Acc time} = \text{Hit time of L1} + \text{Miss rate} \times \text{Miss Penalty L1}$$



Fetch the words in normal order but as the request word of block arrives, send it to the CPU and let this CPU continue execution.

c) Giving ~~the~~ priority to read misses over write

with a right through cache, the most imp. improvement is right buffer of the proper size (i.e. 2).

The simplest way out of ~~the~~ writing is for the read miss ~~until~~ to wait until the right buffer is empty.

d) Merging Right Buffer → This technique involves right buffers with an increase in time efficiency. write through caches rely on write buffers. ~~and~~ If the buffer contains other modified blocks, the addresses can be checked to see if <sup>the</sup> address of <sup>the</sup> new data matches the address of valid write buffer entry. If so new data are combined with entry. (Right Merging).

100
200
300
400

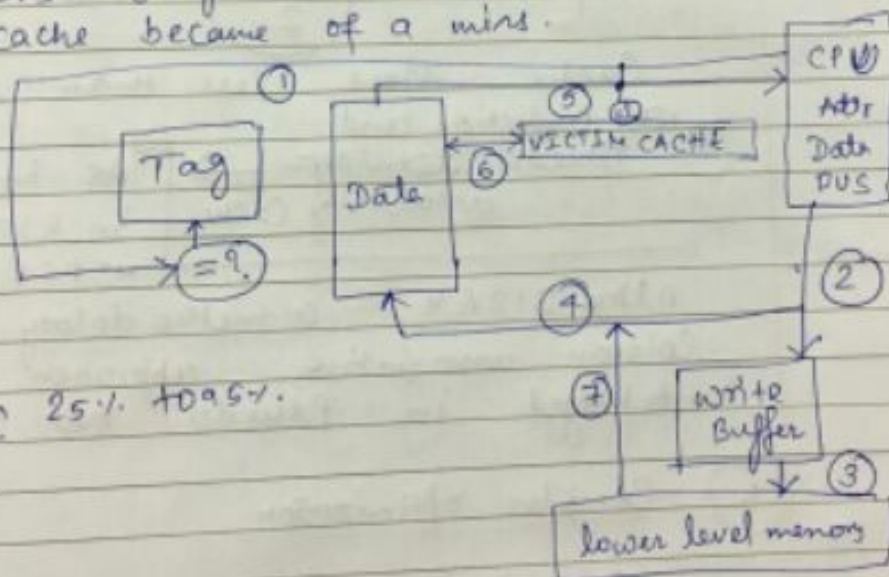
Valid		valid bit	
1	M[100]		0
1	M[200]		0
1	M[300]		0
1	M[400]		0
0	M[500]		0

# There is <sup>no</sup> multiple ~~or~~ writes for same memory location

c) No. of victim caches.

The cache which is suffering from miss ~~operation~~ operation.

① one approach to lower miss penalty is to remember what was discarded in case it is needed again. Since the discarded data has already been fetched, it can be used again at small cost. This is called recycling. ~~The~~ with this activity related cache is referred as a victim cache ~~that~~ which contains only blocks that are discarded from a cache because of a miss.



Maximum 25% to 95%.



## ② Reducing the miss rate

① arranging large blocks with required size

↓  
It is a concept of temporal and spatial locality.

② Arranging larger caches

③ Higher associativity → It is a concept of direct mapped cache of size  $N$  has about the same miss rate as 2 Way set associativity cache of size  $\frac{N}{2}$ . This helped for

Cache sizes less than 128 KB.

④ ~~Way Prediction and Pseudo associativity~~ → Pseudo associativity Cache.

As per the relation in RISC architecture MIPS R4300 and

Alpha 21264, a methodology named Column associative approach has been deployed in Pseudo set.

## ⑤ Compiler optimization

① Loop interchange

② Blocking

# Array can't be accommodated in cache memory

Loop interchange  $\rightarrow$  Some programs have nested loops that access data in memory in non-sequential order/random order. Simply exchanging the nesting of the loops can make the code access flexible. Assuming that the arrays don't fit into the cache. This technique reduces misses by improving spatial locality.  
 Note: Reordering maximizes uses of data in a cache block before it is discarded.

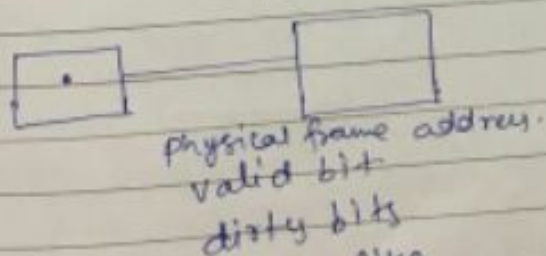
Blocking  $\rightarrow$  This optimization tries to reduce misses through improved temporal locality where the arrays are stored by row major order and column major order to transform the block interchange addresses.

Ranges  $\rightarrow$  find size  
 Segment  $\rightarrow$  virtual size.

Smallest - 1 Byte  
 Largest  $\rightarrow 2^{32}$  byte.

$2^{16}$  to  $2^{32}$  bytes

Virtual memory  $\rightarrow$  TLB  
 Translation looked ahead Buffer



- A TLB entry is like a cash entry where the Tag holds portion of the virtual address and data portion holds a physical frame number.



Valid bit (0,1) and dirty bit.

To change the physical page frame number OS must make sure that old entry is not in the TLB, otherwise system will not work properly.

The dirty ~~bit~~ bit means corresponding page is dirty, not that address translation in TLB is dirty, nor a particular block in data cache is dirty. The OS resets these bits by changing the value in page table and then ~~enter~~ invalidating the corresponding TLB entry.

Page size Choosing the page size is a difficult operation. So OS forces that a larger page size vs smaller size could be created to favour larger size creation. OS always favours large size creation.

## Cache coherence

Two Processors — Two different value

↓  
Same memory location.

	Processor Activity	Bus Activity	Content of CPU A's Cache	CPU B's Cache	Contents of Mem Loc X
OP <sub>1</sub>					0
OP <sub>2</sub>	CPU A reads X	Cache miss for X	0	—	0
OP <sub>3</sub>	CPU B reads X	"	0	0	0
OP <sub>4</sub>	CPU A writes X	write broadcast X	1	1	1
OP <sub>5</sub>	CPU B reads X	—	1	1	1



32 Symmetric Shared memory architecture.

① SSMA consist of several processors with a single physical memory

② SSMA machines usually support the caching of both shared and private data. Private data is used by single processors while shared data is used by multiple processors.

When a private data item is cached, its location is migrated to the cache, reducing the average access time as well as the memory bandwidth required.

Note: Since no other processors use the data, the program behaviour is identical to that in a uniprocessor.

An enforcing based protocol <sup>Snooping base protocol</sup> SBP maintain coherence for multiple processors, where each cache that has a copy of data from a block of physical memory also has a copy of sharing status of the block. and no centralised state has been kept. The cache are usually on a shared memory bus and all cache controller monitor (snoop) on the bus to determine whether or not they have a copy of a block ~~the is~~ that is ~~required~~ requested on the buses.

Implementation The serialization of access enforced by the bus also ~~enforces~~ <sup>forces</sup> serialization of a writer. Since when two processors compete to write to the same location, one must obtain bus access before the other. The first processor to obtain bus access will ~~cost~~ <sup>cost</sup> other processors copy to be invalidated, causing writes to be strictly serialized.

One implication of this scheme is that a write to a shared data item can't complete until it obtains bus access.

### Directory Based Protocol (DBP)

In this the sharing status of a block of physical memory is kept in just one location called the directory.

The drawbacks of (Symmetric-SMA) ~~transfer~~ <sup>transfer</sup>

- ① complex mechanism for ~~transfer~~ <sup>transfer</sup> software cache coherence are very limited
- ② without cache coherence the multiprocessor loses advantage of being able to fetch and use multiple words in a single cache block
- 3> The mechanism for tolerancy the latency such as pre-fetch are more useful when they can fetch multiple words such as a cache-block.

17 A direct that no director of blo to p bottle along direc location go

A that also is an

To write dire of

with the

(a) Sho on Cack

(b) Unc a cop

(c) Ex ce

1. ②



17 A directory keeps the state of every block that may be cached. Information in the directory includes which caches have the copy of block whether it is dirty and so on. To prevent directory from becoming the bottle neck, directory entries can be distributed along with the memory so that different directory accesses can go to different locations just as different memory request go to different memories.

A directory distributed scheme characteristic that a sharing status of block is always in a single non-location. This property is what allows coherence protocol to avoid broadcast.

To implement <sup>the operations like</sup> handling a <sup>read</sup> miss, handling a write and cleaning the cache block, the directory must be able to track the state of each cache block. In simple directory based protocol, these states could be presented

with the following:-

- (a) Shared operations: It illustrates that whether one or more processors have the blocks cached and the value in memory is <sup>up to</sup> update or not.
- (b) Uncached operations: No processor has the a copy of <sup>cached</sup> block.
- (c) Exclusive: Exactly one processor has a copy of cache block and it has ~~set~~ written the block, so the memory copies is out of date. This processor is called owner of the block.  
owner node / Home node.

~~chapter -> 4, 5~~

1, 3, 4, 5,

ch 2

1> chapter - 4

2> chapter - 5

Page No.:

--	--	--



## Multi-threading.

It

- ① fine grained multithreading: → Switches b/w threads on each ~~ex~~ instruction causing the execution of multiple threads to be interleaved. This interleaving often done in round robin fashion, skipping any threads that are stalled/problematic at that particular time.
- ② Coarse grained multithreading → It is alternative of fine grained multithreading. It switches the threads only on costly stalls, such as level 2 cache misses. This change relieves the need to have thread switching be essentially free and is much less likely to ~~slow~~ slow the processor down. (means instruction from other threads <sup>will</sup> only be issued when a thread encounters a costly stall).
- 3) Simultaneous Multithreading → It works on instruction level parallelism.

04/2/2018

Page No.

Thread level Parallelism (TLP) (multiprocessors)

↓ s/w routine

↓ compiler

↓ Instruction exe.

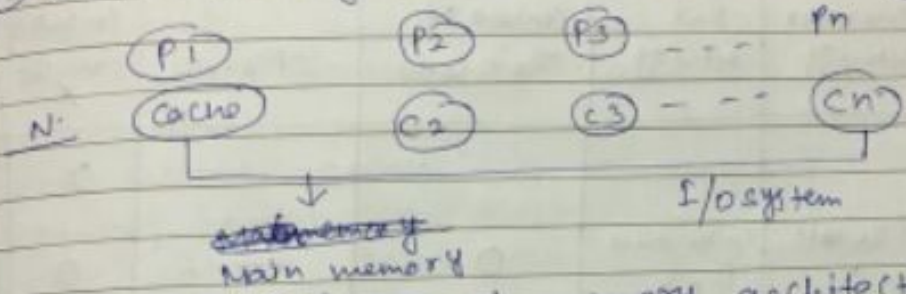
↓ sharing of cache

↓ thread level Parallelism

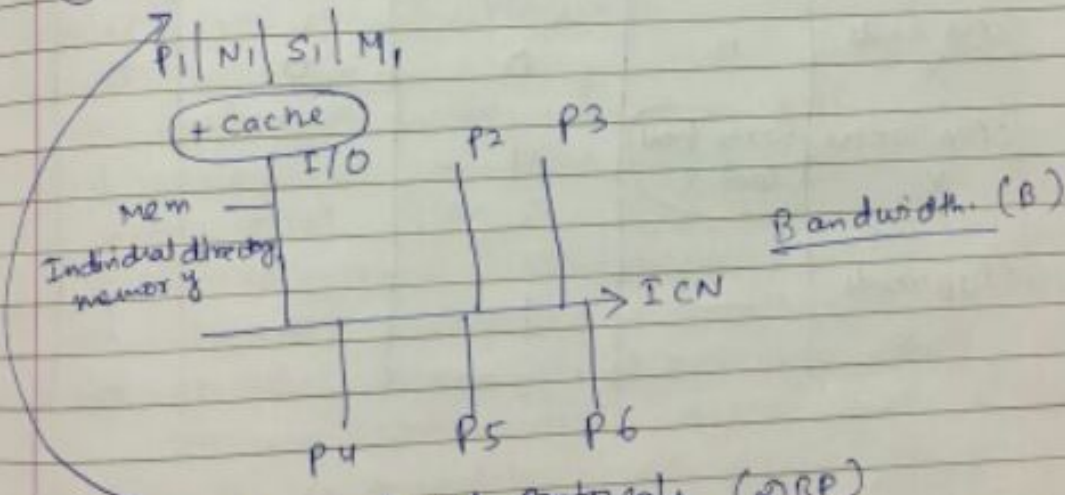
↓ miss/hit

Protocol in SNA

① Shared memory architecture (SMA)



② Distributed shared memory architecture (D-SMA)



Directory based protocols (DBP)

Write Invalidate Protocol (WIP)

Write validate Protocol (WVP)

DBP



## 2) Buffer management

$$\text{Miss Penalty } L1 = \text{Hit time } L2 + \text{Miss rate } L2 \times \text{Miss Penalty } L2$$

~~20. Avg~~

total miss rate  $\rightarrow$   $\frac{\text{The no. of misses in the cache}}{\text{The total no. of memory accesses to that cache.}}$

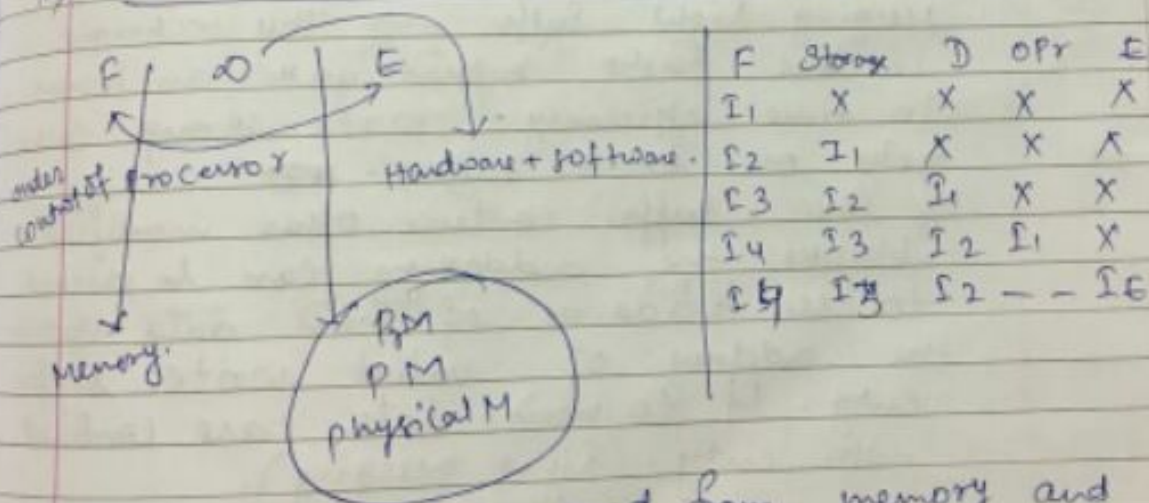
global miss rate  $\rightarrow$   $\frac{\text{The no. of misses in the cache}}{\text{Total no. of memory access generated by the CPU.}}$

Avg memory stall per instruction can be calculated with

$$\text{misses per instruction } L1 \times \text{Hit time } L2 +$$

$$(\text{wrapped word}) \text{ misses per instruction } L2 \times \text{Miss Penalty } L2$$

11) Critical word first and Early restart



Request the miss word from memory and send it to CPU as soon as it arrives; let the CPU continue execution while filling the rest of the words in the block. The accommodated bits of an instruction (Outgoing instruction I<sub>E</sub>) would be referenced w.r.t. the tag of cache associated with its corresponding blocks in memory. Which is called wrapped fetch/Req. word first.