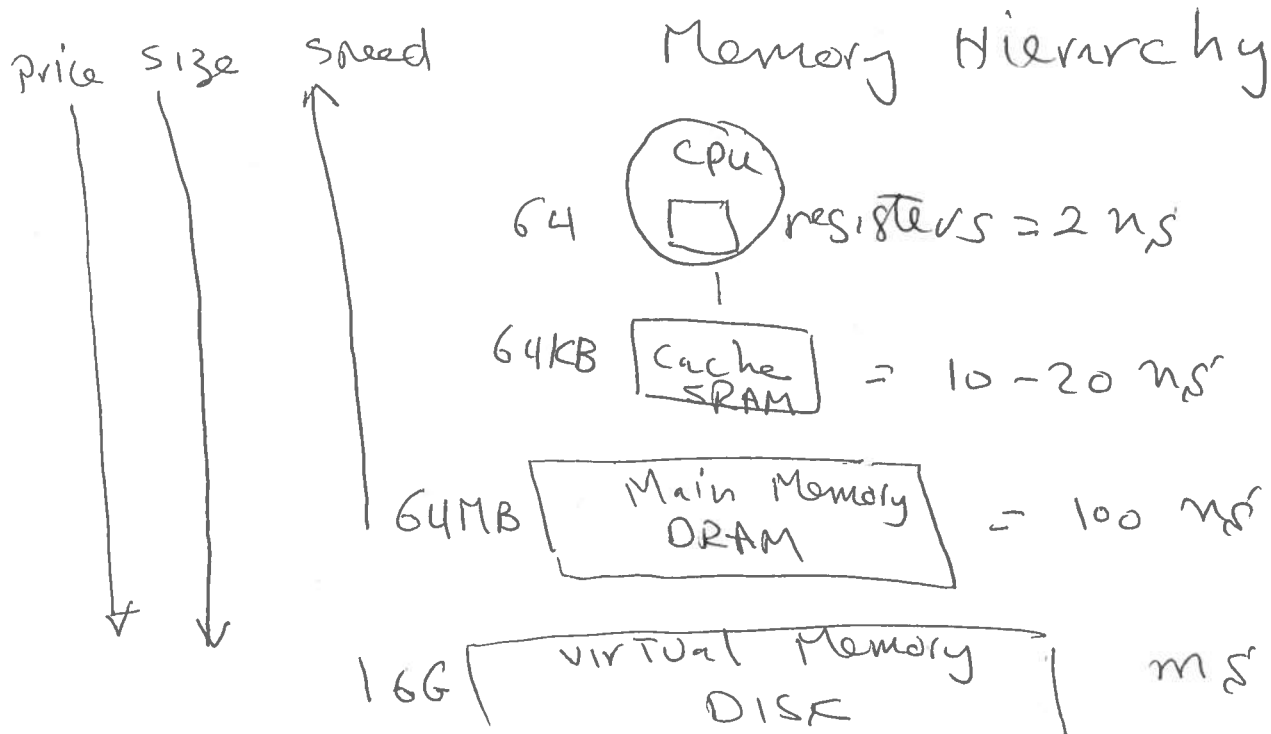


Cache

- 1- Principle of Locality
- 2- Terminology
- 3- Direct Mapped Cache

Cache Concept

- CPU executes 10% of code 90% of the time (locality).
- CPU is very fast (2 ns for 500 MHz) DRAM is slow (100 ns), each memory access takes 50 CPU cycles.
- Small parts of code & Data could be stored in fast SRAM, so CPU will execute at rated speed.
- Cache works because of principle of locality.



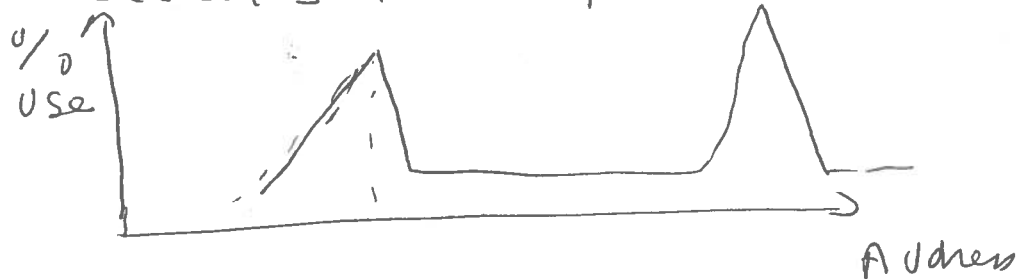
Cache Concept

Principle of Locality

90% of Time executing 10% of Code

Example: For ($i=0; i < 1000; i++$)
 $a[i]=0;$

few instructions executes the loop 1000 times.



If we move mfu references to fast SRAM (Cache), we can improve performance of 90% of CPU time by a factor of \approx DRAM speed / SRAM speed ≈ 10

→ How about the 10%? Amdahl's Law?

Two types of Localities!

1- Spatial locality: (near in space) references close by in space (address) are more likely to be referenced.

$a[i]$ and $a[i+1]$

2- Temporal locality: (in time) items referenced now are more likely to be referenced again soon
 $a[i] = i$ instructions.

Cache Terminology

HIT rate: percentage of references found in cache

Miss rate: $= (1 - \text{hit rate})$; references NOT found in cache and must be retrieved from main memory

Hit Time: Access time to cache + time to determine hit/miss

Miss Penalty: Time to replace a block in cache from main memory + time to deliver reference to processor

Block: Unit of information to be transferred (present or not present) between two levels

Example of a cache parameters:

Cache size = 32 Kbyte, hit rate = 95%,
 Cache time = 10 ns, miss penalty = 100 ns,
 Block size = 32 Bytes.

Direct Mapped Cache

- Mapping, finding a block in cache
- Organization
- Performance

Mapping

Each item in memory maps to one location in cache.

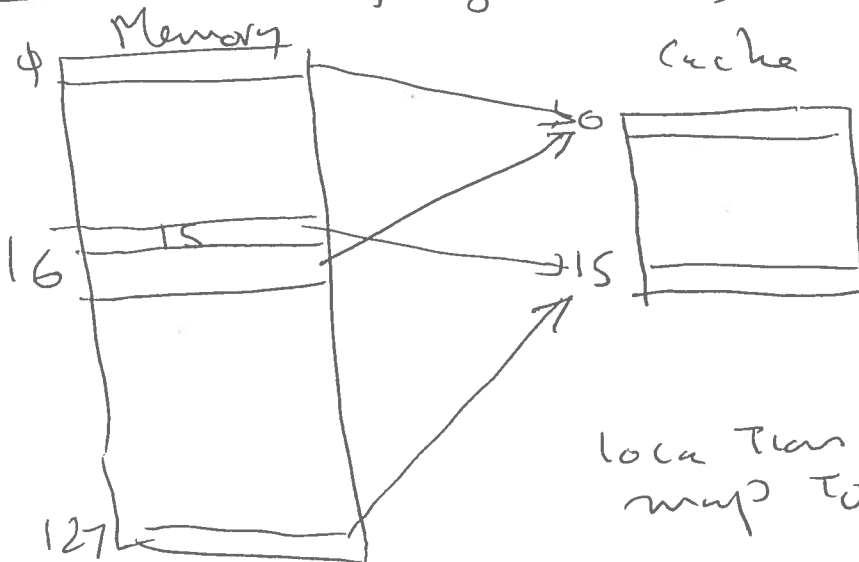
This location = (Memory Address) modulo of number of cache blocks

Finding a block in cache:

Use a tag for each block, and compare new access tag with the stored tag in cache. If it is the same = cache hit

Example

memory size = 128, cache size = 16



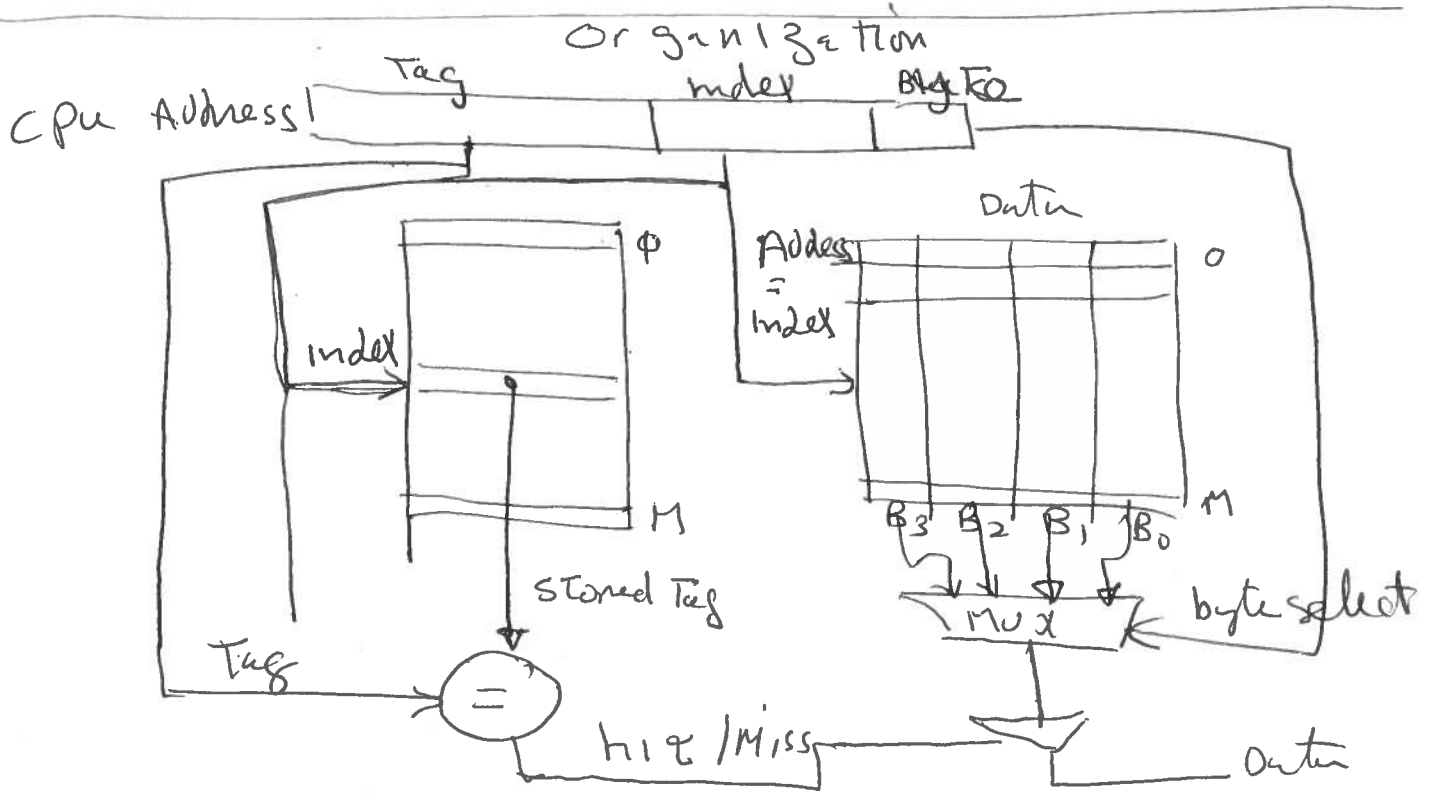
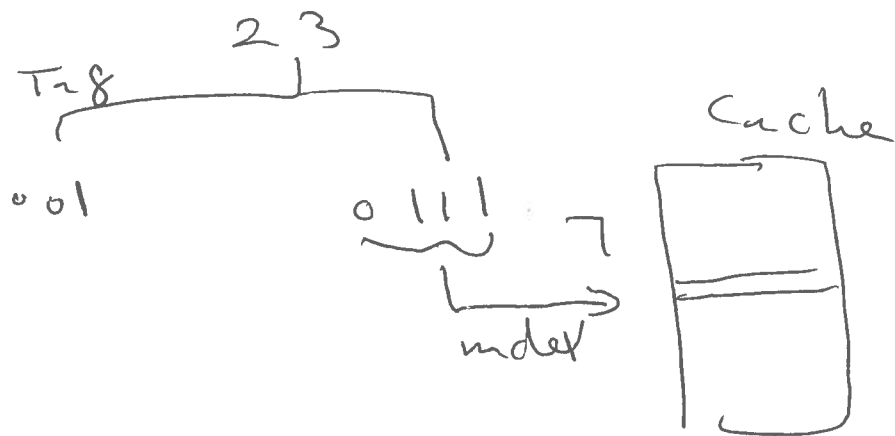
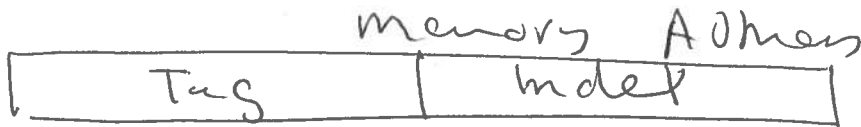
locations 0, 16, 32, ...
map to cache loc 0

locations 15, 31, ... map to 15

of pages or Frames = $\frac{128}{16} = 8$ 5

Address 23 maps to 23 module 16
 = 7

Find if block in cache using
 Tag (Page #)



Operation

if $\text{Tag}(\text{index}) = \text{Tag of current address}$
hit = true
data = Data(index)

both Tag, Data are SRAM cache
memory.

select byte by byte sel

else

hit = false

data = M[Tag | index | byte]

Cache Data(index) = M[Address]

also $\text{Tag}(\text{index}) = \text{Tag new}$

Example what is the cache organization for
1Kbyte direct mapped, blk = 32 byte,
Find Tag, index of loc 14020 hex

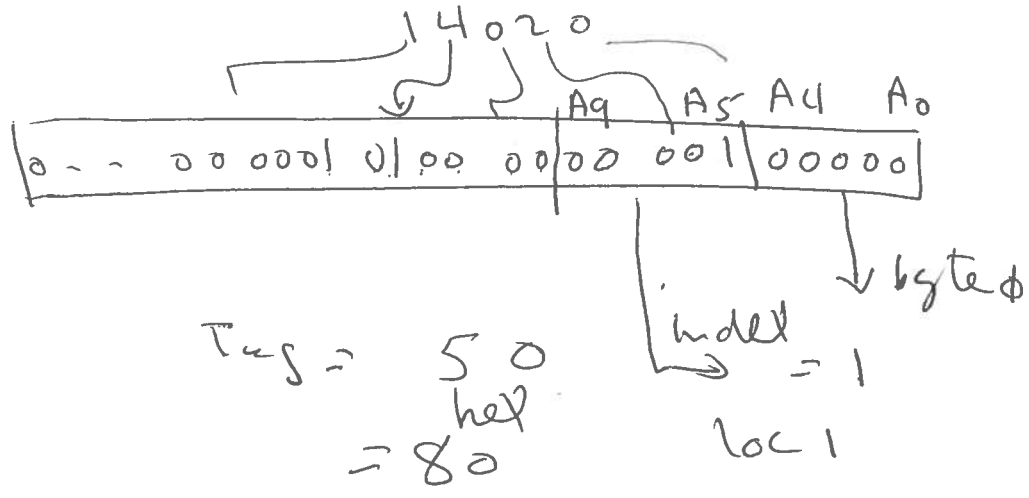
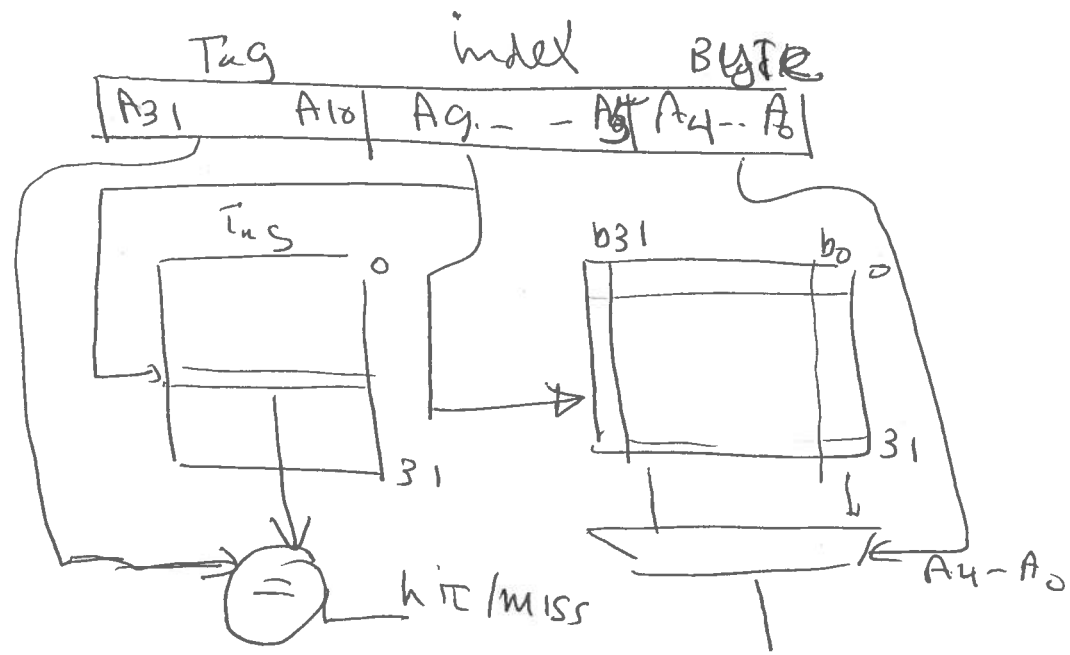
$$\# \text{ of blocks} = \frac{\text{size}}{\text{block size}} = \frac{1024}{32} = 32 \text{ blocks}$$

$$\# \text{ of Address lines for byte sel} = \log_2 32 = 5$$

$$\# \text{ of Address lines for index} = \log_2 32 = 5$$

$$\# \text{ of Address lines for Tag} = 32 - 5 - 5 = 22 \text{ bits}$$

Assuming a 32 bit CPU Address



Cache Performance

Average memory Access Time
 = Hit Time + Miss rate × Miss penalty

Miss rate = (1 - hit rate)

HIT Time = Time to access cache + Time to Find
 HIT/MISS

Miss penalty = time to access Main Memory +
 time to transfer block to cache

Time of transfer = $\frac{\text{Block size}}{\text{bus size}} \times \frac{1}{\text{Bus Speed}}$

Example: Find Average access time for a cache system if hit rate = 90%, cache is 10 ns, processor is 200 MHz, DRAM access time = 100 ns, cache block size = 8 words, BUS width = 1 word, speed = 100 MHz.

$$\begin{aligned} \text{Average Time} &= \text{Cache hit time} \\ &+ \text{Miss rate} \times (\text{DRAM access} + \text{transfer time}) \\ &= 10 + 0.1 \times (100 + 10 \times 8) = 10 + 18 = 28 \text{ ns} \\ &\approx \text{about } 5 \text{ CPU cycles} \end{aligned}$$

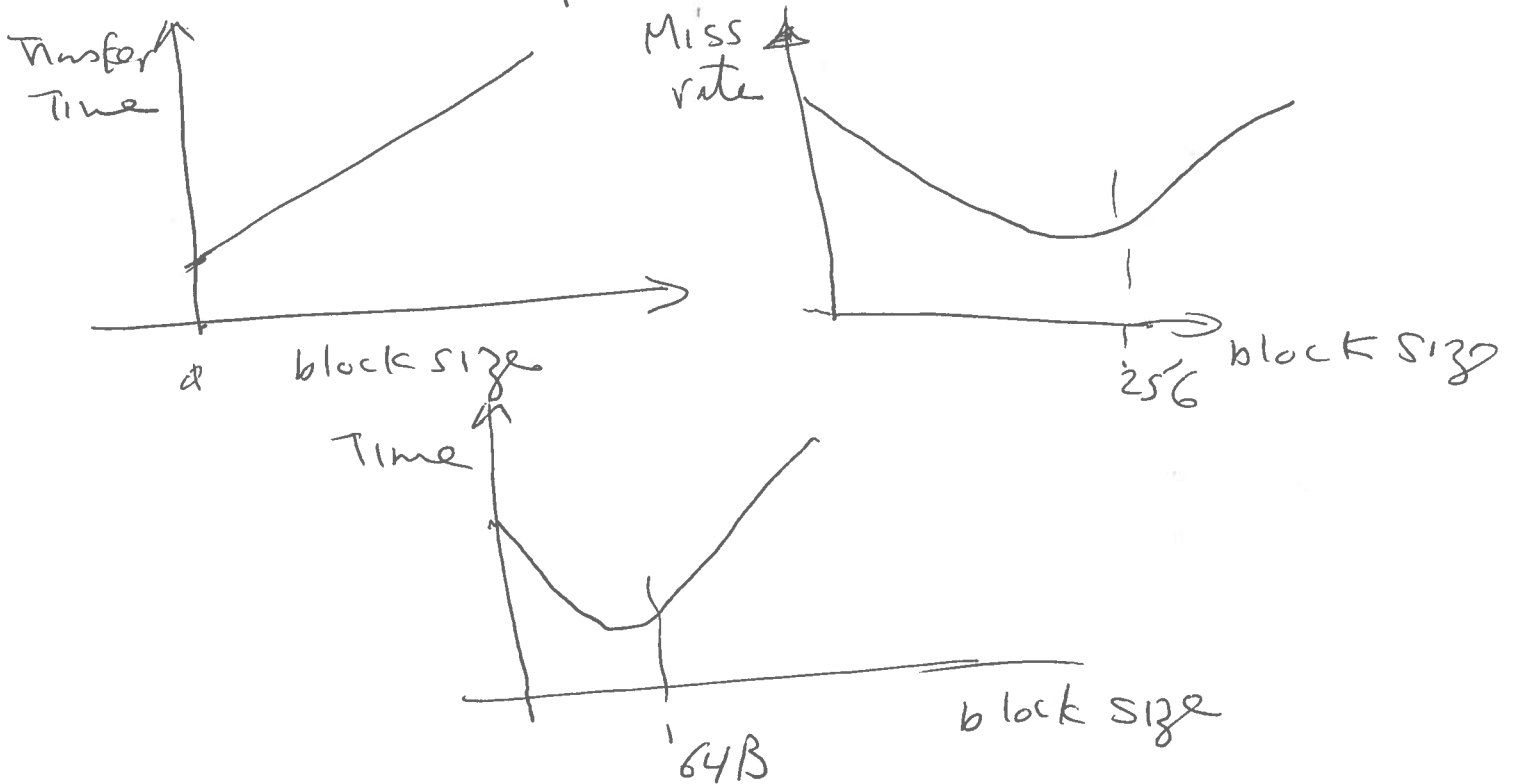
IF instructions, Data have different miss rates

$$\begin{aligned} \text{Time} &= \% \text{ Instruction hit time} + \text{Inst. Miss} \times \text{Memory penalty for Inst.} \\ &+ \% \text{ data} \times (\text{data hit time} + \text{Data Miss} \times \text{Memory penalty}) \end{aligned}$$

→ Memory penalty for Inst might be diff
 than for Data (if I cache, D cache have different Block size)

Block Size Trade-off

- Increasing block size improves hit rate because of spatial locality (to a degree)
- Increasing block size increases transfer time
- There is a point of optimum performance



Two Way Associative Cache

- Motivations
- Mapping
- organization
- Advantages and disadvantages

Two way set Associative

high

Motivations:- \downarrow Conflict misses on direct mapped cache because there is only one location for each address, if taken it has to miss

Solution! Use two way set associative cache. Each address has two locations to map to.

In a loop, if 2 instructions map to same cache location in direct mapped cache, with two way set associative, the 2 instructions will map to diff cache ~~to~~ blocks. ~~in~~

Example Assume Memory size = 128

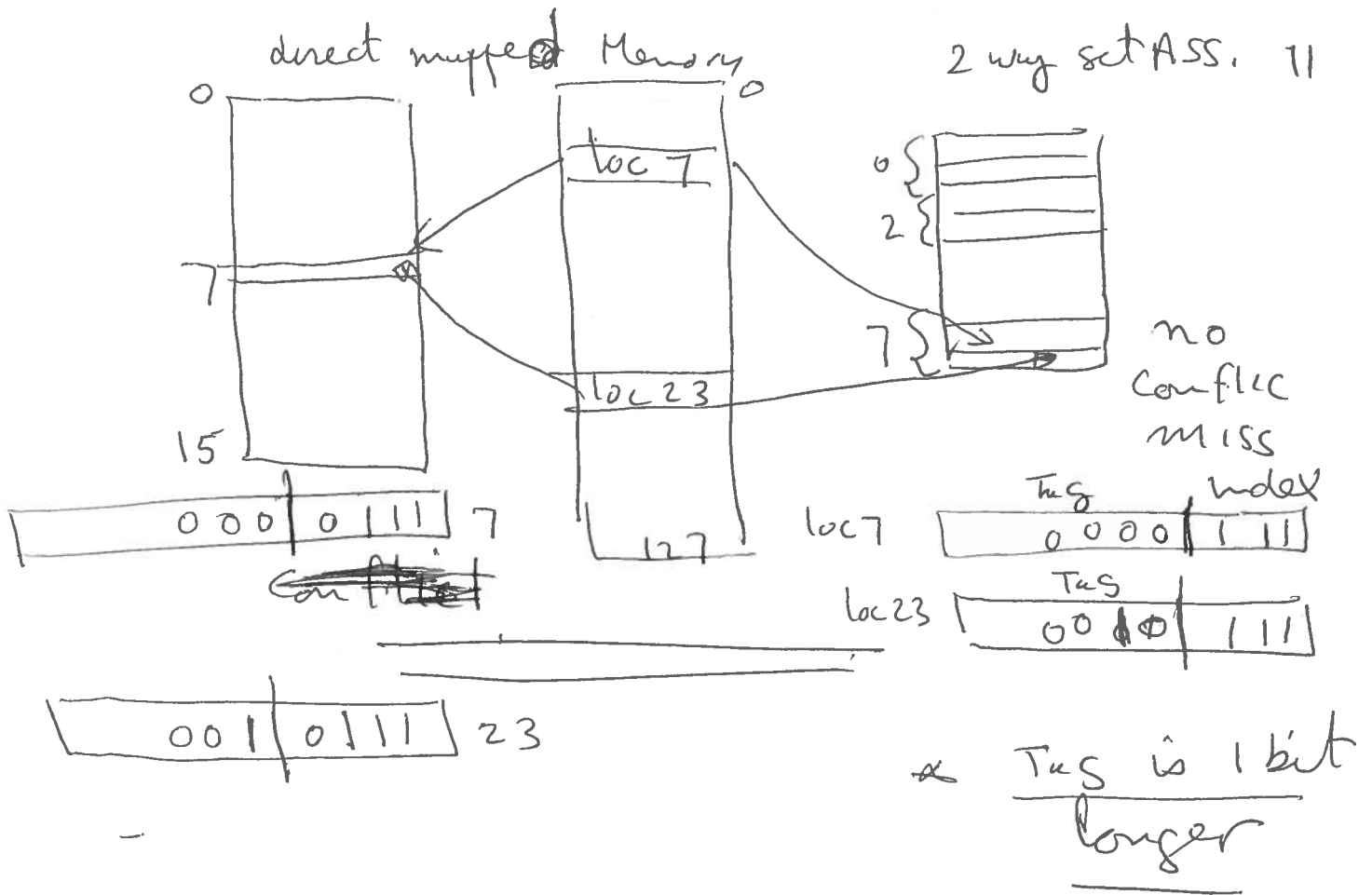
Cache size = 16

Compare direct mapped to 2 way set associative for decoding

loc 7, loc 23

1- direct mapped $\left. \begin{array}{l} \text{Loc 7 maps to } 7 \bmod 16 = 7 \\ \text{Loc 23 maps to } 23 \bmod 16 = 7 \end{array} \right\}$

2- 2 way set ASS. $7 \bmod 8 = 7$ 1st
 $23 \bmod 8 = 7$ 2nd

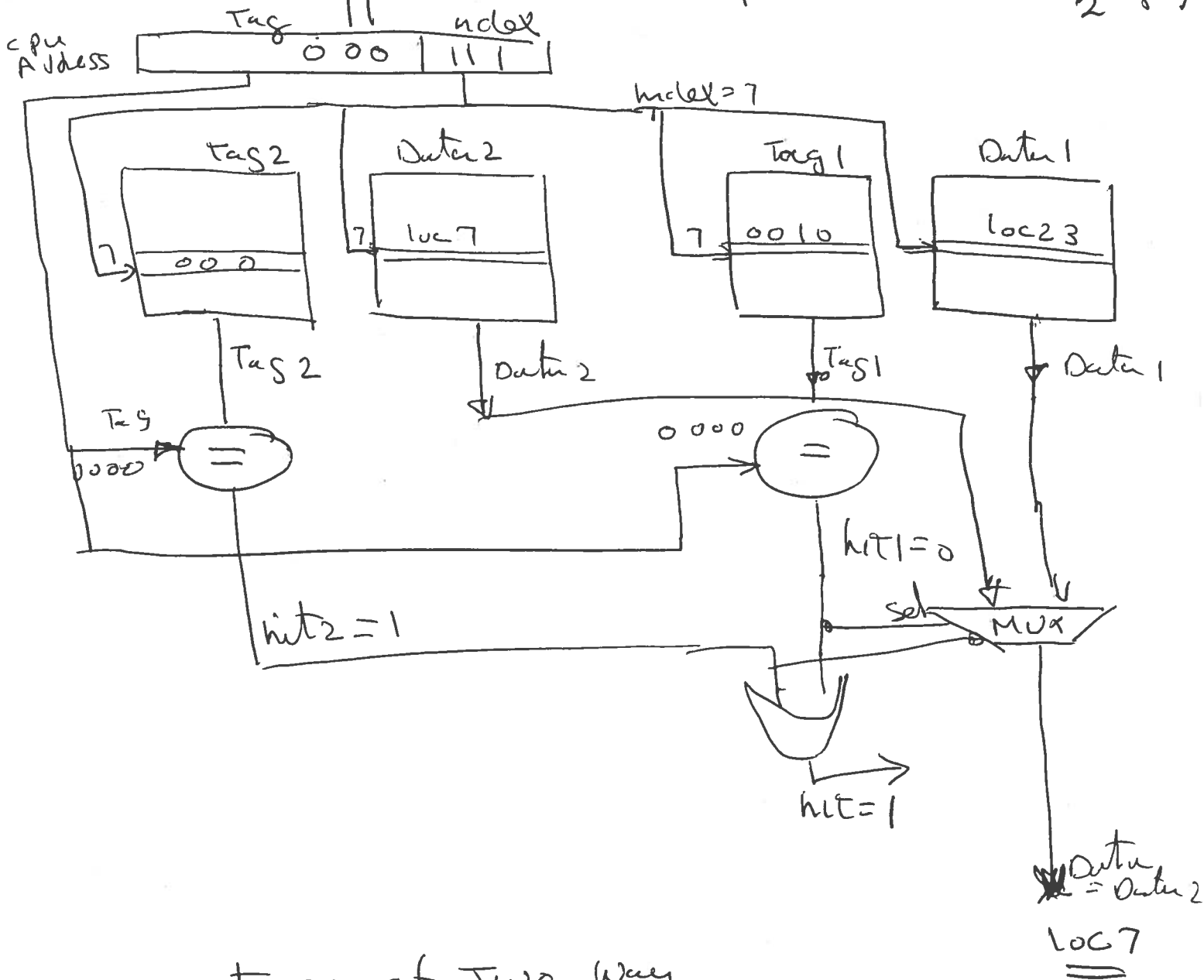


Two way set ASSOCIATIVE
cache organization

2 direct mapped cache in
parallel

Two Way Set Associative organization

2 direct mapped cache in parallel (each 1/2 size)



Advantages of two way set ASSOCIATIVE

= low miss rate (Conflict miss)

Disadvantages of Two Way Set Associative: 13

- 1- Larger Tag size ($1 \text{ bit} \times \# \text{ of loc}$)
- 2- Use of 2 Comparators (one per each tag)
- 3- Search both tags (slow)
- 4- Using a MUX in the critical time path (slower)
- 5- Must wait for hit/miss before can output data from MUX (compared to direct mapped that can output data while finding hit/miss)

Example: Find miss rate of the following sequence: 09, 89, 09, 69, 8 for a cache size = 4 loc., if it is a direct mapped and if it is a two way set Associative.

i- Direct mapped

Access ϕ maps at $\phi \pmod{4} = \phi$
- First time \longrightarrow miss
Cache has loc ϕ at ϕ

B - access 8 maps to $\frac{8}{4} = 2, \phi \rightarrow \phi$ 14
 (miss) cache has loc 8 at ϕ

C - ϕ access maps to loc ϕ
 miss \rightarrow cache will have
 loc ϕ

D - 6 maps to $\frac{6}{4} = 1, 2 \rightarrow$ loc 2
 First time = miss (loc 2 has 6)

E - 8 maps to ϕ , but ϕ has loc
miss

100% miss rate

Two Way Set Association

A - 0 \rightarrow Tag 1, Loc ϕ Tag 1
 $\begin{array}{|c|c|} \hline 000 & 0 \\ \hline \end{array}$

B - 8 \rightarrow loc ϕ of Tag 2
 $\begin{array}{|c|c|} \hline 100 & 0 \\ \hline \end{array}$

C - 0 \rightarrow loc ϕ of Tag 1
 hit \checkmark

D - loc 6 \rightarrow ϕ $\begin{array}{|c|c|} \hline 011 & 0 \\ \hline \end{array}$
 tag = 3
miss

	Loc	Tag 1	Tag 2
A	0	000	-
B	8		100
C	ϕ	000	100
D	6	000	011
E	8	100	011

QUIZ ??

LRU

E - loc 8 \rightarrow ϕ $\begin{array}{|c|c|} \hline 100 & 0 \\ \hline \end{array}$
miss

Tag 2

QUIZ TAG 1

miss rate 80% compared to 10¹⁵%

Replacement policy

- 1- No replacement for direct mapped cache (only one block to be replaced "no choice")
- 2- Need a replacement policy for Associative Cache, Each memory location have more than one choice (N-ASS).

QUESTION? which block to throw out?
for the new coming access.

A- Random replacement.

cost less but it is not effective.

The cache block to be replaced is selected in random

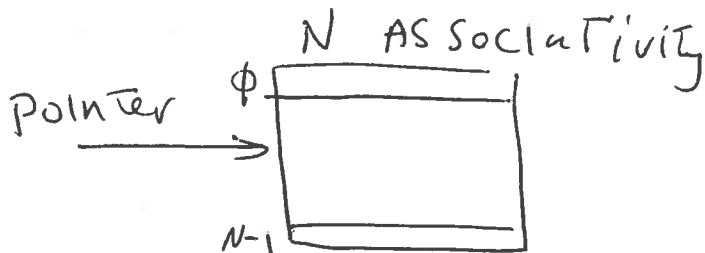
B- Least recently used block LRU
Throw block that has not been used for longest time (temporal locality).

Need Hardware to keep track of access history.

Implementation for an LRU

16

use a pointer that always points to LRU



Protocol: Each time, an access is used, the access moves pointer away to next entry.

If access is not used -- do not move pointer (LRU)

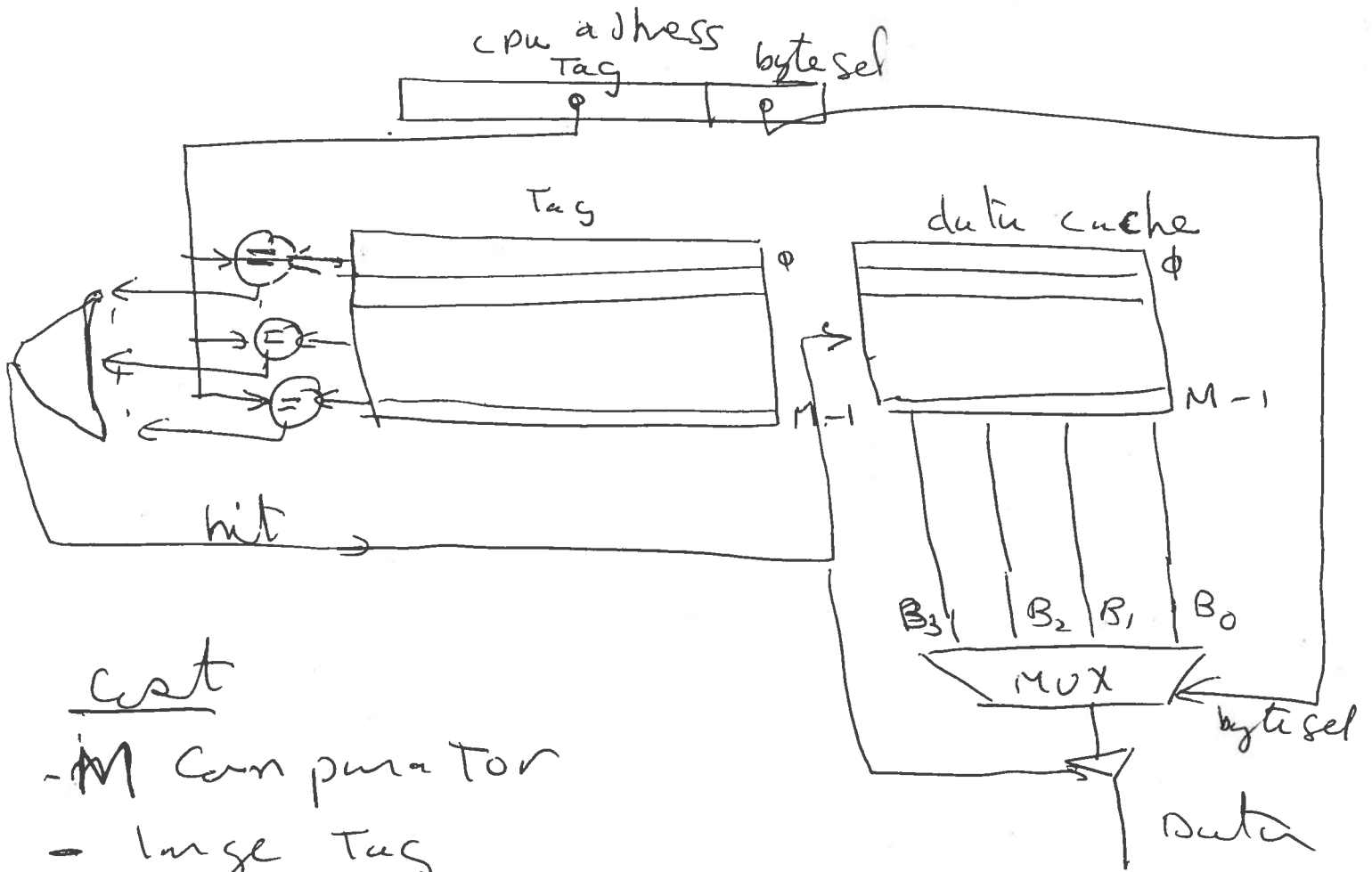
Two way set associativity pointer = 1 bit

Fully Associative Cache

Concept: any memory location can map in any cache location (Conflict miss is reduced to zero).

This means that the index field = ϕ
and tag = length of Memory Address

Fully Associative Cache Organization ¹⁷



cost

- M Comparators
- large Tag
- parallel search of all Comparators (slow)

Example Design Fully Associative cache = 1KB
if Block size = 32 byte.

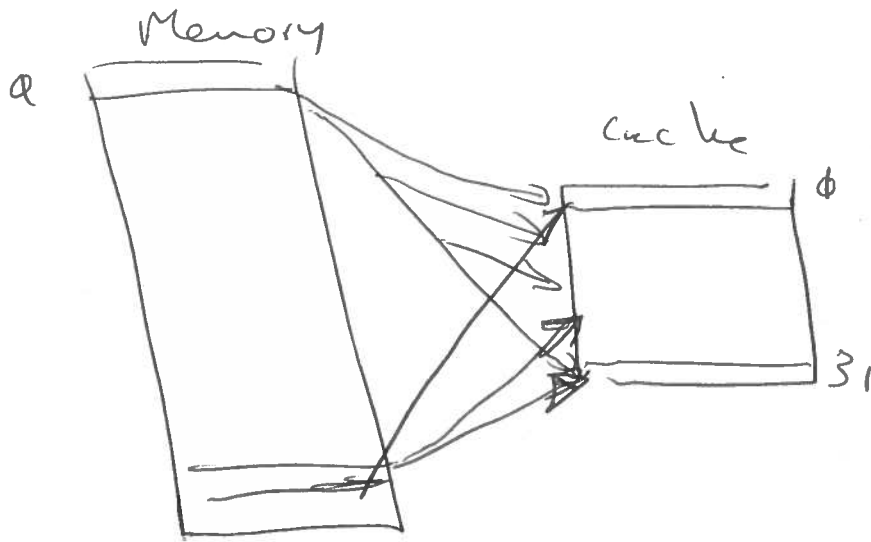
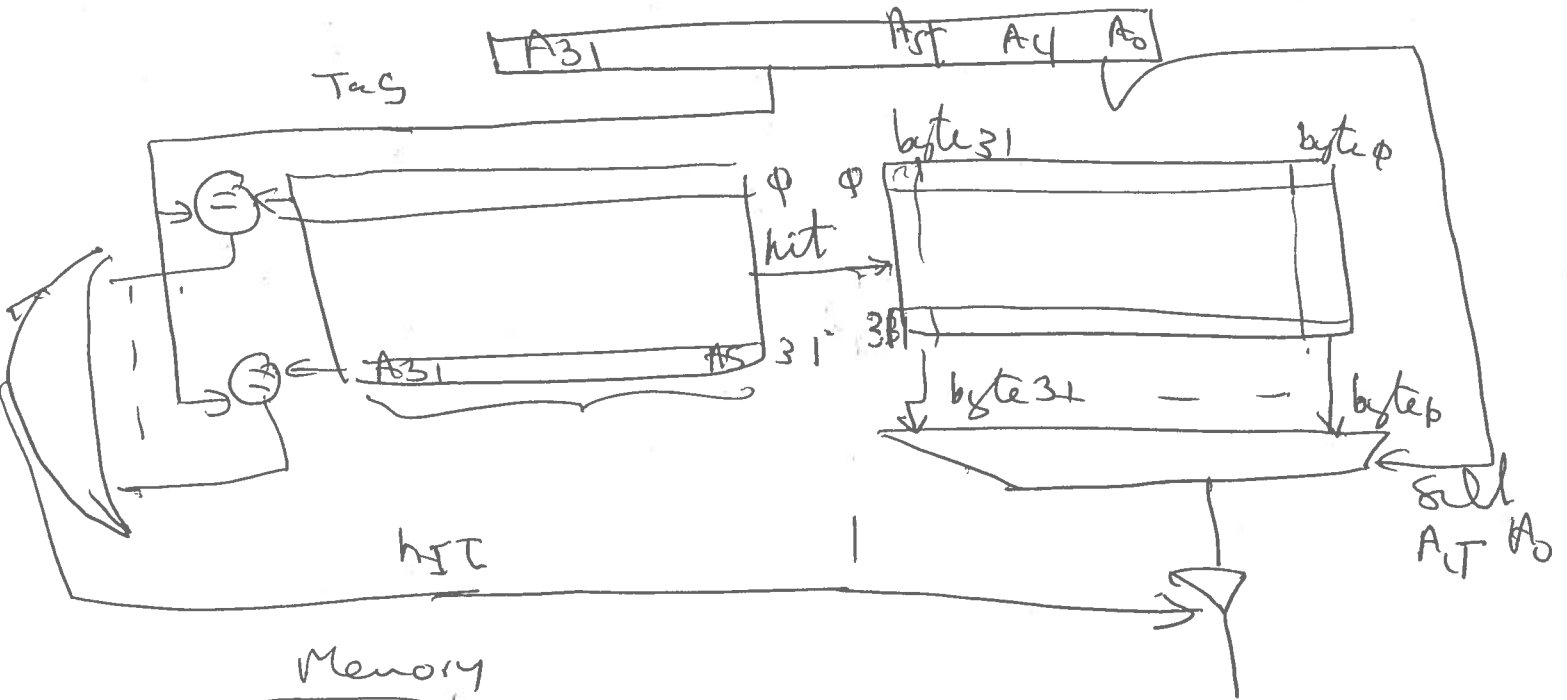
$$\# \text{ of byte sel bits} = \log_2 32 = 5 \text{ bits}$$

$$\text{Tag} = 32 - 5 = 27 \text{ bits}$$

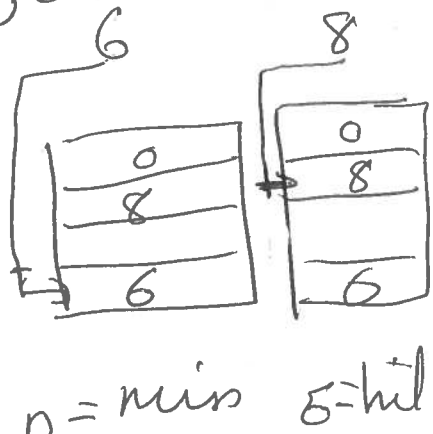
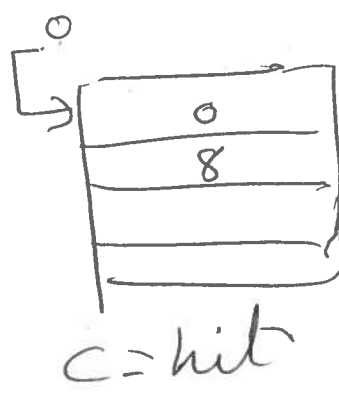
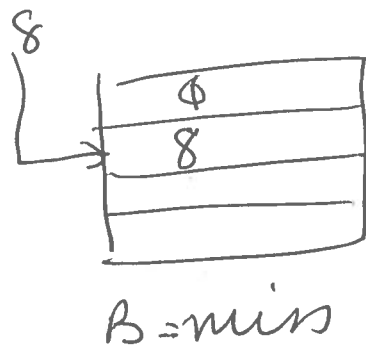
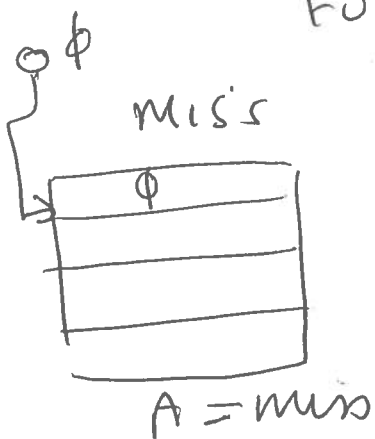
$$\# \text{ of blocks in Cache} = \frac{1024}{32} = 32 \text{ blocks}$$

$$\# \text{ Cache Comparators} = 32$$

$$\text{Tag} = 32 \times 27, \text{ Data} = 32 \times 32 \text{ byte}$$



Example Program uses addresses 0, 8, 0, 6, 8
 cache size = 4 locations, use
 Fully Associative



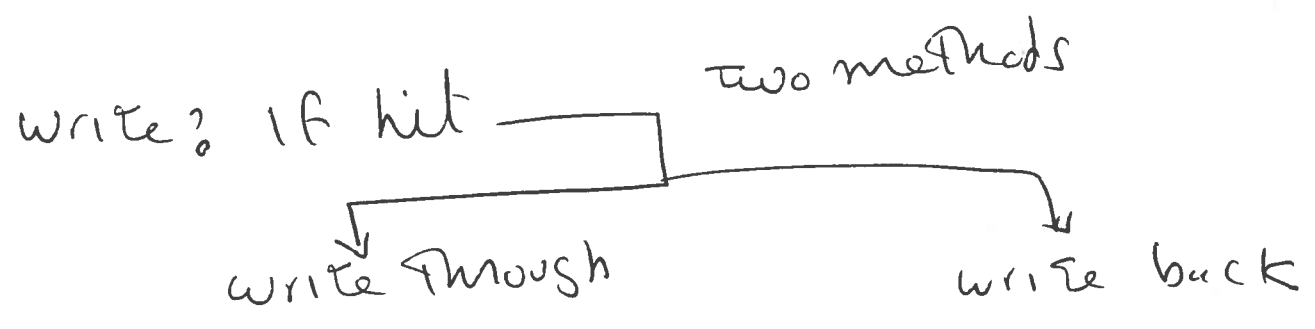
miss = 60% only

Cache write policy

Read: if hit → read from cache

if miss → read from memory
replace (LRU block) by
data from memory

if replaced block is modified (dirty) need
to write it back to memory.



write through: write to cache and memory

Advantages: cache content is the same
(Coherent with memory) in memory.

disadvantage: write to memory is slow

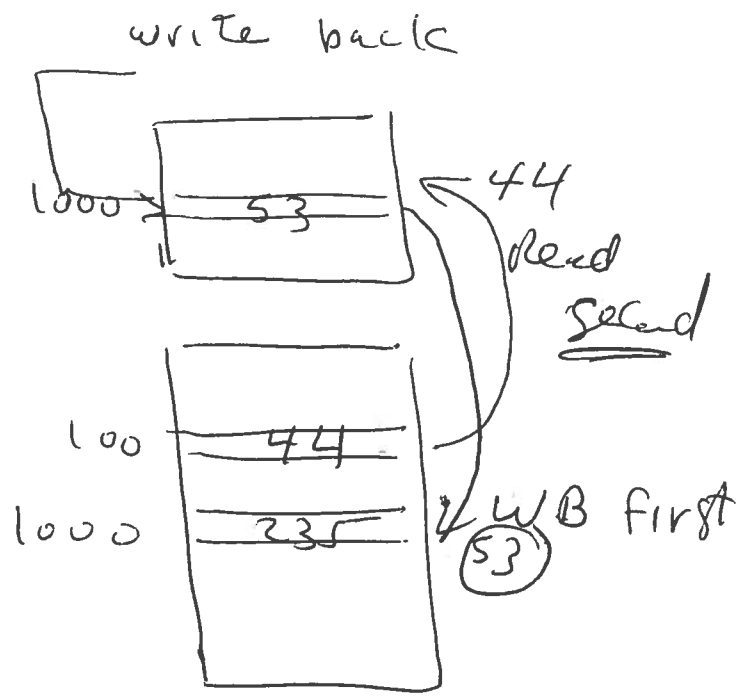
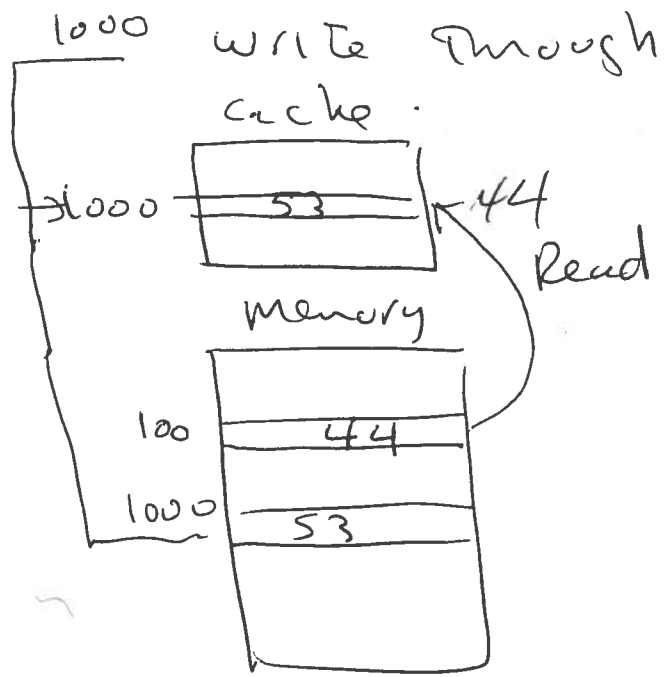
write back: only write to cache, mark
it dirty, so when it is
replaced (in read-miss), it must
be written back to memory.

cost: 1 dirty bit per block

Example?

write 53 at 1000

read 44



write Miss

NWA
 No write allocate

write only to memory
 (block is not in cache)
 Miss rate is worse

WA
 Write allocate
 write to memory,
 read miss to
 move block to
 cache
 improve miss
 rate

Improving Cache Performance

Average memory access time

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

Improving cache performance with improving:

- 1- Hit Time
- 2- Miss rate
- 3- Miss Penalty

- Split cache (I & D) versus unified cache Page 384

Example Compare miss rate of split cache (16KB I + 16KB D) to 32KB unified?

if 16KB I cache miss rate = 1.64%

16KB D cache miss rate = 6.47%

32KB unified miss rate = 1.99%

Assume ~~75~~ Instructions = 75%

Data = 25% of program

Average miss rate for split cache

$$= 0.75 \times 1.64 + 0.25 \times 6.47 = 1.484 + 1.6175 = 3.1015$$

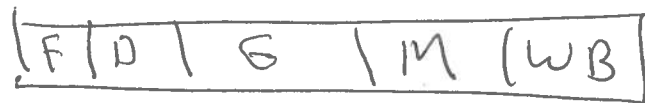
split cache gives worse miss rate than unified.

BUT performance of unified cache might be worse!

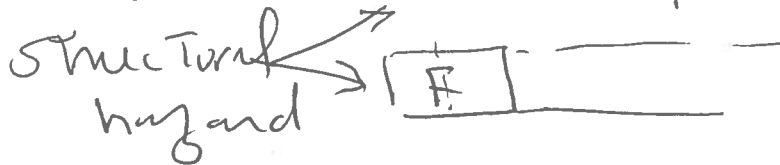
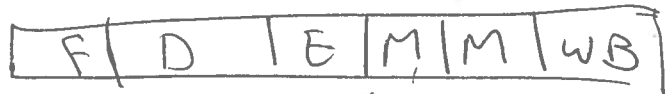
Assume that unified cache hit takes an extra clock cycle, using pipeline

the unified cache leads to a structural hazard (one port). Find Average memory Access Time for each. Assume miss penalty = 50 clock cycle, write is using a buffer with no stalls in writes.

pipeline with Read/Write operation
1 - split



2 - unified (Read/Write) only



Performance of split cache

$$\begin{aligned}
 &= \% \text{Instruction} (\text{hit time} + \text{miss rate} \times \text{miss penalty}) \\
 &+ \% \text{Data} (\text{hit time} + \text{miss rate} \times \text{miss penalty}) \\
 &= 0.75 - (1 + \frac{1.64}{100} \times 50) + 0.25 (1 + \frac{6.47}{100} \times 50) \\
 &= 2.05 \text{ clock}
 \end{aligned}$$

Unified Cache ? \leftarrow not the same \rightarrow ?

$$\begin{aligned}
 &= 0.75 - (1 + 1.99 \times 50) + 0.25 (1 + 1 + 1.99 \times 50) \\
 &= 2.24
 \end{aligned}$$

slower but has better hit rate

Improving cache performance 23

Method 1 Reducing cache misses using
Larger Block size

types of misses?

- 1 - Compulsory - Cold start misses
for first reference misses (Block)
- 2 - Capacity - Cache cannot contain
all blocks (size)
- 3 - Conflict - blocks map to same
location in cache (ASS)

Larger Block size reduce Compulsory misses
because of spatial locality.

Problems increase transfer time and
may increase conflict misses
(few blocks in cache)

Example Assume a 1KB cache with
the following miss rate:

Block size 16B = 15.05% , Block 32B = 13.34%

Block 64B = 13.76% , Block 128B = 16.64%

if memory access time = 40 cycles,

transfer time = 16B for 2 clocks.

Find the Best Block size for memory performance.

Average memory access time =

miss penalty = memory access time + transfer time
Hit time + miss rate * miss penalty

$$\begin{aligned}
 TM \text{ for } 16B &= 1 + .1505 \times (40 + 2) = 7.32 \text{ cycles} \quad \checkmark \checkmark \\
 \checkmark \checkmark 32B &= 1 + .1334 \times (40 + 4) = \underline{6.87} \text{ cycles} \\
 64B &= 1 + .1376 \times (40 + 8) = 7.605 \text{ cycles} \\
 128B &= 1 + .1664 \times (40 + 16) = 10.318 \text{ cycles}
 \end{aligned}$$

Method 2: Higher Associativity

improve conflict miss. For size N direct mapped cache, 2-way set associative of size $\frac{N}{2}$ has about same miss rate ~~at~~^{at} $\frac{1}{2}$

Example assume a 32 KB cache, with the following miss rates

$$\begin{aligned}
 1\text{-way} &= .020, & 2\text{-way} &= .014, & 4\text{-way} &= .013 \\
 8\text{-way} &= .013
 \end{aligned}$$

$$\begin{aligned}
 \text{miss penalty} &= 50 \text{ cycles, hit time for} \\
 1\text{-way} &= 1 \text{ cycle, } 2\text{-way} = 1.1, & 4\text{-way} &= 1.12 \\
 \text{and } 8\text{-way} &= 1.14
 \end{aligned}$$

Find best cache design.

$$\text{Average Memory Access Time} = \text{Hit Time} + \text{Miss rate} \times \text{miss penalty}$$

$$1\text{-way} = 1.0 + .02 \times 50 = 2.0$$

$$2\text{-way} = 1.10 + .014 \times 50 = 1.8$$

$$4\text{-way} = 1.12 + .013 \times 50 = \underline{1.77} \quad \checkmark \checkmark$$

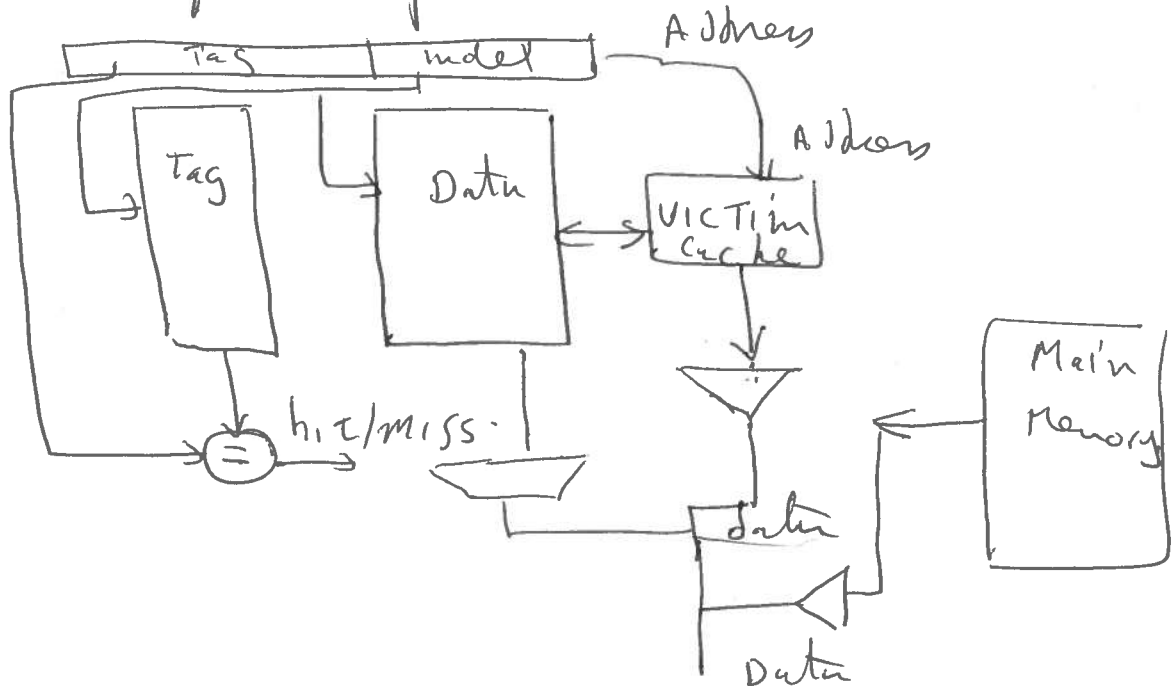
$$8\text{-way} = 1.14 + .013 \times 50 = 1.79$$

Method 3 Using VICTIM Cache

25

- Reduce miss rate without affecting clock cycle time (ASS), or miss penalty (Block).

Add a small fully Associative victim cache for replaced blocks.



ON a cache miss, check VICTIM cache rather than going to lower level (main memory). If VICTIM cache = hit, swap block with one in cache. It will be a cache hit, replaced block on miss is transferred to VICTIM cache.

IT improves conflict miss by 20-95%.

Method 5: Hardware Prefetching of Instructions and Data

28

Instruction Prefetch: on a miss fetch 2 blocks, place one in cache and other in stream buffer.

Instruction stream buffer improves hit rate by 50%.

Data prefetch: use stream buffer and prefetch data at different addresses

- Need accurate prediction for Prefetch.

- Possible Precomputation.

while waiting for Prefetch, processor can proceed and cache can continue supplying data (nonblocking cache).

Method 6: Using Compiler

- Prefetch by inserting instructions to request data before they are needed.

has some overhead, of the inserted instructions must not exceed the benefits of prefetch.

Example: Assume 8KB direct mapped with BLK size = 16B.
Using write back with WA.

assume $a[3][100]$, $b[10][3]$ each element = 8B find miss rate, prefetch instructions to improve miss rate for

Reducing Cache Miss Penalty 26

important as cost of miss penalty increases due to speed gap between CPU / DRAM

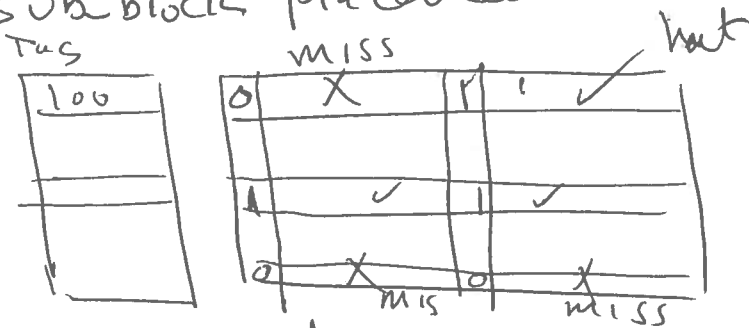
Method 1: Giving priority to Read Miss over writes

write through

using write buffer, and allow read miss to continue even write buffer is not empty as long as there is no conflicts (read same location).

write back on read miss with dirty block to be written back to a buffer, then read memory, then write buffer back to memory.

Method 2: subblocks placement



on miss transfer only sub-blocks, mark it valid. need a bit for each sub-block
For a hit \Rightarrow TAG match and valid bit = 1

reduce transfer time of large blocks.

Reducing Cache Miss Penalty

~~28~~
28

Method 3: Early restart and critical word first.

- No Hardware is required
- only one word in a block is required by CPU

two methods

1- Early restart: CPU waits for requested word, let CPU continue execution when word arrives

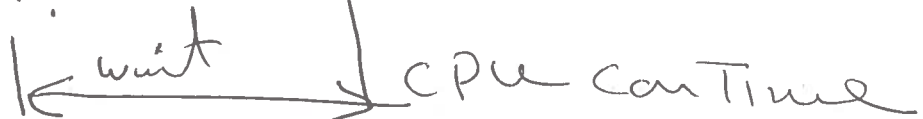
2- critical word first: send requested word first from memory, let CPU continue execution when word arrives (use wrapped fetch).

Example

ld w2

lwo | w1 | w2 | w3

Early restart:



CPU Execution

Critical word first

w2 | w3 | w1 | w0



Method 4: Non Blocking Caches

~~20~~

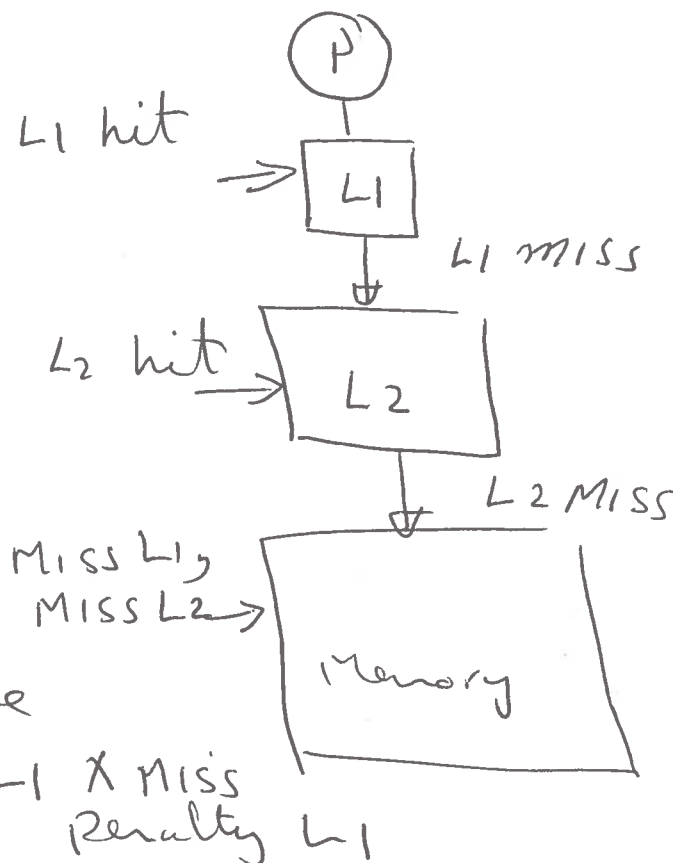
For pipelined machines that allows out of order execution and completion, CPU could continue fetching instructions from instruction cache while waiting for data cache miss. It also allows data cache to supply data while waiting for a miss data to arrive from memory.

- need scoreboard or Tomasulo control.

Method 5: Second Level Caches

on a first level cache miss, can use a faster level 2 cache and not the memory.

only misses from L1 could be found in L2.



Average memory access time

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

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

Average memory access time

$$= \text{Hit Time } L_1 + \text{Miss } L_1 \times \text{Hit Time } L_2 + \text{Miss } L_1 \times \text{Miss } L_2 \times \text{Memory latency}$$

29
30

Local miss rate: Miss L_1
Miss L_2

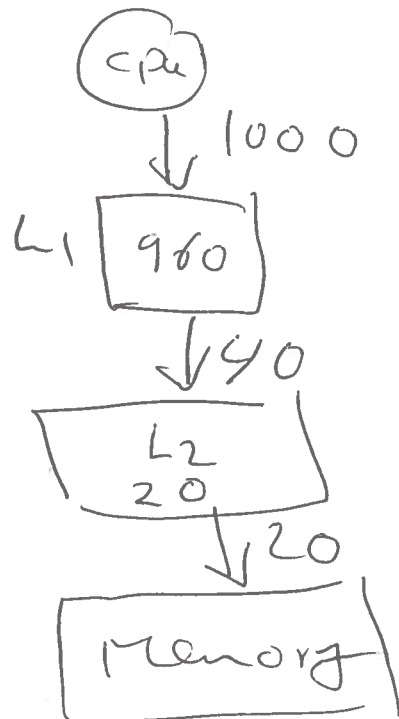
Global miss rate = Miss $L_1 \times$ Miss L_2

Example: Assume 1000 memory references that have 40 misses in L_1 , 20 misses in L_2 , find various miss rates

Answer: Miss $L_1 = \frac{40}{1000} = 4\%$

Miss $L_2 = \frac{20}{40} = 50\%$

Global Miss = $\frac{4}{100} \times \frac{50}{100} = 2\%$



Characteristics of second level cache: ~~28~~ 30

1. Larger than L1
2. Slower than L1.
3. Multi level inclusion property
" data in first level cache exist in second level.
important in multi processors,
can snoop & check second level cache
4. might need larger block size than L1
" inclusion can be maintained at a cost "
5. could use write through L1 if L2 is a write back (speed at L2)

Reducing Hit Time

- affects cpu clock rate (very critical)

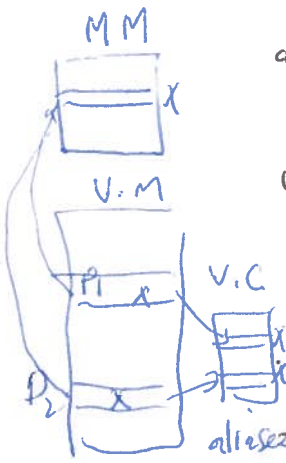
method 1: small single cache

- on chip cache (smaller is faster)
- direct mapped cache is simpler
can overlap tag check with data transmission (only one data set)

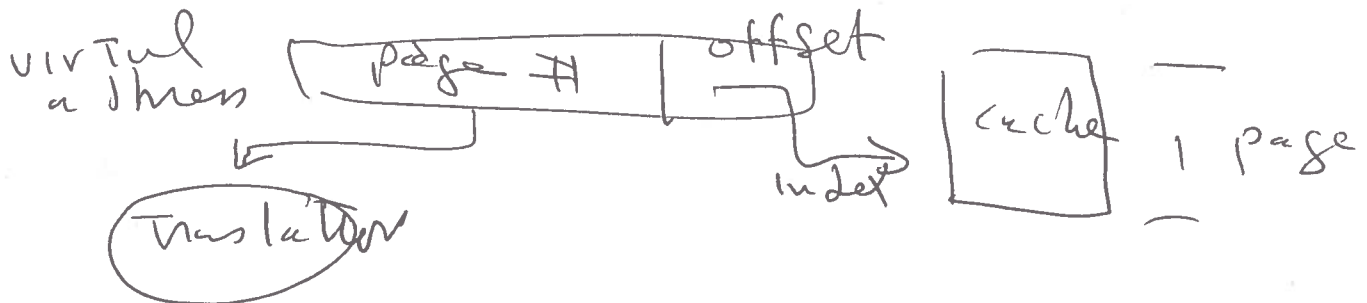
Method 2: Using virtual addresses to³¹ access cache.

- For fast and no translation, use virtual cache
- problem in process switch, must flush cache (cache contains different physical addresses)

- virtual addresses of two different addresses might use the same physical address (aliases), this will result in having two copies in virtual cache for same data. CPU then could modify one copy, the other copy will have the wrong value.



could use both virtual and physical address by using offset index to index the cache



Limitation: cache = page size
Access cache while translation of virtual to physical.