# Cache Memory

---

- What is cache memory?
- Cache Organisation
- Cache Implementation
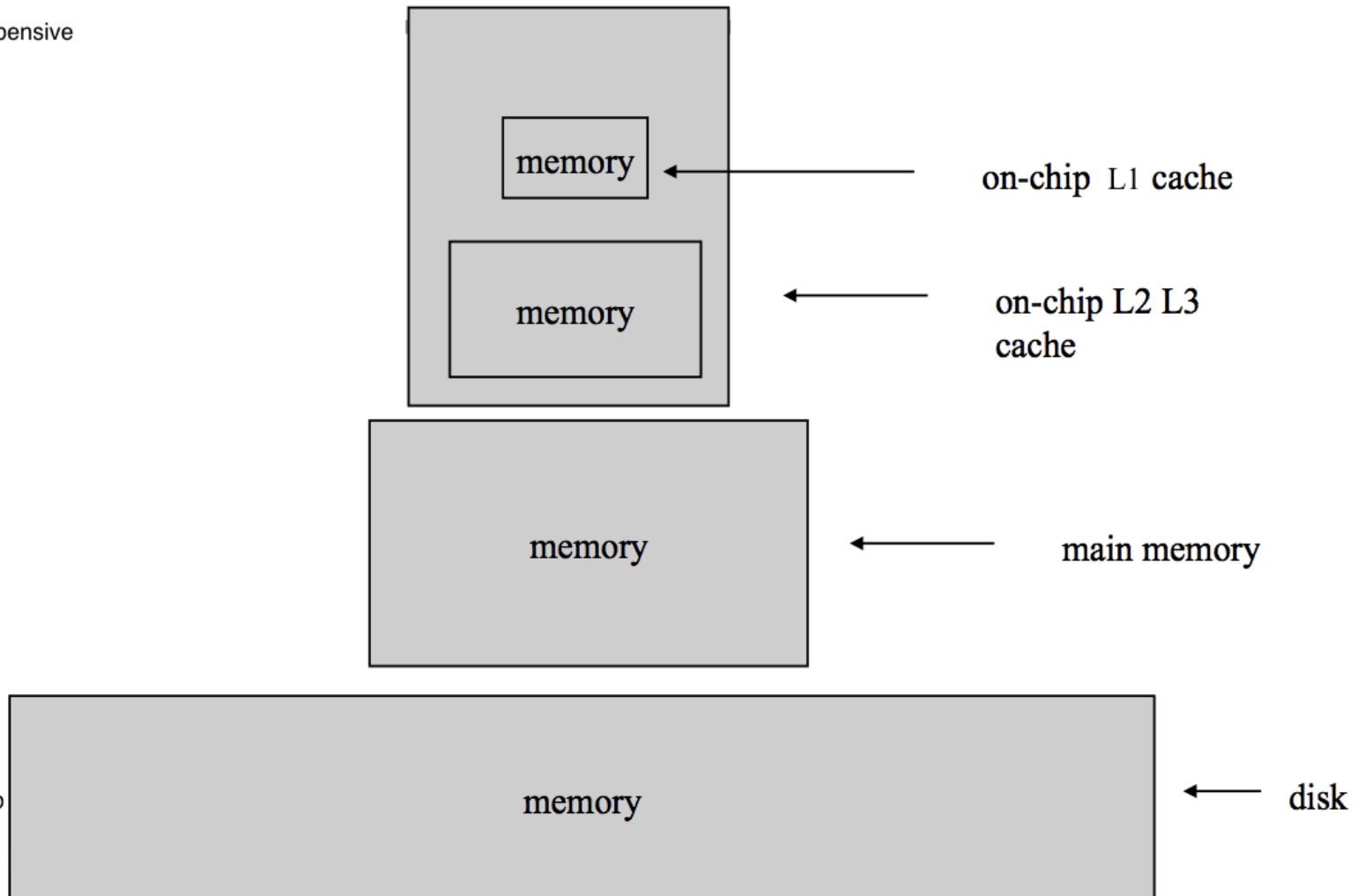- Cache Performance

EEE336/NJP/L13

# Cache Memory

---

- Early processors had an on-chip instruction queue, a modest amount of DRAM and a hard disk

- As processor speeds increased, memory access times did not keep pace and more creative solutions were required for memory interfaces.

- A multi-level memory architecture known as 'cache' provided a solution.

Small / Fast / Expensive

memory ← on-chip L1 cache

memory ← on-chip L2 L3 cache

memory ← main memory
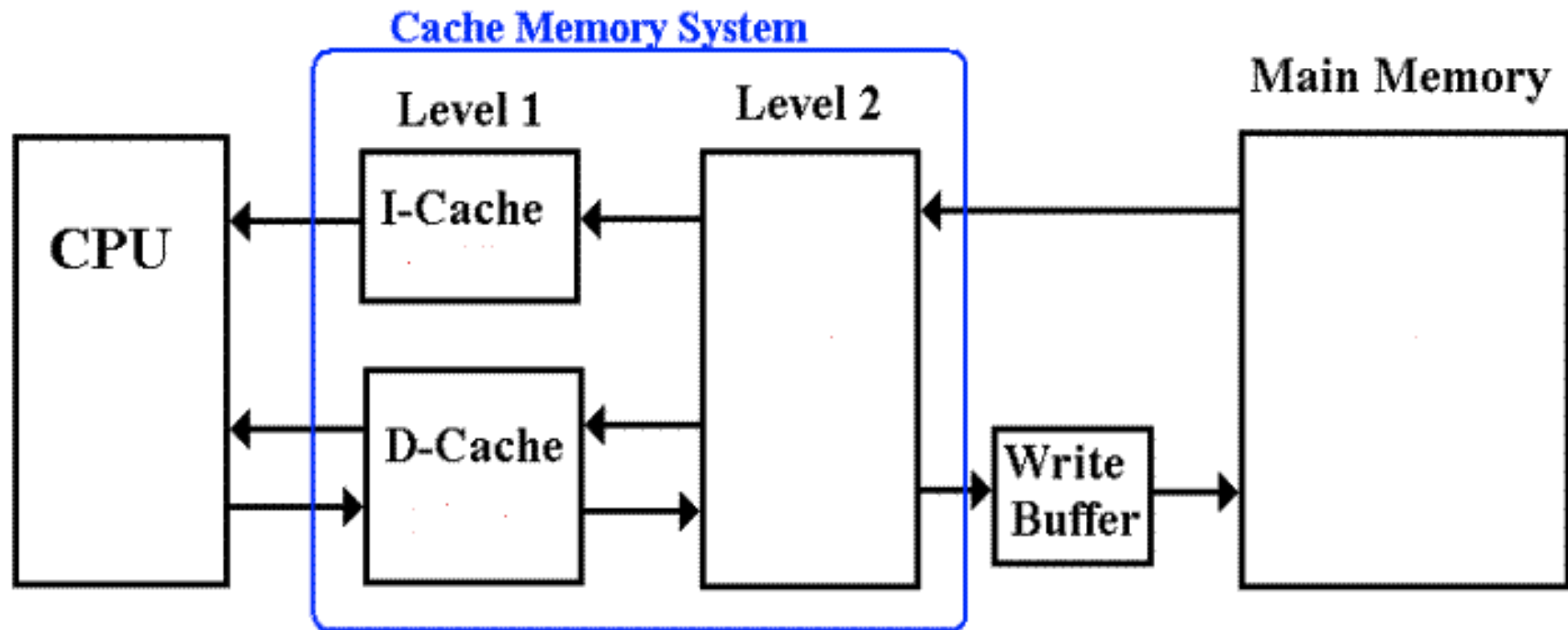
Big / Slow / Cheap     memory ← disk
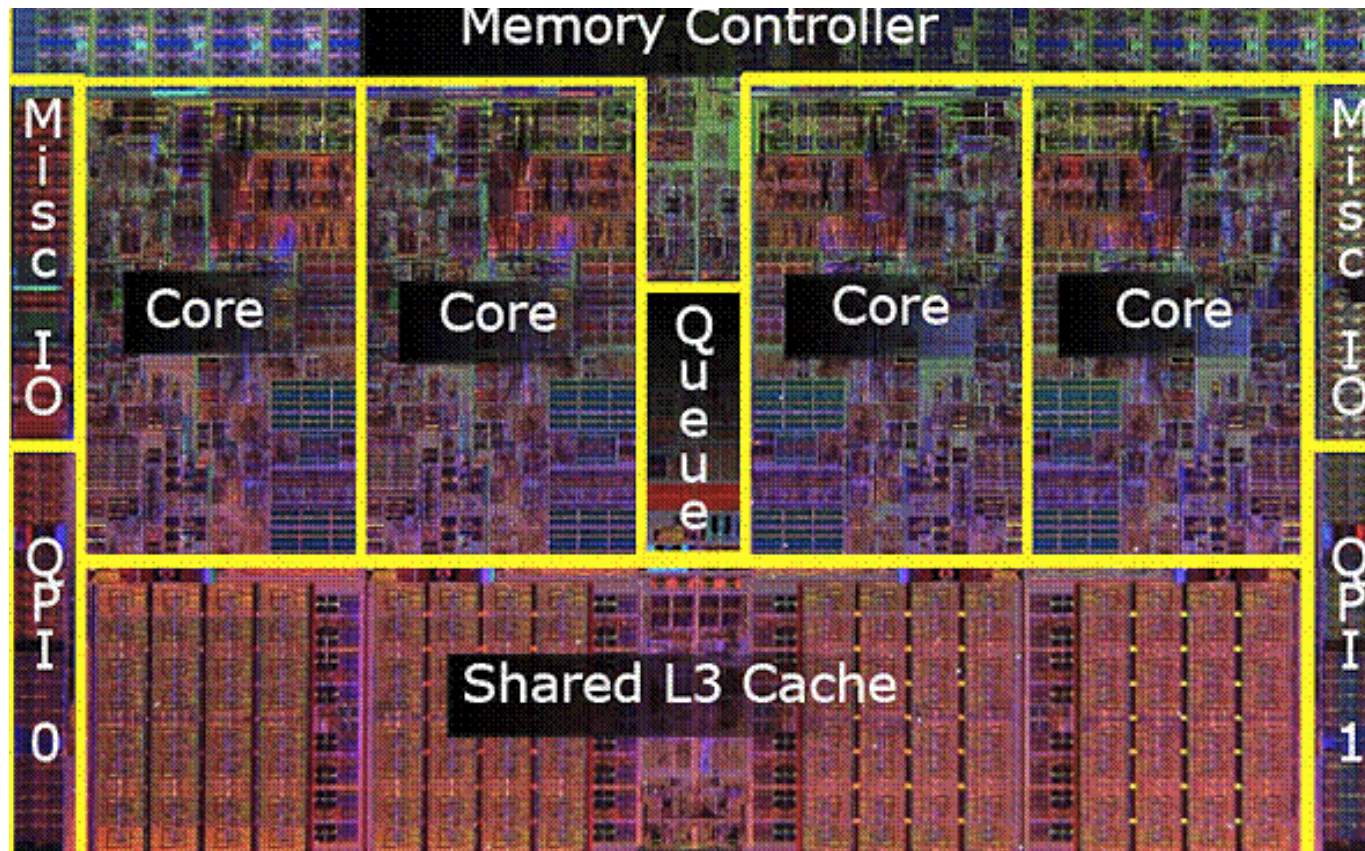
# Cache Memory

---

- Cache memory is a temporary store for data and instructions; it is faster to access than main memory
- L1 cache is a small memory on the processor chip which runs at processor speed
- L2 cache is much larger than L1 but slower. It is accessed if the required data is not found in the L1 cache; aim for L2 to have the data about 90% of the time
- Multi-core chips can have a L3 cache shared by all of the cores

# Cache Memory

**Cache Memory System**

Level 1        Level 2        Main Memory

CPU        I-Cache        
         D-Cache        Write Buffer

# Cache Memory

# Cache Definitions

---

- *cache hit* – cpu requested data is found in the next level cache

- *cache miss* – requested data is not in the cache and must be fetched from  upper levels

- *hit time* -- time to access the cache (time to determine hit/miss + memory access time)

- *miss penalty* -- time to move data from upper levels to cpu

- *hit ratio* -- percentage of time the data is found in the cache

- *miss ratio* -- (1 - hit ratio)

# Types of Organisation

---

- Direct Mapped

- Fully Associative

- Set Associative

- + few others

- Caches are arranged in blocks (or lines) just like lower levels of the hierarchy

# Direct Mapped

---

- Simplest and cheapest

- Lowest performance
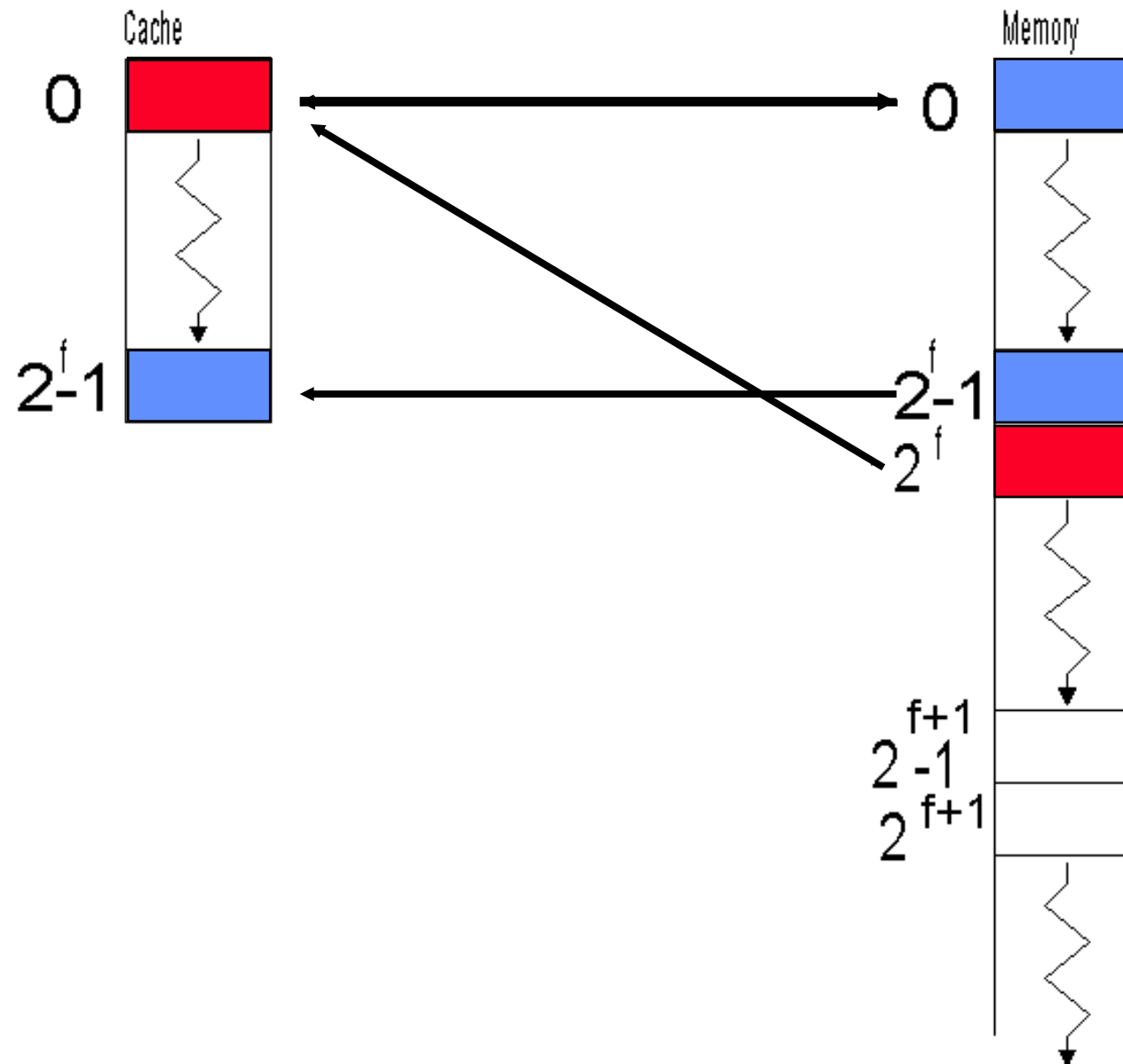
- Some significant drawbacks
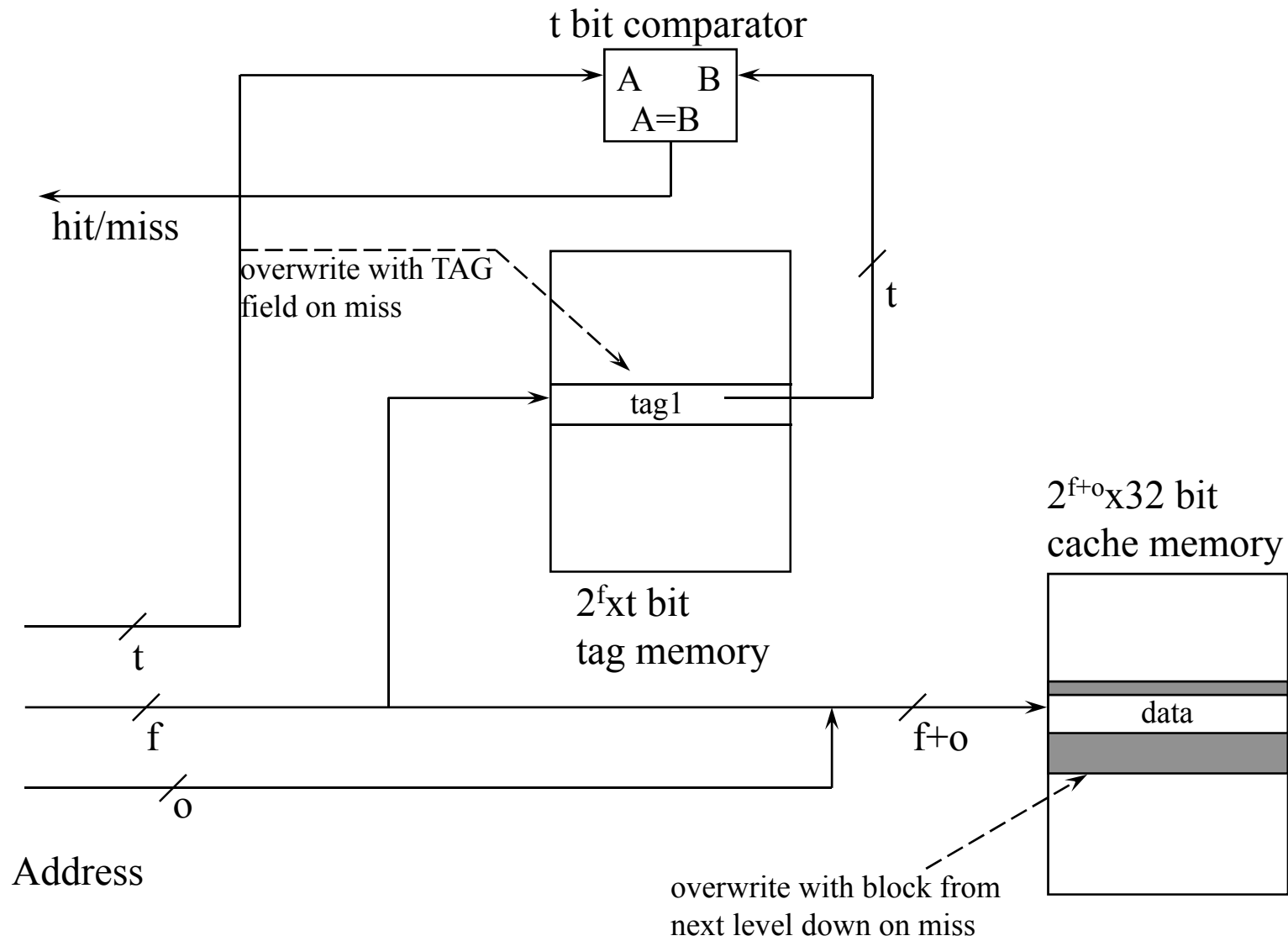
# How does it work

- Address to Mapper organised as follows

| TAG | FRAME | OFFSET |
|-----|-------|--------|

  – OFFSET is offset to the cache block and passes through unchanged

  – FRAME is address of block in cache and passes through unchanged

  – TAG is used to identify whether the requested block is in the cache

t bit comparator

A    B

A=B

hit/miss

overwrite with TAG
field on miss

t

tag1

$2^{f+o}$x32 bit
cache memory

$2^f$xt bit
tag memory

data

t

f

f+o

o

Address

overwrite with block from
next level down on miss

# Advantages/Disadvantages

- Low hardware overhead
- Fast on hit
  - read doesn't add any time
  - On write overall access time is tag memory time + cache memory time
- Block from next level down can only be mapped into a *specific* block

# Fully Associative

- Block from next level down can be mapped into any block in the cache
  - Similar to Primary/Secondary memory interface

# Advantages/Disadvantages

- Very Flexible

- Does not suffer from problem due to blocks which cannot co-exist in the cache

- Needs full TLB to find block in cache

- Access time always incurs the overhead of accessing the TLB

- Big hardware overhead

- Needs memory management (RANDOM ?)

# Set Associative

---

- Combines advantages of both Direct Mapped and Fully Associative

  – Hardware is relatively simple

  – A Practical cache can be constructed which rivals the Fully Associative Cache in performance terms
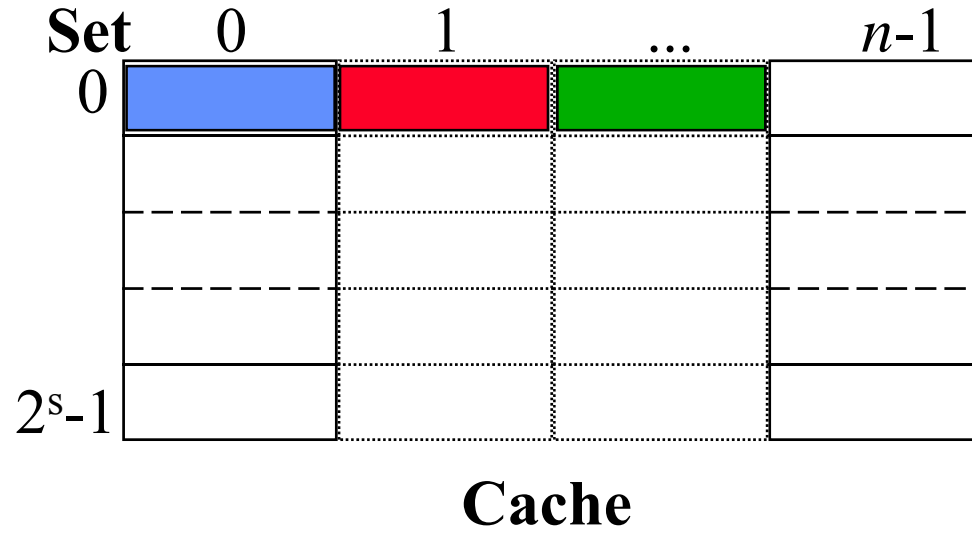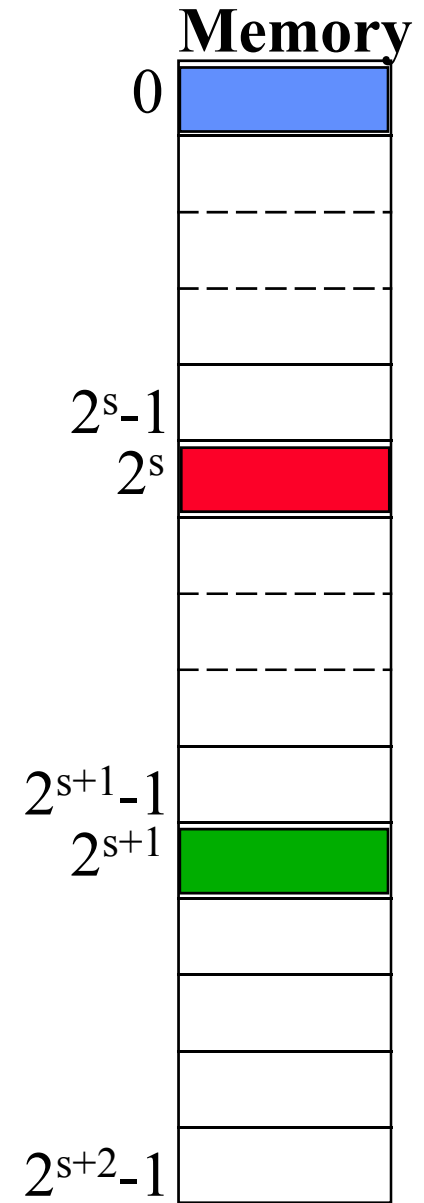
# How does it Work

- Address is as follows

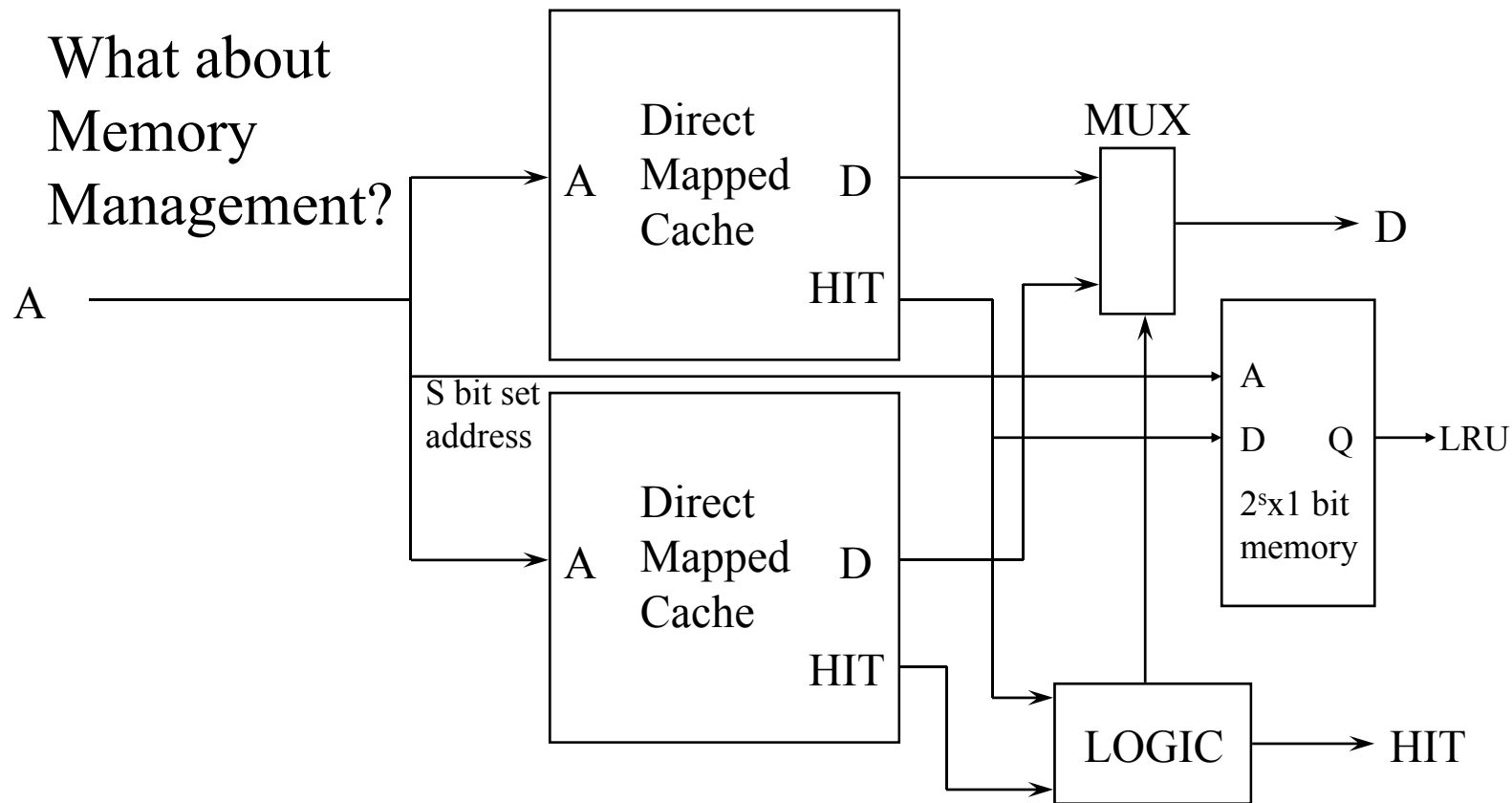| TAG | SET | OFFSET |
|-----|-----|--------|

- Cache is organised into sets
- Each set contains a number of individual blocks or frames - say *n*
- Cache is Direct Mapped with respect to sets (using SET and TAG fields)
- Cache is Fully Associative within a set

Set    0        1      ...    $n$-1

0

$2^s$-1

**Cache**

Memory

0

$2^s$-1
$2^s$

$2^{s+1}$-1
$2^{s+1}$

$2^{s+2}$-1

**n-way Set Associative Cache**

EEE336/NJP/L13

**18**

# How is it Implemented?

What about Memory Management?

A

Direct Mapped Cache

A    D

HIT

MUX

D

S bit set address

Direct Mapped Cache

A    D

HIT

A

D    Q    → LRU

$2^s$x1 bit memory

LOGIC    → HIT

**2-way Set Associative Cache**

# **Complexity w.r.t. *n***

---

- As n becomes bigger, the cache looks more like a Fully Associative Cache

- 2-way set associative is easy just 2 Direct Mapped Caches with a bit of extra hardware to select between them

- What about memory management?
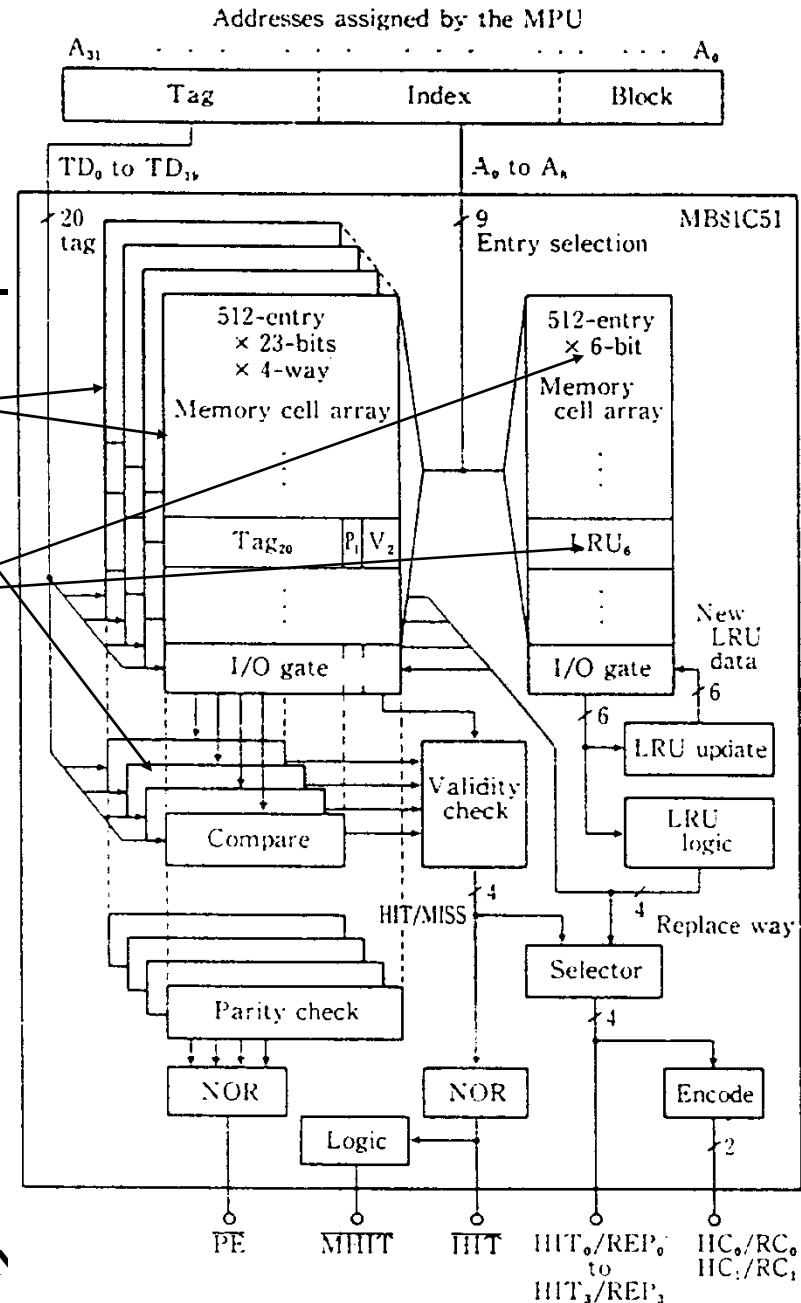
# • Can you explain this?

4-way SA (4 sets of LUTs)

Only one (or none) will HIT

LRU on each set

Why 6 bits / set

- 4 in each set
- 2 bit ID
- A Queue of 3 IDs is all that is necessary = 6 bits / set

Addresses assigned by the MPU

$A_{31}$ . . . . . . . . . . . . . . . $A_0$

| Tag | Index | Block |

$TD_0$ to $TD_{19}$   $A_0$ to $A_8$

20 tag   9 Entry selection   MBS1C51

512-entry × 23-bits × 4-way Memory cell array

512-entry × 6-bit Memory cell array

$Tag_{20}$  $P_1$ $V_2$   $LRU_6$

I/O gate   I/O gate

New LRU data

LRU update

LRU logic

Validity check

Compare

HIT/MISS

Replace way

Parity check

Selector

NOR   NOR   Encode

Logic

PE   MHIT   HIT   HIT$_0$/REP$_0$ to HIT$_3$/REP$_3$   HC$_0$/RC$_0$ HC$_1$/RC$_1$
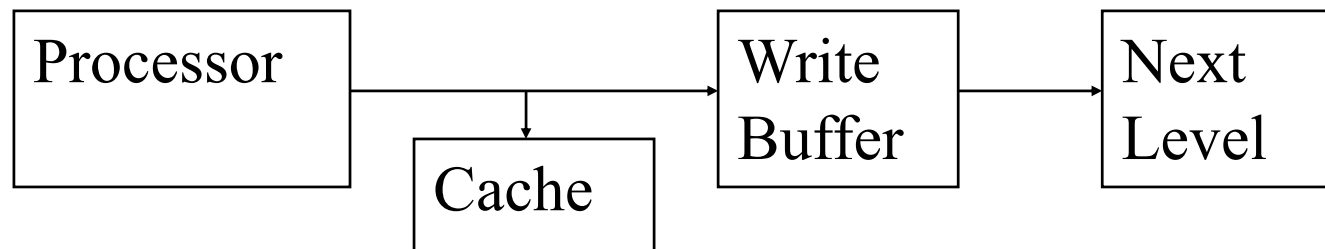
# Other Issues

- Can separate caches for program and data
  - Program cache is never written
  - Caches can be optimised for underlying behaviour of data
  - Disadvantage - decision about balance between program and data is predetermined
- Got to be able to invalidate cache contents
- Some memory must be *uncacheable*

# Writing to a Cache

- Two options on writing to data cache:
    - Write back - When a block is discarded write it back to the next level down
    - Write through - When a datum is written write it to the next level down as well (but don't wait for the write to complete

| Processor | | Cache | Write Buffer | Next Level |

# Write-Through

- Means that blocks do not need to be written back but:

  - Can cause the processor to wait if a write is required before the previous write is finished

  - Two further options with write-through when there is a miss on a write

    - bring the block into the cache - With Allocate
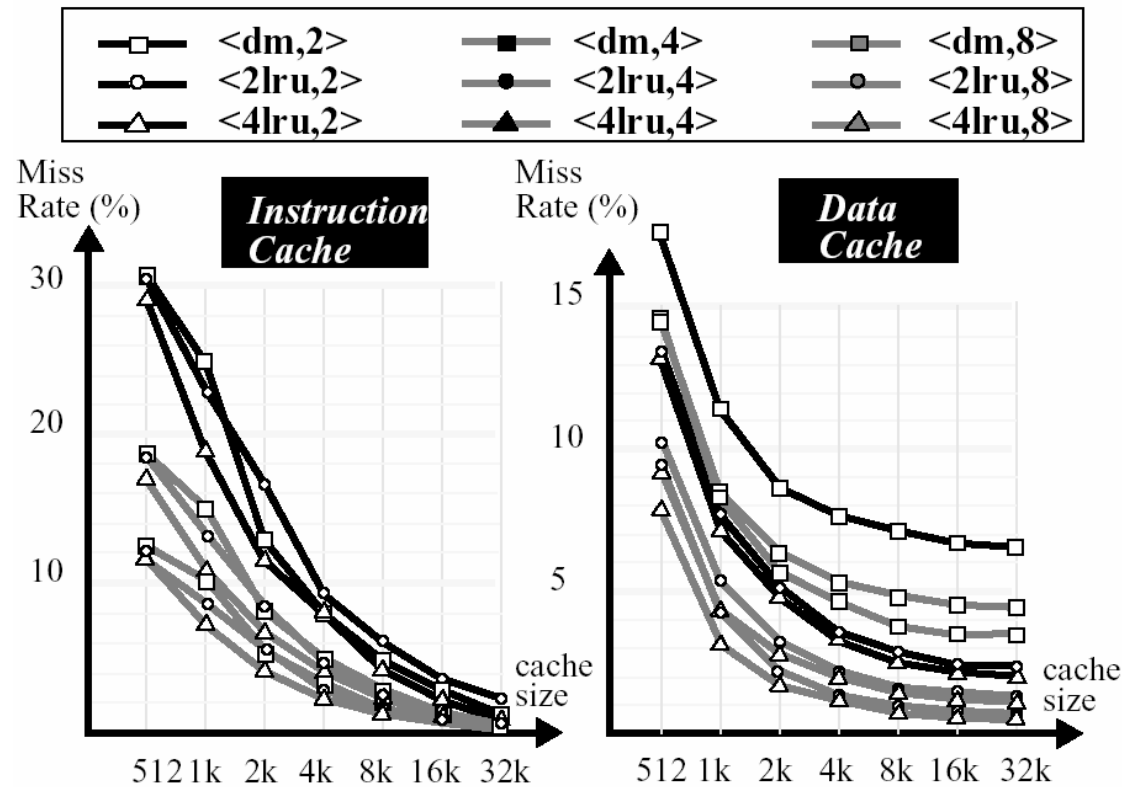    - don't bring the block - With No Allocate

# Cache Performance

---

- There are analytical models but, generally, performance must be derived empirically for representative work-loads

- Performance depends on work-load, size, type, line-size, etc.

# Case Study (L1 Cache)



C-L Su, A M Despain, Cache Design Trade-Offs for Power and Performance Optimisation: A Case Study. Proc. ISLPED, pp63-68, 1995

# Execution Trace Cache

- Intended for instructions
- They do not store an image of the memory
- They store the sequence of instructions that have been executed
- This implies a greater efficiency
- Instructions can be partly decoded before being put in the TEC

# Example

- ## Code snippet

  ```
  L1    move x,y
        cmp y,#0
        je L2
        inst 1
        inst 2
        inst 3
  L2    cmp y,#9
        add x,#1
        jlt L3
        inst 4
        inst 5
  L3    in x,23
        jmp L1
  ```

## Normal cache contents

| move x,y | cmp y,#0 | je L2    | inst 1   |
|----------|----------|----------|----------|
| inst 2   | inst 3   | cmp y,#9 | add x,#1 |
| jlt L3   | inst 4   | in x,23  | jmp L1   |

## Execution trace cache contents

| move x,y | cmp y,#0 | je L2   | cmp y,#9 |
|----------|----------|---------|----------|
| add x,#1 | jlt L3   | in x,23 | jmp L1   |