# Chapter 2
# (Starting with Appendix B)

Instructor: Josep Torrellas

CS433

# Memory Hierarchy Design
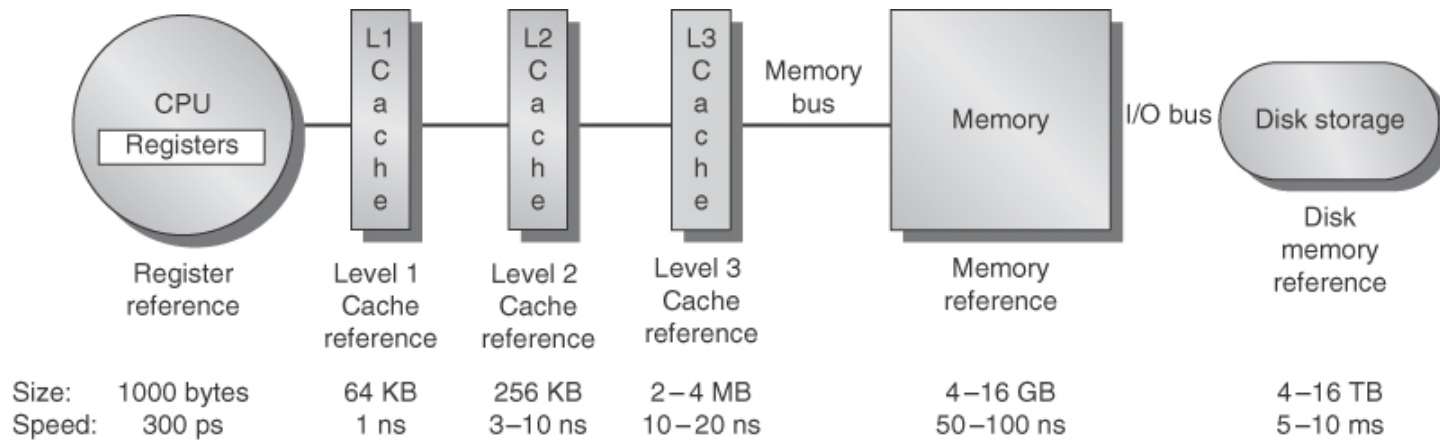
1980 microprocessor → no cache

1995 → 2 level of caches
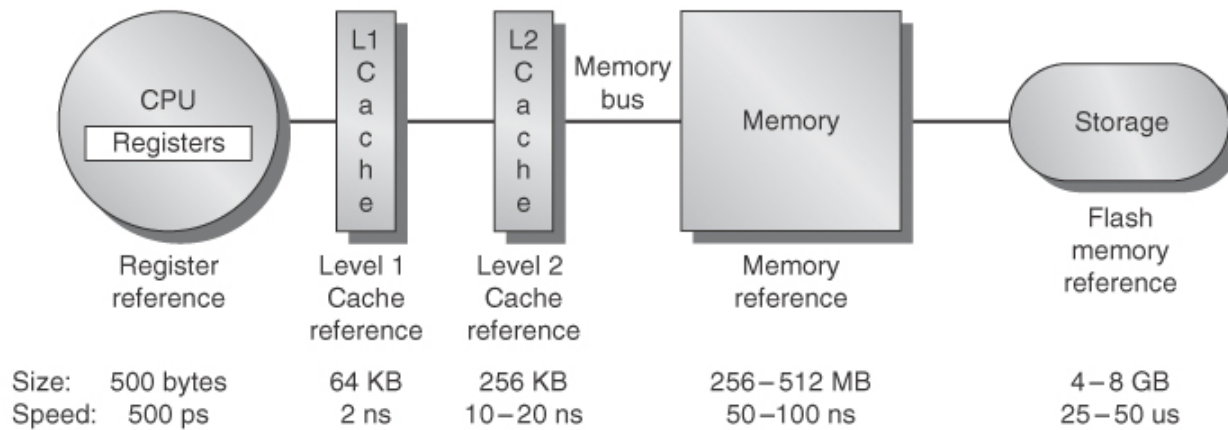
<u>Why ?</u>

- Programmers want more memory
- Principle of locality
- Processor - memory performance gap
- Smaller hardware is faster
- faster ⇒ more expensive

# Memory Hierarchy

CPU
Registers

L1 Cache

L2 Cache

L3 Cache

Memory bus

Memory

I/O bus

Disk storage

Register reference

Level 1 Cache reference

Level 2 Cache reference

Level 3 Cache reference

Memory reference

Disk memory reference

| | | | | | |
|---|---|---|---|---|---|
| Size: | 1000 bytes | 64 KB | 256 KB | 2−4 MB | 4−16 GB | 4−16 TB |
| Speed: | 300 ps | 1 ns | 3−10 ns | 10−20 ns | 50−100 ns | 5−10 ms |

(a) Memory hierarchy for server

CPU
Registers

L1 Cache

L2 Cache

Memory bus

Memory

Storage

Register reference

Level 1 Cache reference

Level 2 Cache reference

Memory reference

Flash memory reference

| | | | | | |
|---|---|---|---|---|---|
| Size: | 500 bytes | 64 KB | 256 KB | 256−512 MB | 4−8 GB |
| Speed: | 500 ps | 2 ns | 10−20 ns | 50−100 ns | 25−50 us |

(b) Memory hierarchy for a personal mobile device

- Figure 2.1

3

# Cache Blocks (Lines)

Block $\rightarrow$ minimum unit of information that can be present in the cache
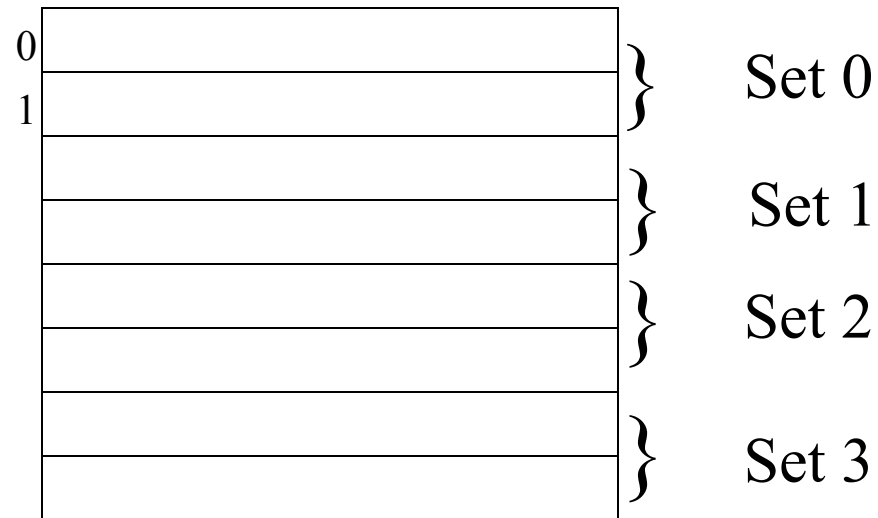
Hit / Miss

Block placement

Block identification

Block replacement

Write strategy

# Block Placement

Block #



0 — 1 — } Set 0

} Set 1

} Set 2

} Set 3

Block address in memory : 10
Direct mapped : (Block addr) Mod (# blocks in cache)
Block # : 10 MOD 8  = 2
Set Associative : (Block Addr) Mod (# sets in cache)
Set # : 10 MOD 4  =  2
Fully Associative : Anywhere

# Block Placement



Fully associative:
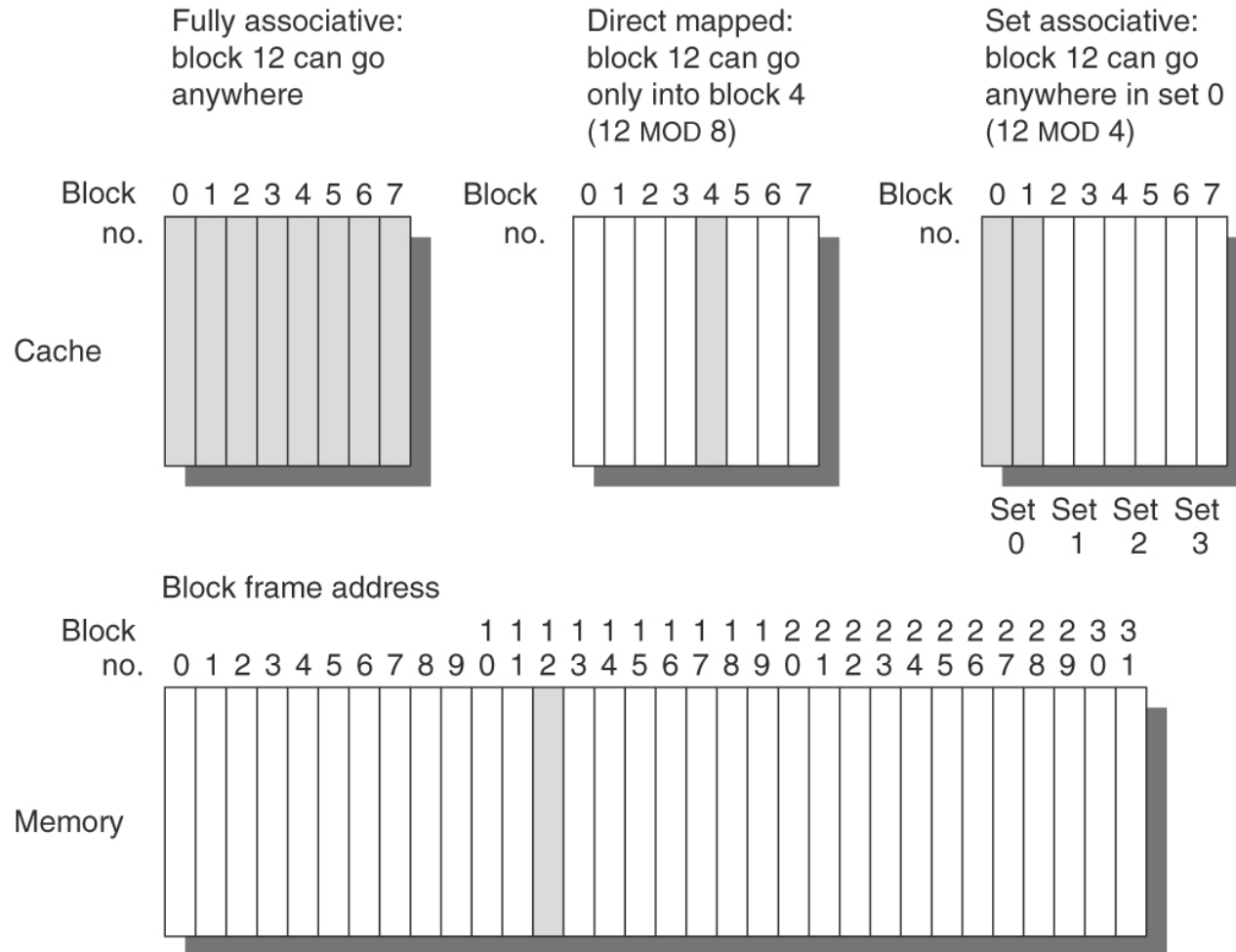block 12 can go
anywhere

Direct mapped:
block 12 can go
only into block 4
(12 MOD 8)

Set associative:
block 12 can go
anywhere in set 0
(12 MOD 4)

Figure B.2

# Block Identification

| Block Address | | Block |
|---|---|---|
| Tag | Index | Offset |

- Index and block offset not stored

- Tag

- Valid bit

- Associativity ↑ ⟹ Size of tag ↑

# Block Replacement

- Random
- Least recently used (LRU)

# Write Policy

Writes form  25% of data cache traffic

　　　　　　 7% of total memory traffic

- Make common case fast $\Rightarrow$ optimize for reads
- Also processor waits for reads to complete; need not wait for writes
- Easy to make reads fast :

  read data in parallel with reading and comparing tag
- Writes : Block cannot be modified until tag matches
- Size of write

# Two Basic Write Policies

- Write through
- Write back
- Dirty bit saves writes
- Write back $\Rightarrow$ less memory bandwidth
- Write through $\Rightarrow$ read miss never results in write to lower level. Easier to implement
- Write stall during write through
  - use a write buffer

# Write Miss Policies

On a write miss :

- Write allocate (fetch on write)

- No-write allocate (write around)


- Write miss policy and write policy are independent

- Write back caches generally use write allocate

- Write through caches often use no write allocate

# Cache Performance

Avg Memory Access Time = Hit time + Miss rate * Miss Penalty

Example : Compare a 32-KB unified cache with a 16KB I-cache and 16 KB D-cache.

Given : Miss rates are 1.99% , 0.64% and 6.47% respectively . 75% of memory accesses are instruction references.

Hit     1 clock cycle

Miss    50 clock cycles

Load/Store access in unified cache :  adds 1 extra clock cycle due to structural hazard

For the split cache,overall miss rate is

(75% * 0.64%) + (25% * 6.47%) = 2.10%

Average memory access time split

= 75% (1+ 0.64% * 50) + 25% * (1+6.47% *50)

= 2.05

Avg memory access time unified

= 75% * (1+1.99% * 50) + 25% *(1+**1**+1.99% * 50)

= 2.24

$$\text{CPU time} = \left[ \text{CPU Execution Clock cycles} + \text{Memory stall Clock cycles} \right] * \text{clock cycle time}$$

Assumption : All memory stalls are due to cache misses.

$$\left[ \text{Memory stall clock cycles} \right] = \text{Reads * Read miss rate * Read miss Penalty}$$
$$+ \text{Writes + Write miss Rate * Write miss Penalty}$$
$$= \text{Memory accesses * Miss rate * Miss Penalty}$$

Consider a machine with :

Cache miss penalty = 50 cycles

$CPI_{execution} = 2.0$

Miss Rate = 2%

Memory References per inst = 1.33

$$CPU\ time = IC * ( CPI_{execution} + \frac{Memory\ stall\ clock\ cycles}{Instruction} ) * Clock\ cycle\ time$$

$$= IC * ( 2.0 + (1.33 * 2\% * 50 ) ) * Clock\ cycle\ time$$
$$= IC * 3.33 * Clock\ Cycle\ time$$

Impact of cache:   1.67x

If no cache:         68.5   30x

If perfect CPIexecution: 2.33x

Cache misses have double - barreled impact on a CPU with
   low  CPI and fast clock

1.  Lower CPI $_{\text{execution}}$ $\rightarrow$ relative impact
     of FIXED number of cache miss cycles increases

2.  Identical memory hierarchy $\rightarrow$ CPU with faster clock sees
                                     more stall cycles

Note: in ooo execution processors: part of the memory access
   latency is overlapped with computation!