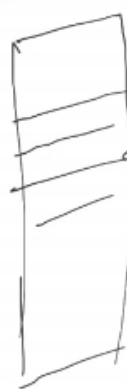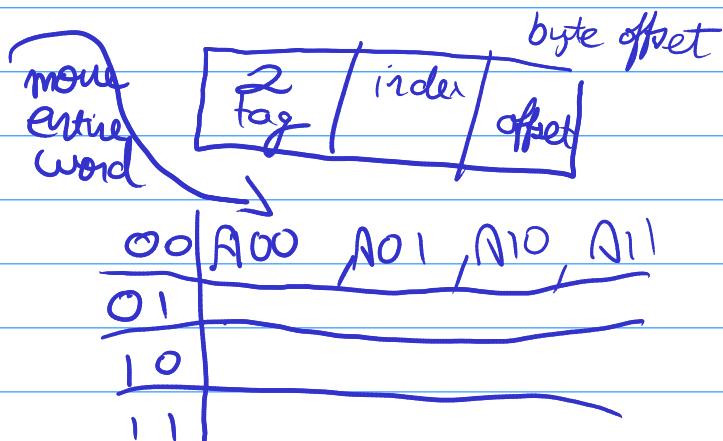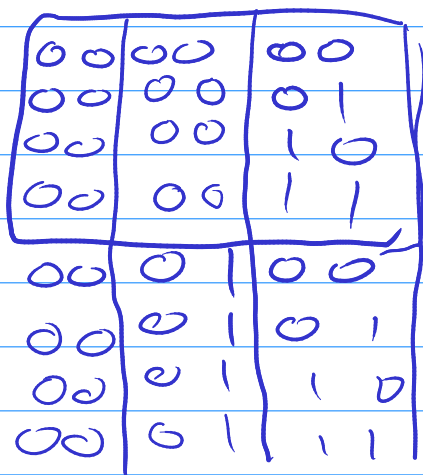Line
=
Block

we bring data from physical memory to the cache memory in blocks rather than as bytes.

if 4 bytes → word

① case 1 → transfer words from PM to CM

② case 2 → transfer blocks (group of bytes) from phy mem to CM →

move entire word

| 2 Tag | index | byte offset offset |
|---|---|---|

| | | |
|---|---|---|
| 00 | A00 | A01, A10, A11 |
| 01 | | |
| 10 | | |
| 11 | | |

generalized location

| k-n-2 | n | 2 | for k bits |
|---|---|---|---|

Case 2    Cache line → block of memory location in PM that is in CM

4 words 1 block

| k-n-4 | n | 2 | 2 |
|---|---|---|---|

index     word offset

→ byte offset

$Line = Block$

↓ Set associative

↓ direct mapped Cache

| tag | word1 | word2 | word3 | word |

| b0 | b1 | b2 | b3 |

# Cache and Main Memory Structure



Figure 4.4 Cache/Main-Memory Structure

$PM - M$ blocks

1 - block - k words

1 word - 4 bytes

$2^n$ addressable words

$2^n/K$ blocks

Cache contains C lines or slots       $C << M$

Cache   Line 0  ←  Block 0    eg 4k words in each block

Line 1  ←  Block 1

...

Line C-1    Block$(2^n/k - 1)$

# Direct Mapping

| Tag s-r | Line or Slot r | Word w |
|---------|----------------|--------|
| 8 | 14 | 2 |

*(handwritten: → byte offset, index)*

*(handwritten: Set associative mapping)*

- 24 bit address *PM*
- 2 bit word identifier (4 byte block)
- 22 bit block identifier
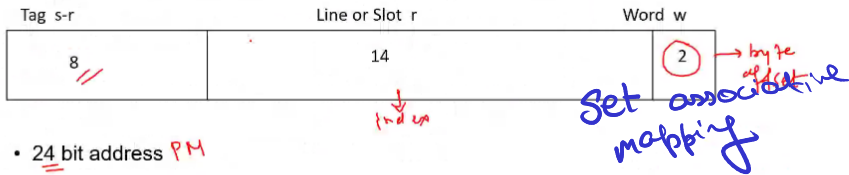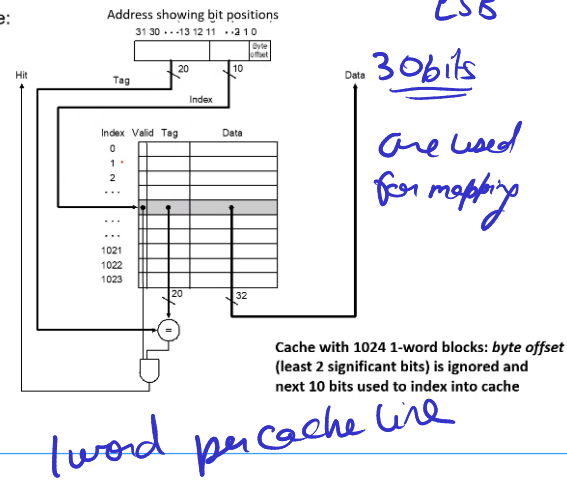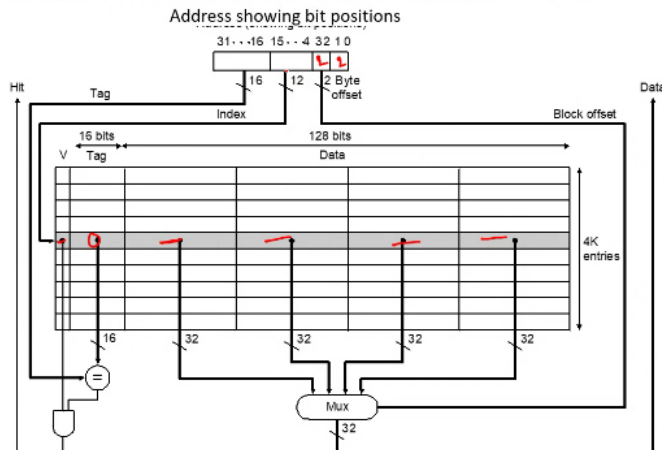  - 8 bit tag (=22-14)
  - 14 bit slot or line
- No two blocks in the same line have the same Tag field
- Check contents of cache by finding line and checking Tag

*(handwritten: different tags cannot map into some place)*

- MIPS style:

Address showing bit positions
31 30 ··· 13 12 11 ··· 2 1 0



**Cache with 1024 1-word blocks:** *byte offset* (least 2 significant bits) is ignored and next 10 bits used to index into cache

*(handwritten top: 32 bits, 2 bits are byte offset LSB)*
*(handwritten: 30 bits are used for mapping)*
*(handwritten: 1 word per cache line)*

# Direct Mapped Cache: Taking Advantage of Spatial Locality

- Taking advantage of spatial locality with *larger* blocks:



**Cache with 4K 4-word blocks:** *byte offset* (least 2 significant bits) is ignored, next 2 bits are *block offset*, and the next 12 bits are used to index into cache

*(handwritten: Spatial locality may cause dependence b/w words in block ⇒)*
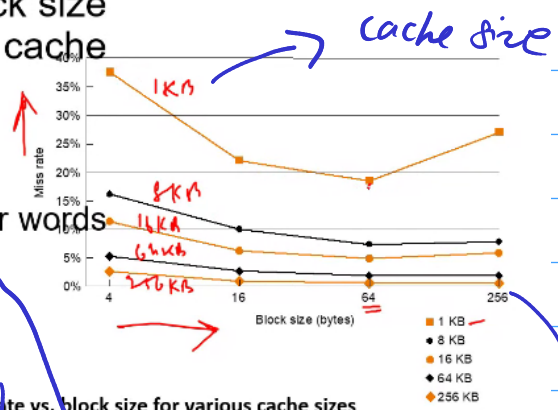
# Direct Mapped Cache: Taking Advantage of Spatial Locality

- *Cache replacement* in large (multiword) blocks:
  - word *read miss*: read entire block from main memory
  - word *write miss*: *cannot* simply write word and tag! *Why*?!
  - writing in a *write-through* cache:
    - if *write hit*, i.e., tag of requested address and and cache entry are equal, continue as for 1-word blocks by replacing word and writing block to both cache and memory
    - if *write miss*, i.e., tags are unequal, fetch block from memory, replace word that caused miss, and write block to both cache and memory
    - therefore, unlike case of 1-word blocks, a write miss with a multiword block causes a memory read

- Miss rate falls at first with increasing block size as expected, but, as block size becomes a large fraction of total cache size, miss rate may go up because
  - there are few blocks
  - competition for blocks increases
  - blocks get ejected before most of their words are accessed (*thrashing* in cache)

Too many block replacements due to temporal locality



**Miss rate vs. block size for various cache sizes**

multiprogrammings

high cache ↓
⟹ power consumption ↑
cost ↑
space is small

128 kB cache data

1 w block size

32 bits space

total bits?

$32k$ words $= 2^{15}$

entry data $\quad$ tag $\quad$ valid

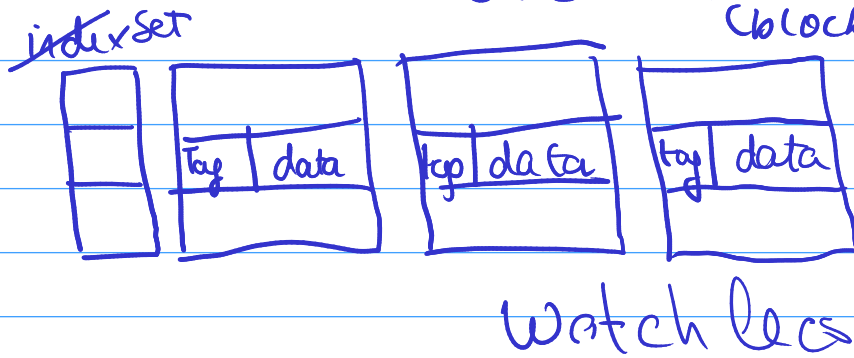$32 + (32 - 15 - 2) + 1 = 48 \text{ bits}$

$48 \times 2^{15} = (1.5 \times 32) \times 2^{15}$

$= 1.5 \times 2^{20} \text{ bits}$

data bits = 1 Mbits

total cache size / actual cache data = 1.5

# Set associative mapping

In set associative mapped cache, same indices in the cache can store > 1 line (block)

index/set

| | | Tag | data | | Tag | data | | tag | data |

watch lecs

main memory — 128 bytes. locations are byte addressable

CM — 32 bytes
BS — 4 bytes
Set — 2way set associative

$$\frac{128}{4} = 32$$

| tag | data | | tag | data |

Lines: $\frac{32}{4} = 8$ lines

| 4 | BS with | 4 |

no. of sets = $\frac{\text{no. of lines}}{2 \text{ no. of line/set}} = 4$

BS → size: 2 bits for Block offset

| tag | set no. | bytes block offset |
|-----|---------|--------------------|
| 3 | 2 | 2 |

of index

Now, to which # (associated) in same set do we place

Step 1   find the set no.
Step 2   Compare the tag (valid bit first)

n way   set associative
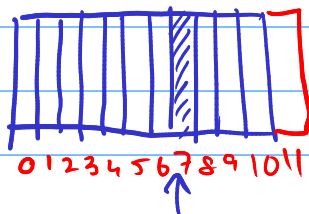        = n comparators

Each comparator has no. of bits in tag

# Decreasing Miss Rates with Associative Block Placement

1 comparators

- *Direct mapped*: one *unique* cache location for each memory block
  - cache block address = memory block address *mod* cache size

- *Fully associative*: each memory block can locate *anywhere* in cache
  - *all* cache entries are searched (in parallel) to locate block    → comparator = all locations (no. of)

n-comparators

- *Set associative*: each memory block can place in a *unique set* of cache locations
  – if the set is of size n it is n-way set-associative
  - cache set address = memory block address *mod*  number of sets in cache
  - all cache entries in the corresponding set are searched (*in parallel*) to locate block

- Increasing degree of associativity
  - *reduces miss rate*
  - *increases hit time* because of the parallel search and then fetch

0 1 2 3 4 5 6 7 8 9 10 11

direct

fully set associative
0   1   2   3

→ n-way associative (Compromise)
it's like a small fully associative cache

# Decreasing Miss Rates with Associative Block Placement

One-way set associative
(direct mapped)

Block   Tag   Data
0
1
2
3
4
5
6
7

Two-way set associative

Set   Tag   Data   Tag   Data
0
1
2
3

Four-way set associative

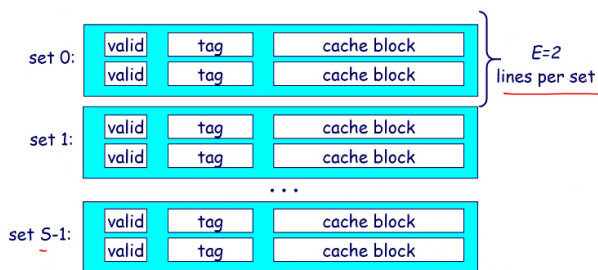Set   Tag   Data   Tag   Data   Tag   Data   Tag   Data
0
1

Eight-way set associative (fully associative)

Tag  Data  Tag  Data  Tag  Data  Tag  Data  Tag  Data  Tag  Data  Tag  Data  Tag  Data

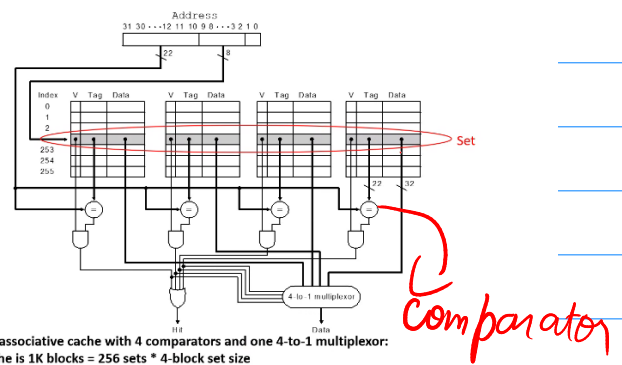**Configurations of an 8-block cache with different degrees of associativity**

## Example: Set Associative Cache
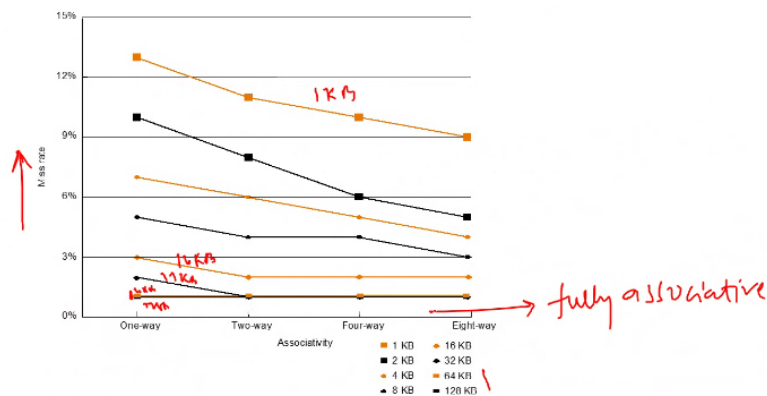
**Characterized by more than one line per set**

set 0:
| valid | tag | cache block |
| valid | tag | cache block |

E=2
lines per set

set 1:
| valid | tag | cache block |
| valid | tag | cache block |

. . .

set S-1:
| valid | tag | cache block |
| valid | tag | cache block |

E-way associative cache

## Implementation of a Set-Associative Cache

Address
31 30 · · · 12 11 10 9 8 · · · 3 2 1 0

22          8

Index  V  Tag  Data        V  Tag  Data        V  Tag  Data        V  Tag  Data
0
1
2

253
254
255                                                                          Set

                                                            22      32

                                    4-to-1 multiplexor

Hit                                          Data

*Comparator*

4-way set-associative cache with 4 comparators and one 4-to-1 multiplexor:
size of cache is 1K blocks = 256 sets * 4-block set size

# Performance with Set-Associative Caches



15%

12%

1K M

9%

Miss rate

6%

3%            16 KB
              32 KB
              64 KB        fully associative
0%
       One-way   Two-way   Four-way   Eight-way
       Associativity

| ■ 1 KB | ◆ 16 KB |
| ■ 2 KB | ▲ 32 KB |
| ◆ 4 KB | ■ 64 KB |
| ▲ 8 KB | ■ 128 KB |

Miss rates for each of eight cache sizes with increasing associativity:
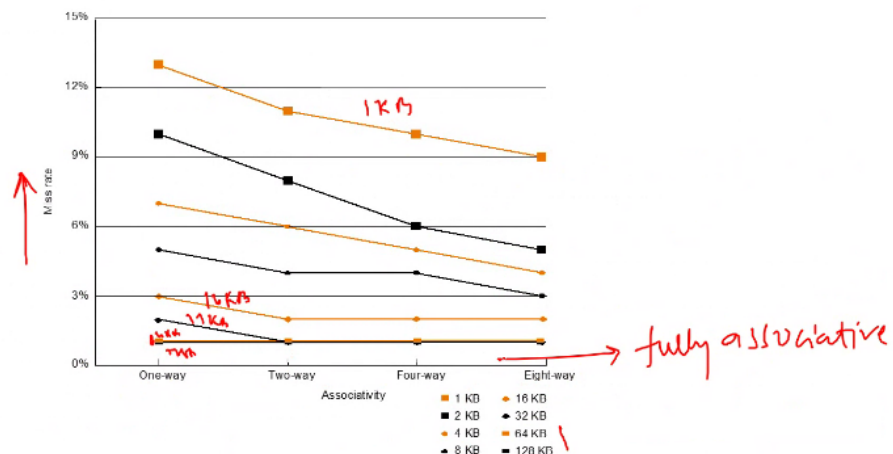data generated from SPEC92 benchmarks with 32 byte block size for all caches

# Performance with Set-Associative Caches



Miss rates for each of eight cache sizes with increasing associativity:
data generated from SPEC92 benchmarks with 32 byte block size for all caches

# Example Problems

- *Find the number of misses for a cache with four 1-word blocks given the following sequence of memory block accesses:*  0, 8, 0, 6, 8,  ← mem. references.
  *for each of the following cache configurations*
  1. direct mapped
  2. 2-way set associative (use LRU replacement policy)  Least recently used
  3. fully associative

- Note about LRU replacement
  - in a 2-way set associative cache LRU replacement can be implemented with one bit at each set whose value indicates the mostly recently referenced block