

# Lecture - 27

## Associative Caches

\* Problem with direct<sup>①</sup> mapped caches

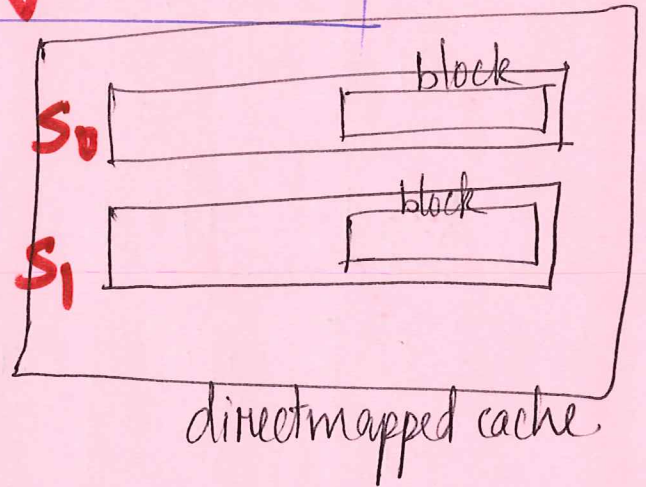
"Conflict misses"

a, b - int arrays of size 4.

```
for (i=0; i < 8; i++)
    sum += x[i] * y[i];
```

Cache block size = 8 bytes

Cache size = 16 bytes



Elem.	Addr	Set index
x[0]	0	0
x[1]	4	0
x[2]	8	1
x[3]	12	1
y[0]	16	0
y[1]	20	0
y[2]	24	1
y[3]	28	1

block 0

x[4]

block 1

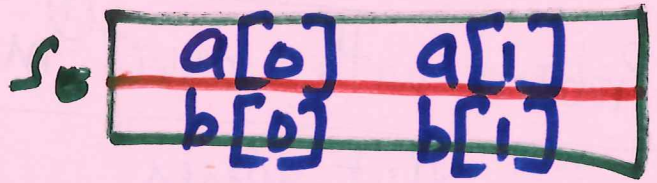
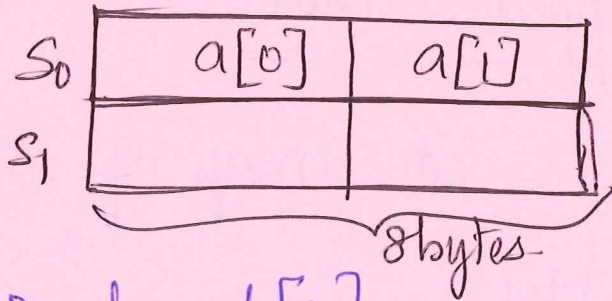
x[6]

block 2

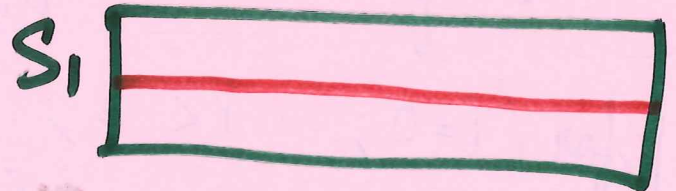
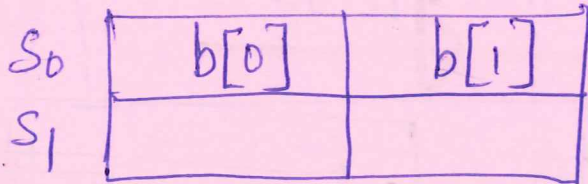
block 3

(2)

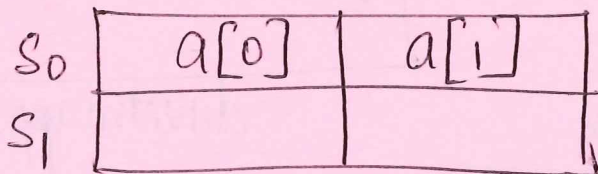
1. Read  $a[0]$  .



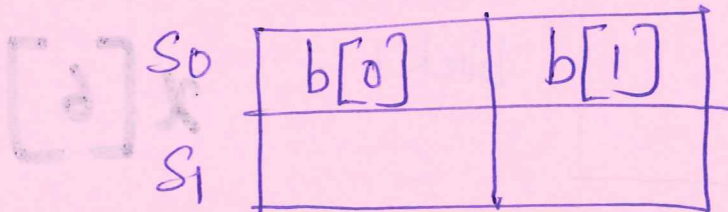
2. Read  $b[0]$



3. Read  $a[1]$



4. Read  $b[1]$



⇒ Trashing!

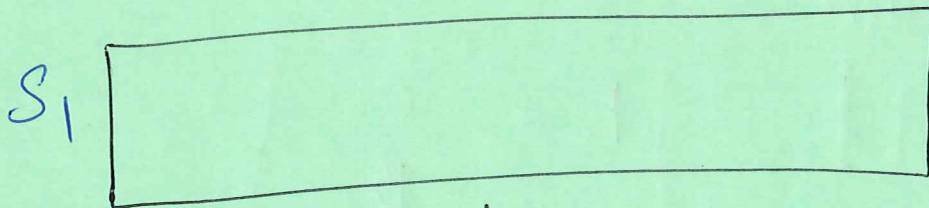
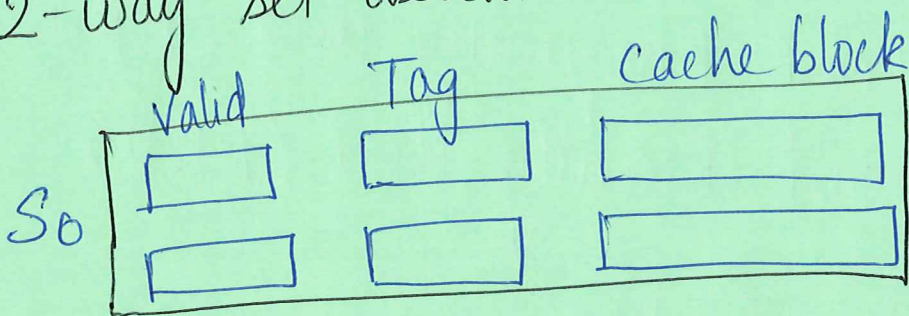
# \* Set Associative <sup>(3)</sup> caches

Each set holds more than one cache line.

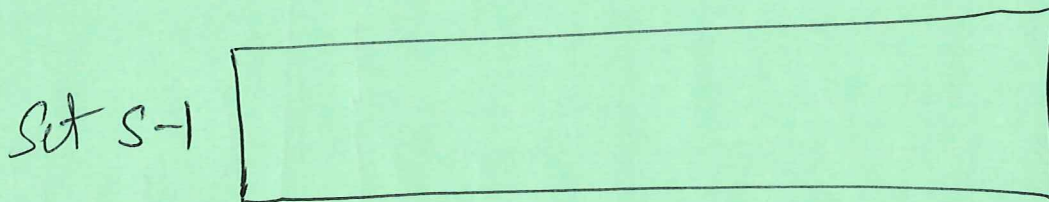
$$1 < E < \frac{C}{B}$$

E-way set associative cache.

eg. 2-way set associative cache.



⋮



1. Set selection - same as direct-mapped cache. (4)
2. Line matching and word selection. (X)
3. Line replacement on Misses.

### Replacement policies

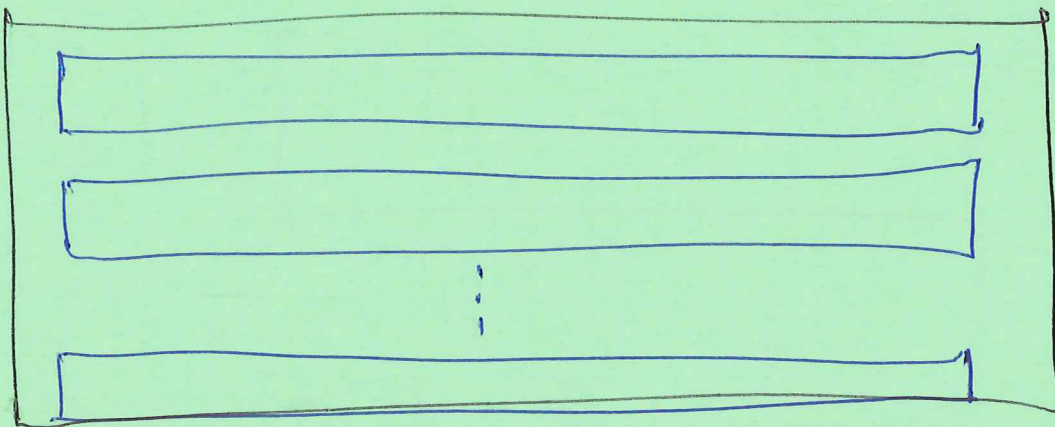
1. Random
2. Least Frequently used (LFU)
3. Least Recently Used (LRU).

### \* Fully Associative Caches

Only 1 Set!

$$E = \frac{C}{B}$$

So

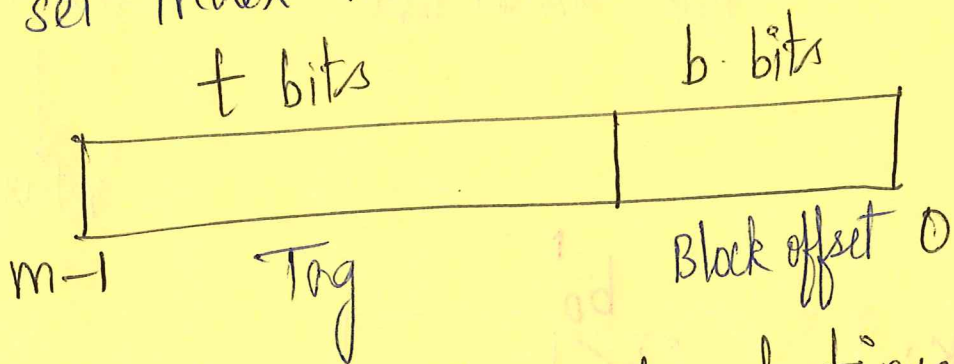


$E = \frac{C}{B}$  lines  
in the one  
and only  
set.

# 1. Set Selection (5)

trivial  $\because$  only one set.

$\therefore$  no set index bits in the address



# 2. Line matching and word selection.

!!! to set associative cache:

Fully associative caches  $\rightarrow$  appropriate for small caches

$\because$  tag matching in a large cache takes a long time

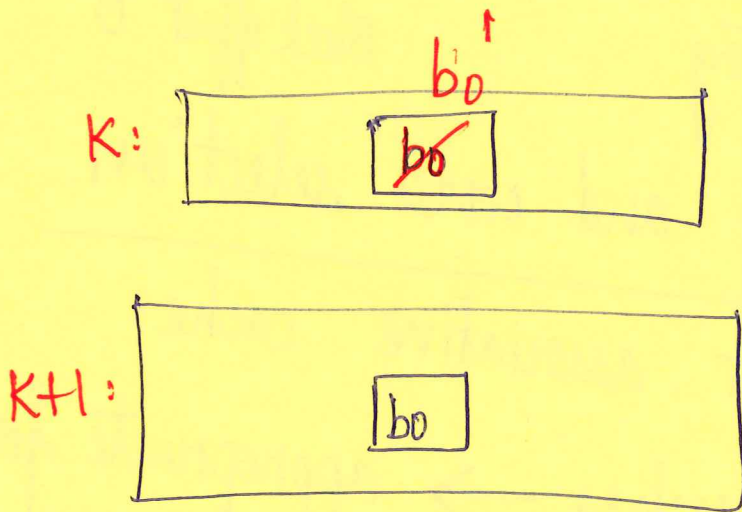
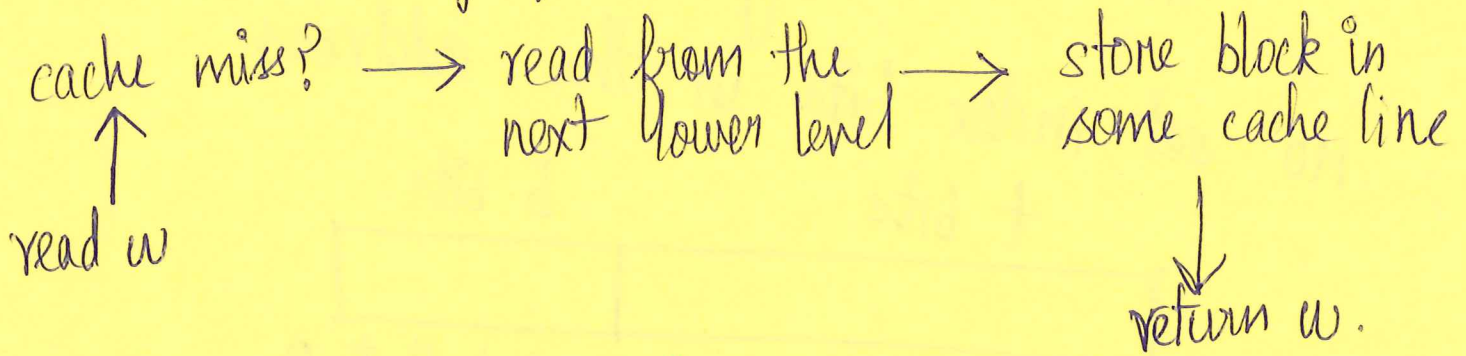
# HW problem 6.30

HW: Why the middle bits are used for set-index?  
Why not the higher order bits?

# Writes

(6)

Reads - straightforward.



## Writes

### Write through

- \* more bus traffic
- \* simple (no dirty bit)
- \* writes block to level K+1 immediately

### Write-back

- \* less bus traffic
- \* dirty bit needed.
- \* writes block to level K+1 only when the block is evicted

7

Writes

no-write allocate

writes directly in main memory if block not present in cache

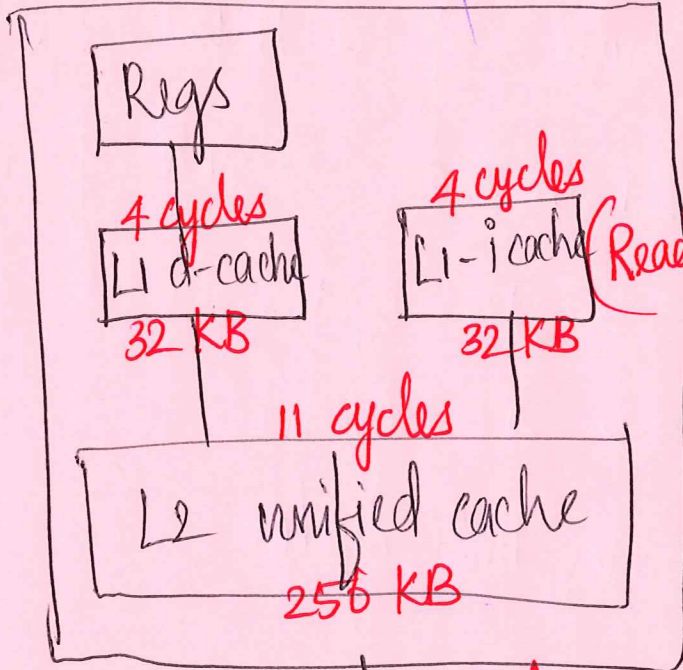
write-allocate loads block in cache and writes in cache

write through cache

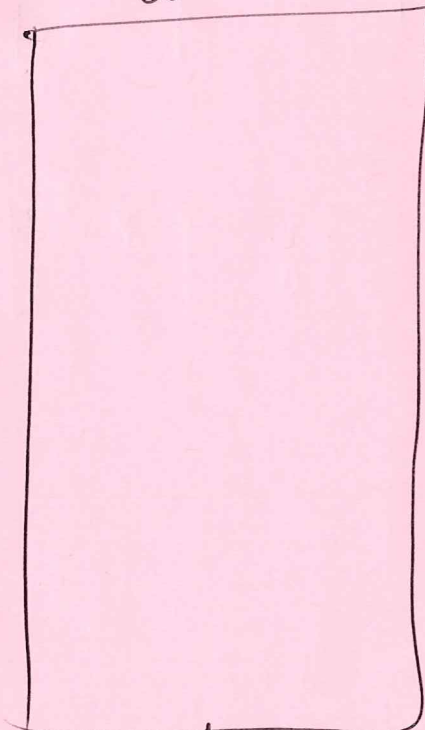
write-back caches

### \* Real Cache Hierarchy

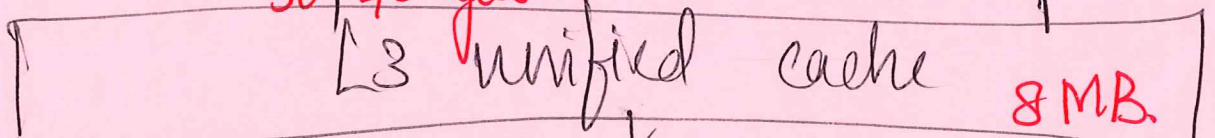
Core 0



Core 3



30-40 cycles



# Performance impact <sup>8</sup> of cache parameters

1. miss rate =  $\frac{\# \text{ misses}}{\# \text{ references}}$

2. hit rate =  $1 - \text{miss rate}$ .

3. hit time = time to deliver a word in the cache to the CPU

=  $t_{\text{set selection}} + t_{\text{line identification}} + t_{\text{word identification}}$

4. Miss penalty = any additional time required because of a miss. (eg. choosing a victim for eviction)

eg. penalty for L1 misses served from:

1. L2  $\rightarrow$  10 cycles

2. L3  $\rightarrow$  40 cycles

3. L4  $\rightarrow$  100 cycles.



## Impact of cache size (9)

large cache  $\Rightarrow$  high hit rate.

"  $\Rightarrow$  high hit time

$\therefore$  Faster caches are usually smaller.

## Impact of block size

larger blocks  $\Rightarrow$  increase hit rate  
(for programs with good spatial locality).

"  $\Rightarrow$  decrease hit rate  
(for programs with good temporal locality)

$\therefore$  For a given cache size,

larger blocks  $\Rightarrow$  fewer cache lines.

larger blocks

$\Rightarrow$  negative impact on  
"miss penalty"

$\therefore$  larger blocks  $\Rightarrow$  larger transfer times.

# Impact of Associativity (E)

larger  $E \Rightarrow$  increase hit rate  
 $\because$  conflict misses are reduced.

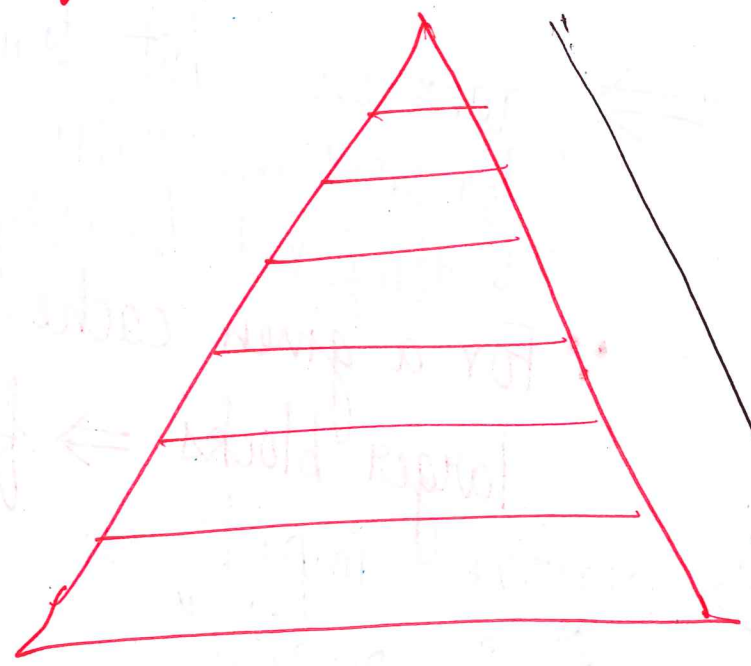
choice of E:  
hit time vs miss penalty.

Intel i7

$L_1, L_2 \rightarrow E = 8$   
 $L_3 \rightarrow E = 16$

miss penalty increases  
~~hit time~~

# Impact of Write Strategy



use of write-back caches increases.

$\because$  transfer time increases.