

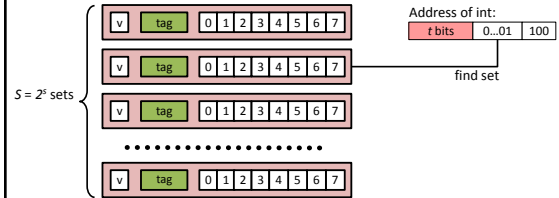
# Lecture 21

## Cache II

CPSC 275  
Introduction to Computer Systems

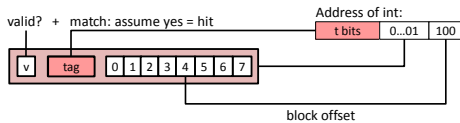
### Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set  
Assume: cache block size 8 bytes



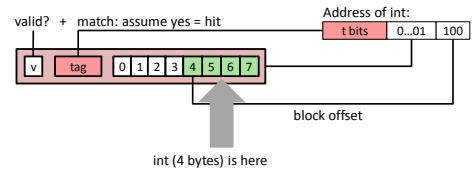
### Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set  
Assume: cache block size 8 bytes



### Example: Direct Mapped Cache (E = 1)

Direct mapped: One line per set  
Assume: cache block size 8 bytes



If no match, an old line is evicted and replaced.

### Direct-Mapped Cache Simulation

t=1	s=2	b=1
x	xx	x

M=16 byte addresses (total size of memory)

B=2 bytes/block

S=4 sets

E=1 block/set

Address trace:

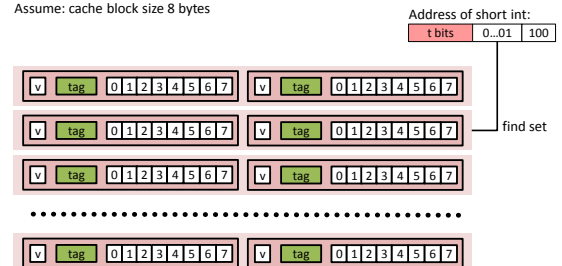
0	[0000] <sub>2</sub>	miss
1	[0001] <sub>2</sub>	hit
7	[0111] <sub>2</sub>	miss
8	[1000] <sub>2</sub>	miss
0	[0000] <sub>2</sub>	miss

	v	Tag	Block
Set 0	1	0	M[0-1]
Set 1			
Set 2			
Set 3	1	0	M[6-7]

### 2-way Set Associative Cache

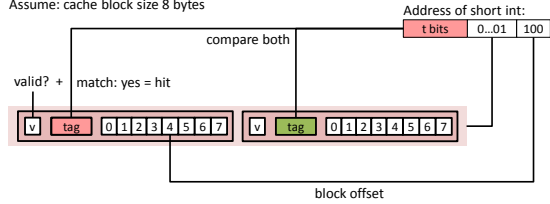
E = 2: Two lines per set

Assume: cache block size 8 bytes



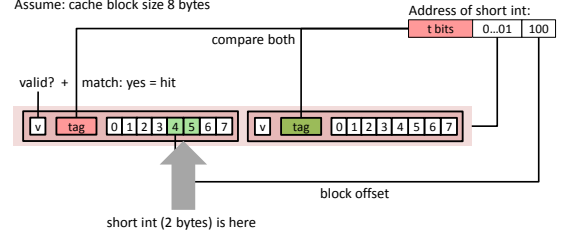
## 2-way Set Associative Cache

E = 2: Two lines per set  
Assume: cache block size 8 bytes



## 2-way Set Associative Cache

E = 2: Two lines per set  
Assume: cache block size 8 bytes



No match:

- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

## 2-Way Set Associative Cache Simulation

t=2	s=1	b=1
xx	x	x

M=16 byte addresses, B=2 bytes/block,  
S=2 sets, E=2 blocks/set

Address trace (reads, one byte per read):

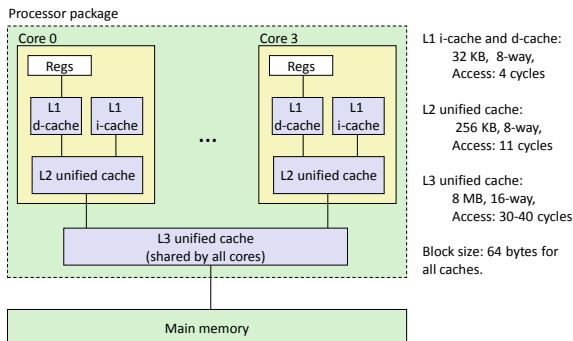
0	[0000 <sub>2</sub> ]	miss
1	[0001 <sub>2</sub> ]	hit
7	[0111 <sub>2</sub> ]	miss
8	[1000 <sub>2</sub> ]	miss
0	[0000 <sub>2</sub> ]	hit

	v	Tag	Block
Set 0	1	00	M[0-1]
	1	10	M[8-9]
Set 1	1	01	M[6-7]
	0		

## What about Writes?

- Multiple copies of data exist:
  - L1, L2, Main Memory, Disk
- What to do on a write-hit?
  - **Write-through** (write immediately to memory)
  - **Write-back** (defer write to memory until replacement of line)
    - Need a *dirty (modify)* bit (line different from memory or not)
- What to do on a write-miss?
  - **Write-allocate** (load into cache, update line in cache)
    - Good if more writes to the location follow
  - **No-write-allocate** (writes immediately to memory)
- Typical
  - Write-through + No-write-allocate
  - Write-back + Write-allocate

## Intel Core i7 Cache Hierarchy



## Cache Performance Metrics

- Miss Rate
  - Fraction of memory references not found in cache (misses / accesses)
  - = 1 – hit rate
  - Typical numbers (in percentages):
    - 3-10% for L1
    - can be quite small (e.g., < 1%) for L2, depending on size, etc.
- Hit Time
  - Time to deliver a line in the cache to the processor
    - includes time to determine whether the line is in the cache
  - Typical numbers:
    - 1-2 clock cycle for L1
    - 5-20 clock cycles for L2
- Miss Penalty
  - Additional time required because of a miss
    - Typically 50-200 cycles for main memory

## Lets think about those numbers

- Huge difference between a hit and a miss
- Would you believe 99% hits is twice as good as 97%?
  - Consider:
    - cache hit time of 1 cycle
    - miss penalty of 100 cycles
  - Average memory access time:
    - 97% hits:  $1 \text{ cycle} + 0.03 * 100 \text{ cycles} = 4 \text{ cycles}$
    - 99% hits:  $1 \text{ cycle} + 0.01 * 100 \text{ cycles} = 2 \text{ cycles}$

## Practice Problems

- Read CSaPP Sec. 6.4 and try the following problems:  
6.11, 6.12, and 6.13