

15-213

"The course that gives CMU its Zip!"

Cache Memories October 6, 2006

Topics

- Generic cache memory organization
- Direct mapped caches
- Set associative caches
- Impact of caches on performance
- The memory mountain

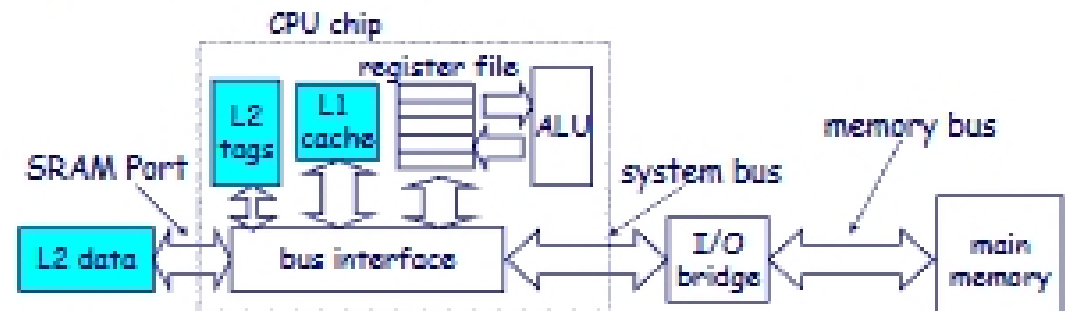
Cache Memories

Cache memories are small, fast SRAM-based memories managed automatically in hardware.

- Hold frequently accessed blocks of main memory

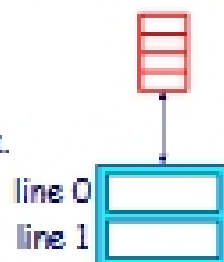
CPU looks first for data in L1, then in L2, then in main memory.

Typical system structure:



Inserting an L1 Cache Between the CPU and Main Memory

The transfer unit between the CPU **register file** and the **cache** is a 4-byte block.



The tiny, very fast CPU **register file** has room for four 4-byte words.

The small fast **L1 cache** has room for two 4-word blocks.

The transfer unit between the **cache** and **main memory** is a 4-word block (16 bytes).



The big slow **main memory** has room for many 4-word blocks.

General Organization of a Cache

Cache is an array of sets.

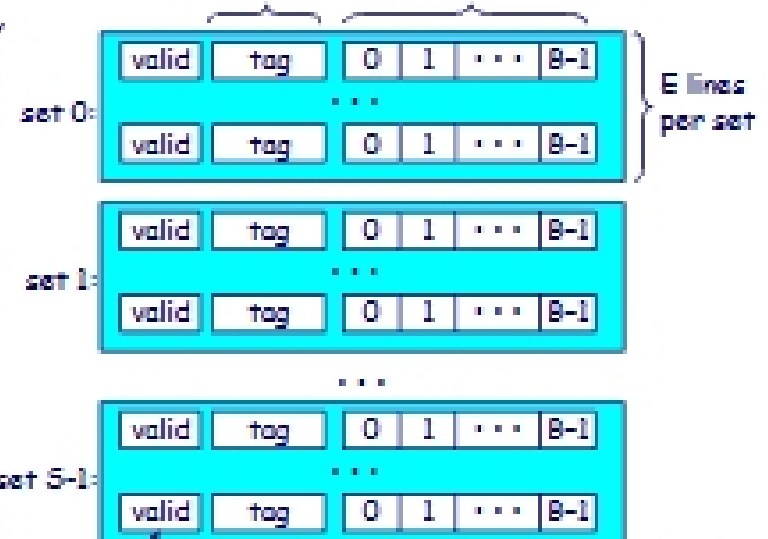
Each set contains one or more lines.

Each line holds a block of data.

$$S = 2^e \text{ sets}$$

t tag bits per line

$B = 2^b$ bytes per cache block

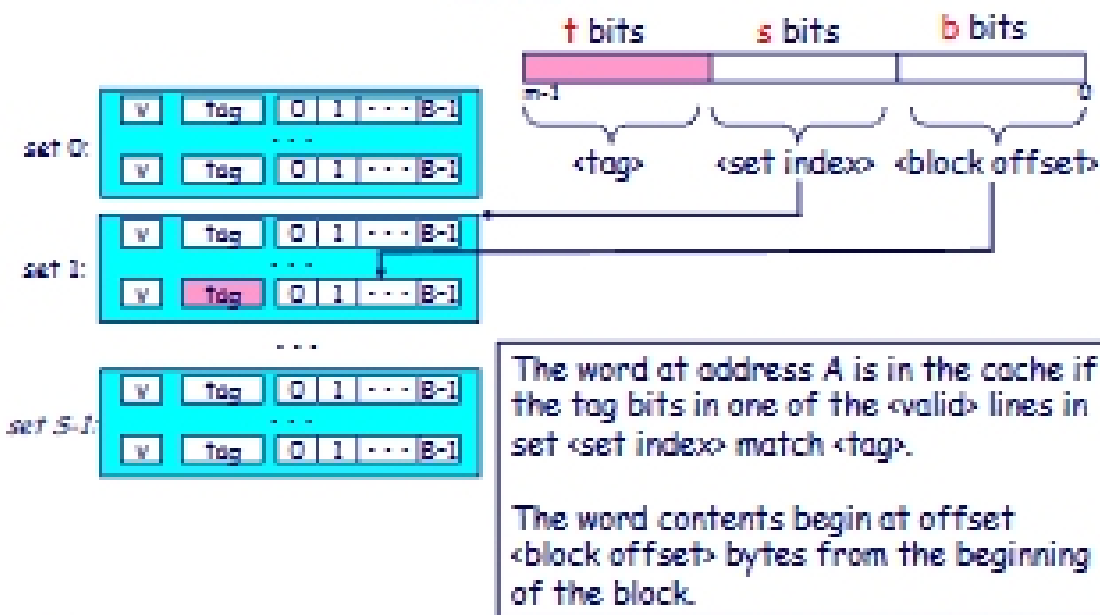


1 valid bit per line

$$\text{Cache size: } C = B \times E \times S \text{ data bytes}$$

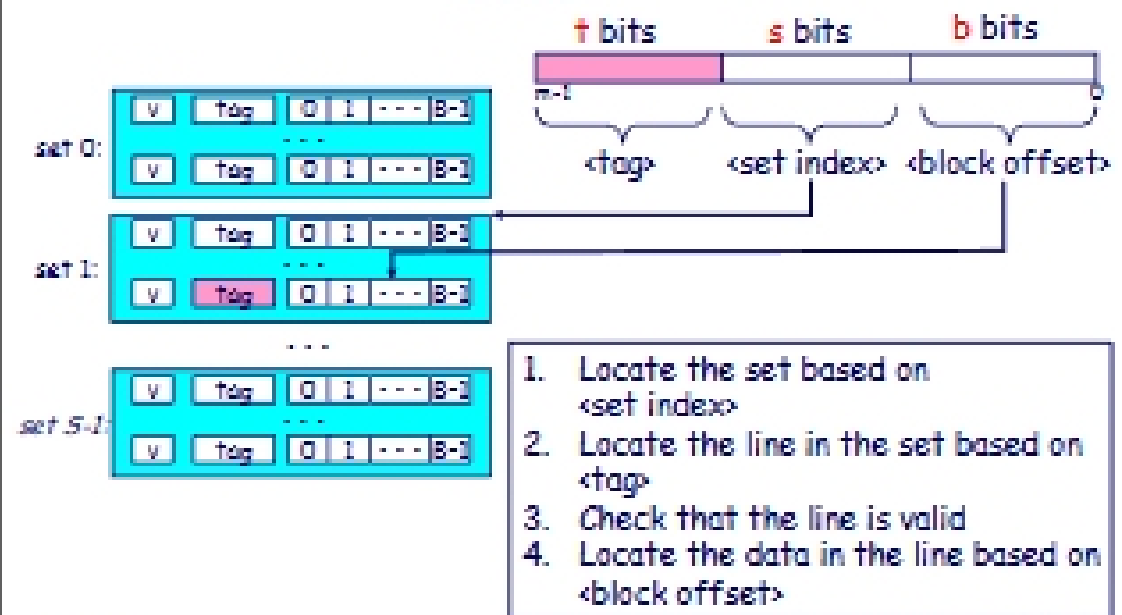
Addressing Caches

Address A:



Addressing Caches

Address A:



Direct-Mapped Cache

Simplest kind of cache, easy to build
(only 1 tag compare required per access)

Characterized by exactly one line per set.

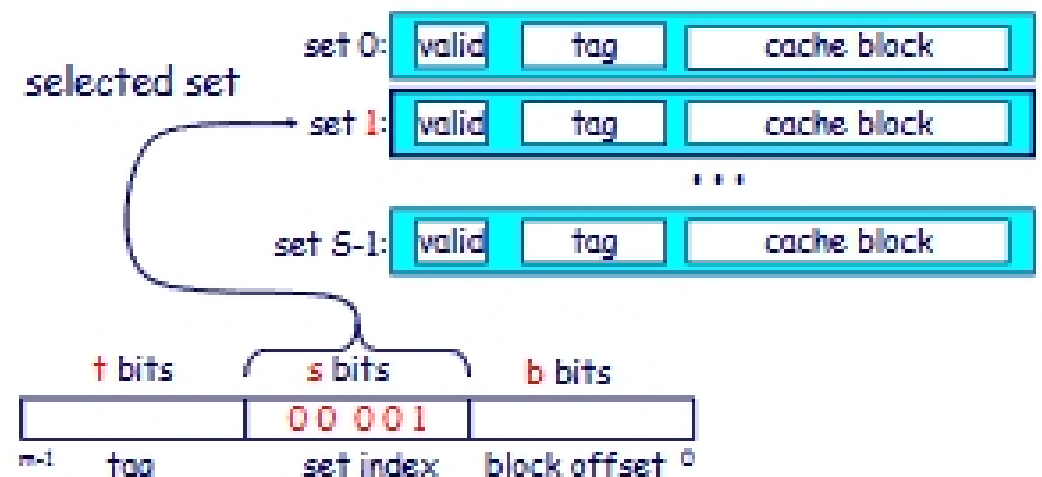


Cache size: $C = B \times S$ data bytes

Accessing Direct-Mapped Caches

Set selection

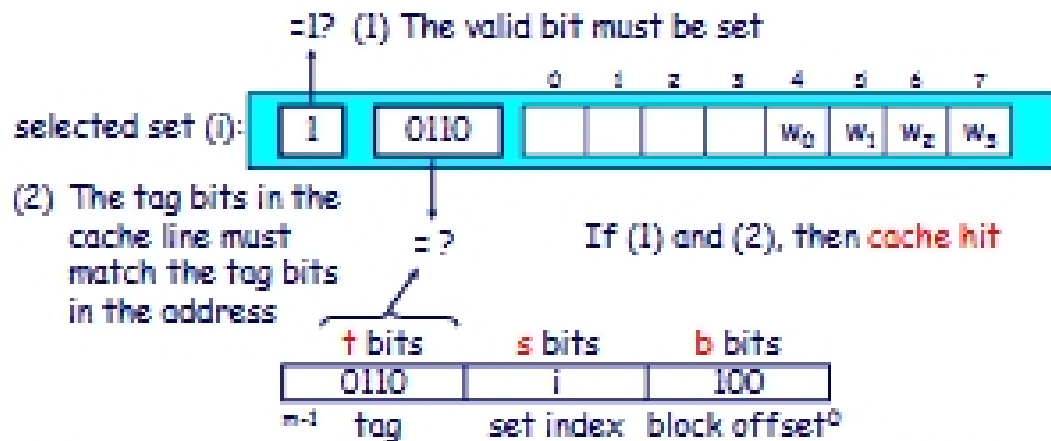
- Use the set index bits to determine the set of interest.



Accessing Direct-Mapped Caches

Line matching and word selection

- **Line matching:** Find a valid line in the selected set with a matching tag
- **Word selection:** Then extract the word



Accessing Direct-Mapped Caches

Line matching and word selection

- **Line matching:** Find a valid line in the selected set with a matching tag
- **Word selection:** Then extract the word



Direct-Mapped Cache Simulation

$M=16$ byte addresses, $B=2$ bytes/block,
 $S=4$ sets, $E=1$ entry/set

$t=1$ $s=2$ $b=1$

X	XX	X
---	----	---

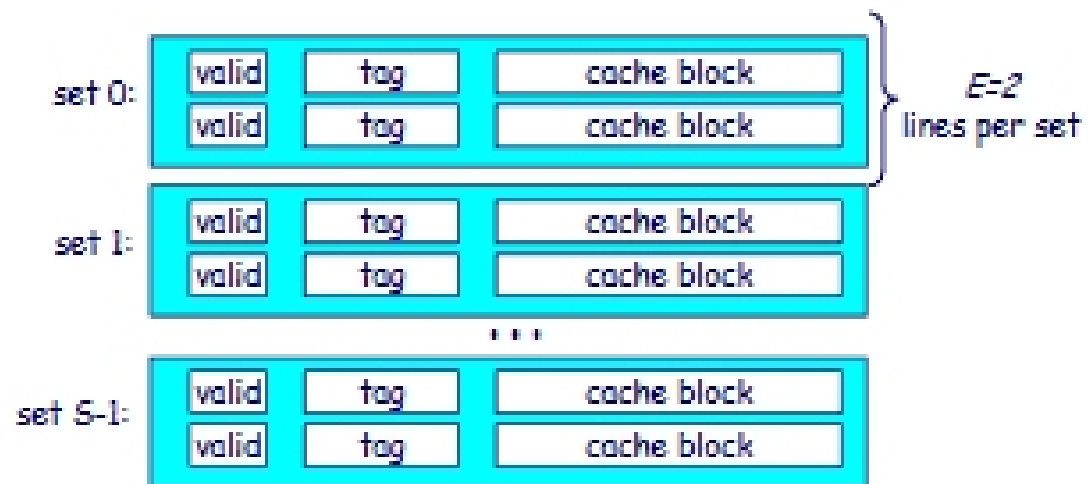
Address trace (reads):

0	[0000 ₂],	miss
1	[0001 ₂],	hit
7	[0111 ₂],	miss
8	[1000 ₂],	miss
0	[0000 ₂]	miss

v	tag	data
1	0	$M[0-1]$
1	0	$M[6-7]$

Set Associative Caches

Characterized by more than one line per set



E-way associative cache