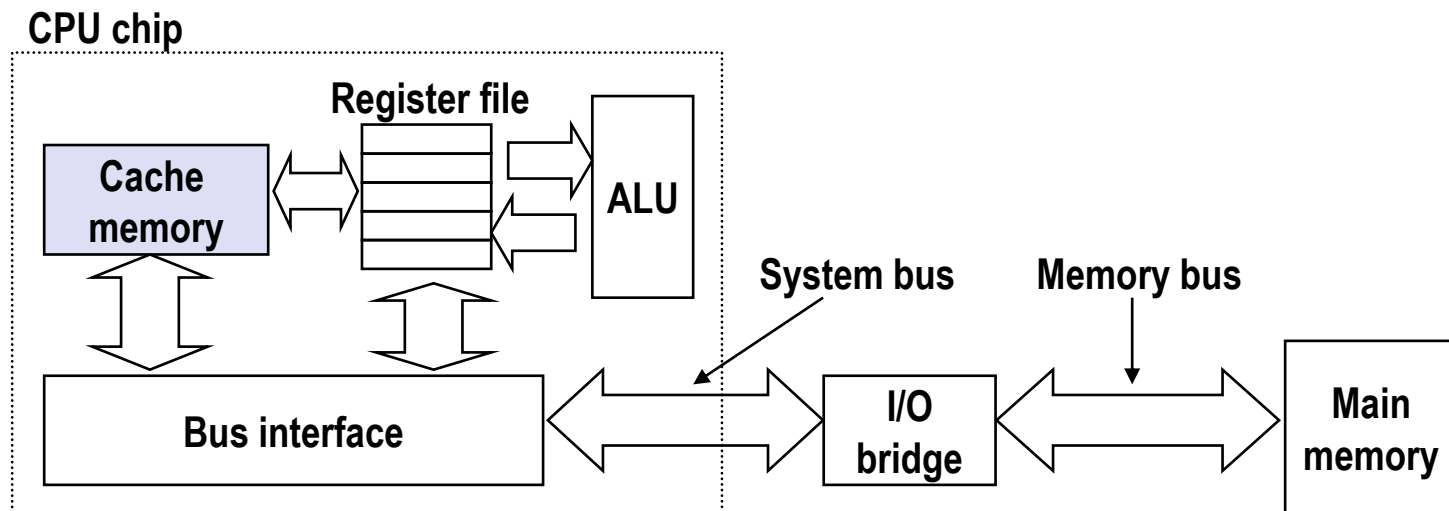# The Memory Hierarchy:
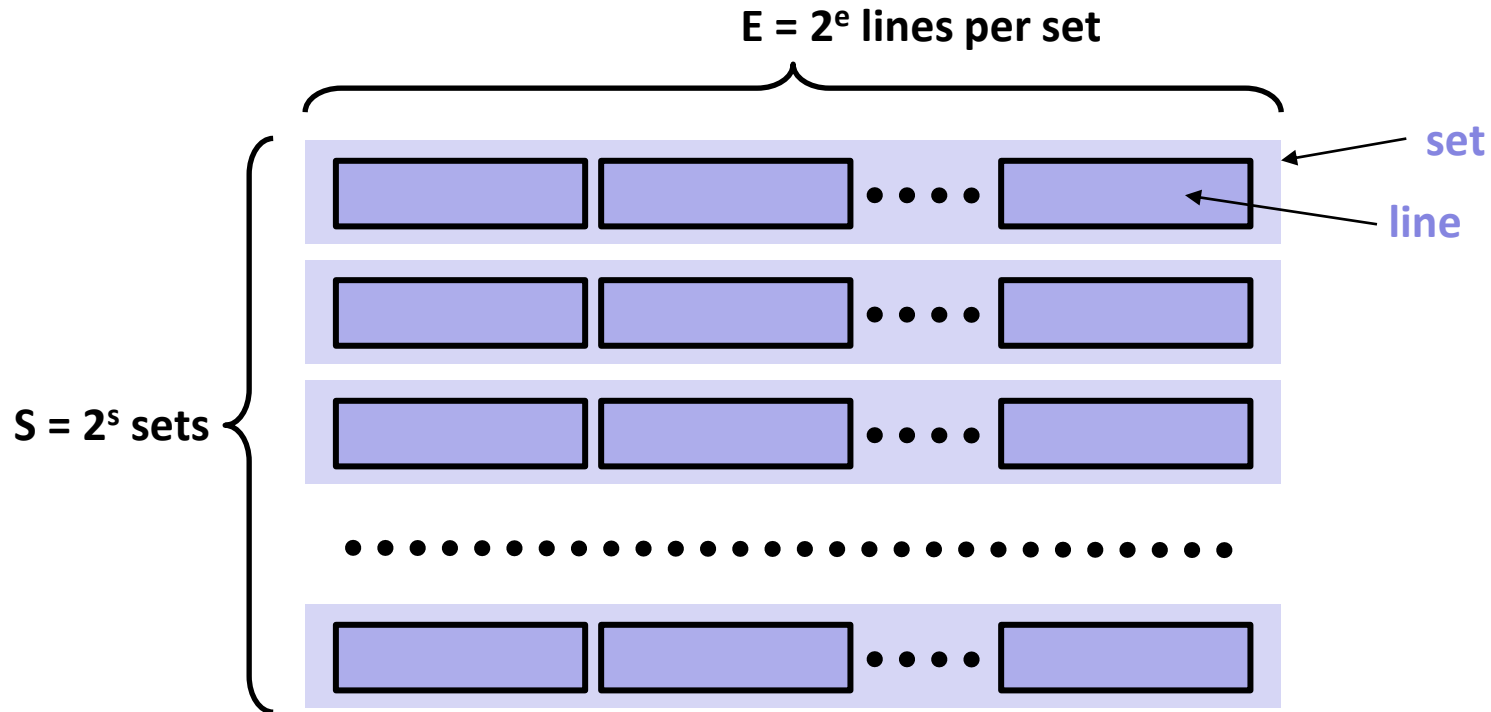## *Cache memory organization & operation*

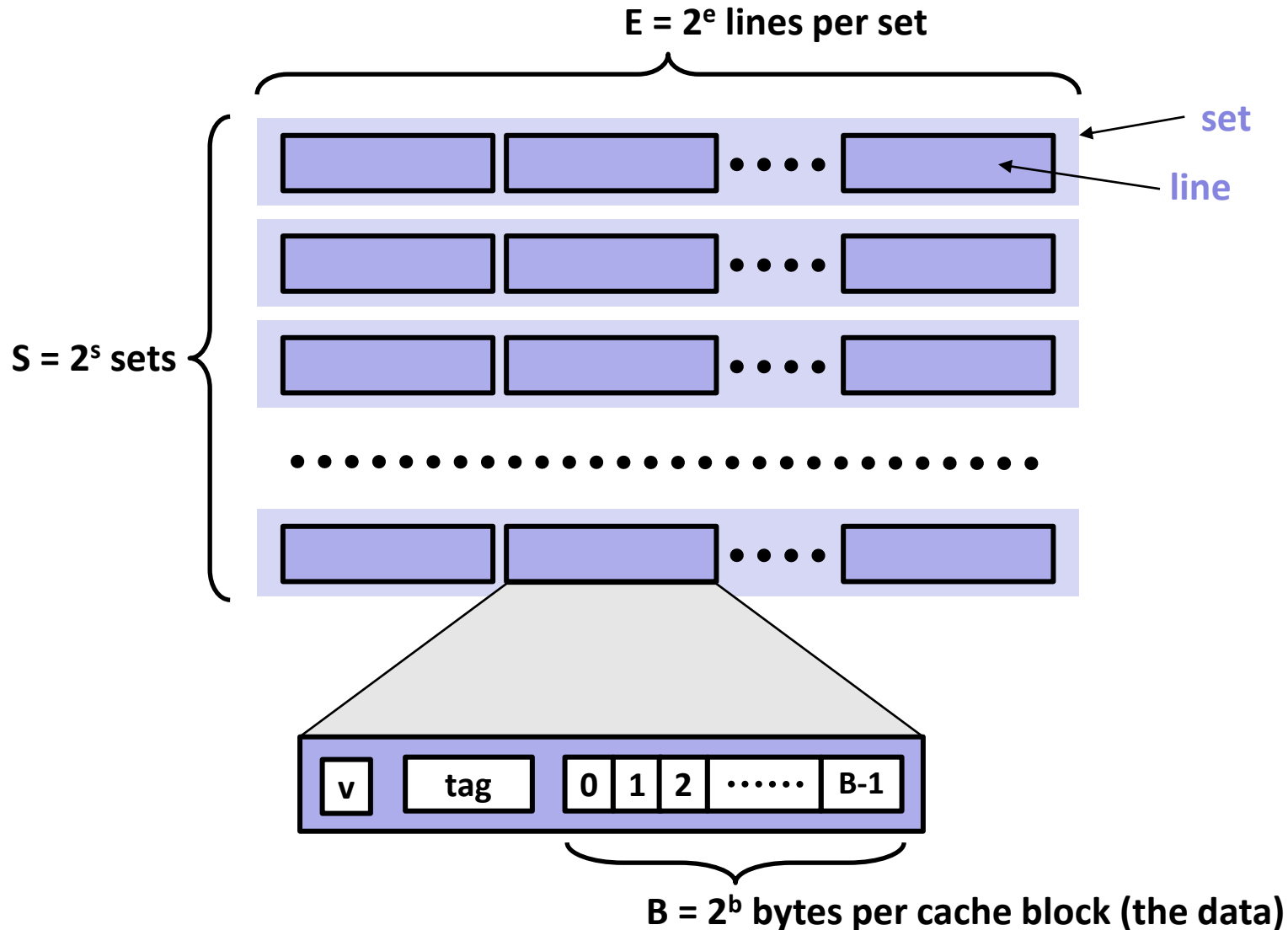# Cache Memories

- **<span style="color:red">Cache memories</span> are small, fast SRAM-based memories managed automatically in hardware**
  - Hold frequently accessed blocks of main memory
- **CPU looks first for data in cache**
- **Typical system structure:**



Prof. Michaël Cadilhac (based on slides by Bryant and O'Hallaron, CMU)

2

# General Cache Organization (S, E, B)

E = $2^e$ lines per set

S = $2^s$ sets

set

line

# General Cache Organization (S, E, B)

$E = 2^e$ lines per set



$S = 2^s$ sets

set

line

v | tag | 0 | 1 | 2 | ······ | B-1

$B = 2^b$ bytes per cache block (the data)

# General Cache Organization (S, E, B)

$E = 2^e$ lines per set



*set*

*line*

$S = 2^s$ sets

**valid bit**

v | tag | 0 | 1 | 2 | ······ | B-1

$B = 2^b$ bytes per cache block (the data)

*Cache size:*
*C = S x E x B data bytes*

# Cache Read

$E = 2^e$ lines per set

$S = 2^s$ sets

**Address of word:**

| t bits | s bits | b bits |
|--------|--------|--------|

tag      set index     block offset

# Cache Read

- *Locate set*

**E = $2^e$ lines per set**

**S = $2^s$ sets**

**Address of word:**

| t bits | s bits | b bits |
|--------|--------|--------|

tag    set index    block offset

# Cache Read

- *Locate set*
- *Check if any line in set has matching tag*
- *Yes + line valid: hit*

$E = 2^e$ lines per set

$S = 2^s$ sets

**Address of word:**

| t bits | s bits | b bits |
|--------|--------|--------|

tag      set index      block offset

| v | tag | 0 | 1 | 2 | ⋯⋯ | B-1 |
|---|-----|---|---|---|-----|-----|

**valid bit**

$B = 2^b$ **bytes per cache block (the data)**

Prof. Michaël Cadilhac (based on slides by Bryant and O'Hallaron, CMU)

8

# Cache Read

- *Locate set*
- *Check if any line in set has matching tag*
- *Yes + line valid: hit*
- *Locate data starting at offset*

$E = 2^e$ lines per set

$S = 2^s$ sets

**Address of word:**

| t bits | s bits | b bits |
|--------|--------|--------|

tag    set index    block offset

data begins at this offset

| v | tag | 0 | 1 | 2 | ...... | B-1 |

**valid bit**

$B = 2^b$ bytes per cache block (the data)

# Example: Direct Mapped Cache (E = 1)

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



**S = $2^s$ sets**

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

**Address of int:**

| t bits | 0...01 | 100 |

# Example: Direct Mapped Cache (E = 1)

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



**Address of int:**

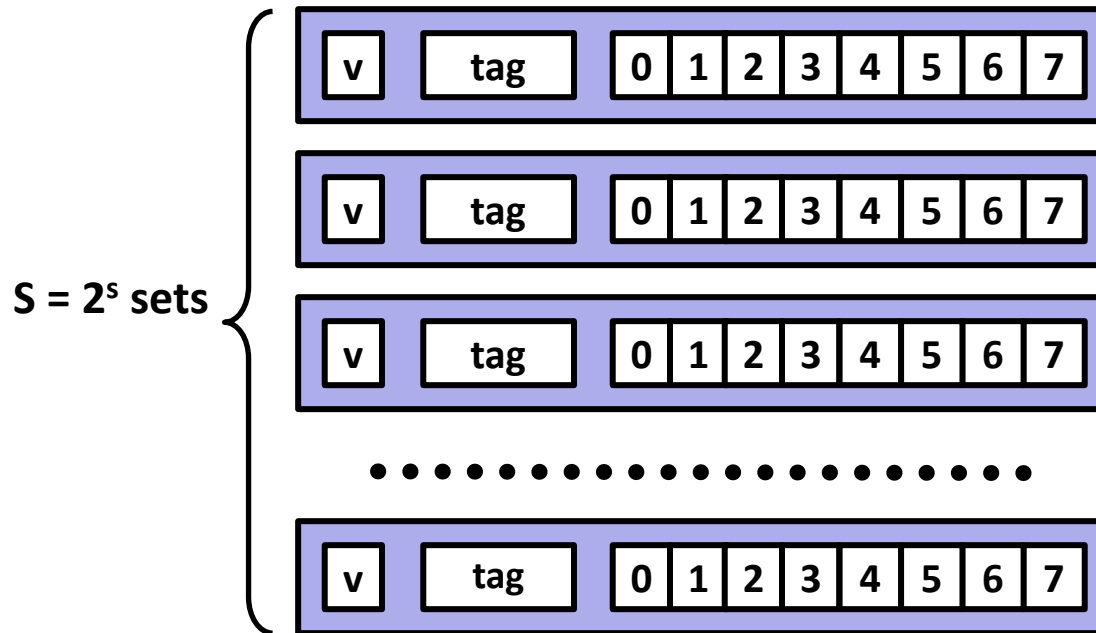| t bits | 0...01 | 100 |
|---|---|---|

**find set**

$S = 2^s$ **sets**

# Example: Direct Mapped Cache (E = 1)

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

**Address of int:**

| t bits | 0...01 | 100 |
|--------|--------|-----|

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|-----|---|---|---|---|---|---|---|---|

Prof. Michaël Cadilhac (based on slides by Bryant and O'Hallaron, CMU)

12

# Example: Direct Mapped Cache (E = 1)

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



valid?   +   match: assume yes = hit

Address of int:

| t bits | 0...01 | 100 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

# Example: Direct Mapped Cache (E = 1)

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

**valid?** + **match: assume yes = hit**

**Address of int:**

| t bits | 0...01 | 100 |
|--------|--------|-----|

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|-----|---|---|---|---|---|---|---|---|

**block offset**

# Example: Direct Mapped Cache (E = 1)

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

**valid?** + **match: assume yes = hit**

**Address of int:**

| t bits | 0...01 | 100 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

**block offset**

**int (4 Bytes) is here**

Prof. Michaël Cadilhac (based on slides by Bryant and O'Hallaron, CMU)

15

# Example: Direct Mapped Cache (E = 1)

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

**valid?** + **match: assume yes = hit**

**Address of int:**

| t bits | 0...01 | 100 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

**block offset**

↑

**int (4 Bytes) is here**

**If tag doesn't match: old line is evicted and replaced**

# Direct-Mapped Cache Simulation

| t=1 | s=2 | b=1 |
|:---:|:---:|:---:|
| x | xx | x |

M=16 bytes (4-bit addresses), B=2 bytes/block, S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| | |
|---|---|
| 0 | [$0000_2$], |
| 1 | [$0001_2$], |
| 7 | [$0111_2$], |
| 8 | [$1000_2$], |
| 0 | [$0000_2$] |

| | v | Tag | Block |
|---|:---:|:---:|:---:|
| Set 0 | 0 | ? | ? |
| Set 1 | | | |
| Set 2 | | | |
| Set 3 | | | |

# Direct-Mapped Cache Simulation

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

M=16 bytes (4-bit addresses), B=2 bytes/block, S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| | | |
|---|---|---|
| 0 | [$0000_2$], | miss |
| 1 | [$0001_2$], | |
| 7 | [$0111_2$], | |
| 8 | [$1000_2$], | |
| 0 | [$0000_2$] | |

| | v | Tag | Block |
|-------|---|-----|-------|
| Set 0 | 0 | ? | ? |
| Set 1 | | | |
| Set 2 | | | |
| Set 3 | | | |

# Direct-Mapped Cache Simulation

| t=1 | s=2 | b=1 |
|:---:|:---:|:---:|
| x | xx | x |

M=16 bytes (4-bit addresses), B=2 bytes/block, S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| | | |
|:---:|:---:|:---:|
| 0 | $[0\underline{00}0_2]$, | miss |
| 1 | $[0\underline{00}1_2]$, | |
| 7 | $[0\underline{11}1_2]$, | |
| 8 | $[1\underline{00}0_2]$, | |
| 0 | $[0\underline{00}0_2]$ | |

|  | v | Tag | Block |
|---|:---:|:---:|:---:|
| **Set 0** | 1 | 0 | M[0-1] |
| **Set 1** | | | |
| **Set 2** | | | |
| **Set 3** | | | |

# Direct-Mapped Cache Simulation

| t=1 | s=2 | b=1 |
|:---:|:---:|:---:|
| x | xx | x |

M=16 bytes (4-bit addresses), B=2 bytes/block, S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| 0 | [$0\underline{00}0_2$], | miss |
|---|---|---|
| 1 | [$0\underline{00}1_2$], | hit |
| 7 | [$0\underline{11}1_2$], | |
| 8 | [$1\underline{00}0_2$], | |
| 0 | [$0\underline{00}0_2$] | |

| | v | Tag | Block |
|---|:---:|:---:|:---:|
| **Set 0** | 1 | 0 | M[0-1] |
| **Set 1** | | | |
| **Set 2** | | | |
| **Set 3** | | | |

# Direct-Mapped Cache Simulation

t=1    s=2    b=1

| x | xx | x |
|---|----|---|

M=16 bytes (4-bit addresses), B=2 bytes/block,
S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| 0 | $[0\underline{00}0_2]$, | miss |
|---|---|---|
| 1 | $[0\underline{00}1_2]$, | hit |
| 7 | $[0\underline{11}1_2]$, | miss |
| 8 | $[1\underline{00}0_2]$, | |
| 0 | $[0\underline{00}0_2]$ | |

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

# Direct-Mapped Cache Simulation

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

M=16 bytes (4-bit addresses), B=2 bytes/block, S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| 0 | $[0\underline{00}0_2]$, | miss |
|---|---|---|
| 1 | $[0\underline{00}1_2]$, | hit |
| 7 | $[0\underline{11}1_2]$, | miss |
| 8 | $[1\underline{00}0_2]$, | |
| 0 | $[0\underline{00}0_2]$ | |

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

# Direct-Mapped Cache Simulation

| t=1 | s=2 | b=1 |
|:---:|:---:|:---:|
| x | xx | x |

M=16 bytes (4-bit addresses), B=2 bytes/block, S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| 0 | [$0\underline{00}0_2$], | miss |
|---|---|---|
| 1 | [$0\underline{00}1_2$], | hit |
| 7 | [$0\underline{11}1_2$], | miss |
| 8 | [$1\underline{00}0_2$], | miss |
| 0 | [$0\underline{00}0_2$] | |

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

# Direct-Mapped Cache Simulation

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

M=16 bytes (4-bit addresses), B=2 bytes/block,
S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| 0 | [$0\underline{00}0_2$], | miss |
| 1 | [$0\underline{00}1_2$], | hit |
| 7 | [$0\underline{11}1_2$], | miss |
| 8 | [$1\underline{00}0_2$], | miss |
| 0 | [$0\underline{00}0_2$] | |

|  | v | Tag | Block |
|------|---|-----|-------|
| Set 0 | 1 | 1 | M[8-9] |
| Set 1 |  |  |  |
| Set 2 |  |  |  |
| Set 3 | 1 | 0 | M[6-7] |

# Direct-Mapped Cache Simulation

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

M=16 bytes (4-bit addresses), B=2 bytes/block, S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| 0 | $[0\underline{00}0_2]$, | miss |
|---|---|---|
| 1 | $[0\underline{00}1_2]$, | hit |
| 7 | $[0\underline{11}1_2]$, | miss |
| 8 | $[1\underline{00}0_2]$, | miss |
| 0 | $[0\underline{00}0_2]$ | miss |

|  | v | Tag | Block |
|---|---|---|---|
| Set 0 | 1 | 1 | M[8-9] |
| Set 1 |  |  |  |
| Set 2 |  |  |  |
| Set 3 | 1 | 0 | M[6-7] |

# Direct-Mapped Cache Simulation

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

M=16 bytes (4-bit addresses), B=2 bytes/block,
S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| 0 | [0$\underline{00}$0$_2$], | miss |
|---|---------------------------|------|
| 1 | [0$\underline{00}$1$_2$], | hit |
| 7 | [0$\underline{11}$1$_2$], | miss |
| 8 | [1$\underline{00}$0$_2$], | miss |
| 0 | [0$\underline{00}$0$_2$] | miss |

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

# E-way Set Associative Cache (Here: E = 2)

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

**Address of short int:**

| t bits | 0…01 | 100 |
|--------|------|-----|

# E-way Set Associative Cache (Here: E = 2)

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

**Address of short int:**

| t bits | 0…01 | 100 |
|---|---|---|



find set

# E-way Set Associative Cache (Here: E = 2)

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

**Address of short int:**

| t bits | 0...01 | 100 |
|--------|--------|-----|

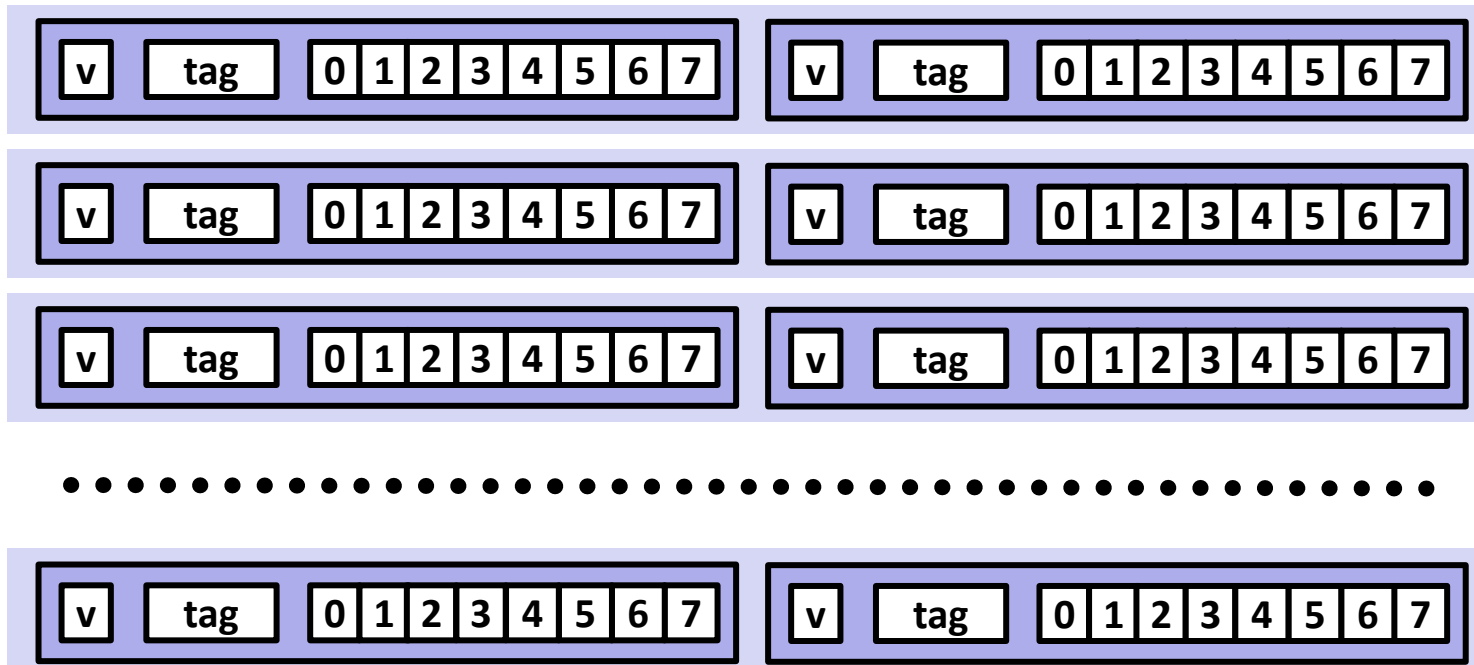| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

# E-way Set Associative Cache (Here: E = 2)

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

**Address of short int:**

| t bits | 0...01 | 100 |

**compare both**

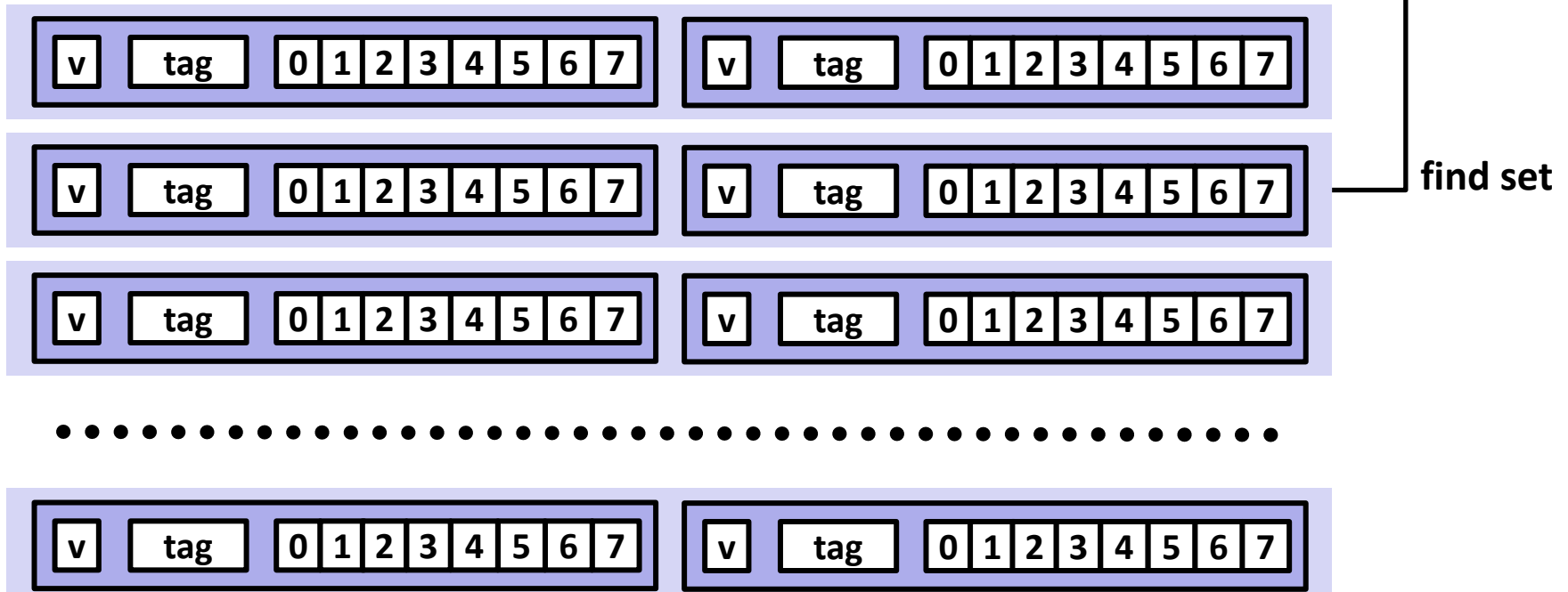| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

# E-way Set Associative Cache (Here: E = 2)

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

**Address of short int:**

| t bits | 0…01 | 100 |
|--------|------|-----|

**compare both**

**valid?  +  match: yes = hit**

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|-----|---|---|---|---|---|---|---|---|

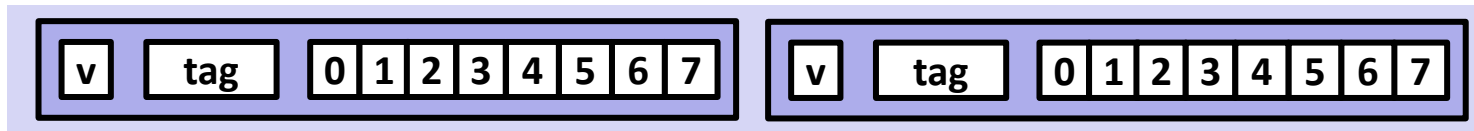| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|-----|---|---|---|---|---|---|---|---|

# E-way Set Associative Cache (Here: E = 2)

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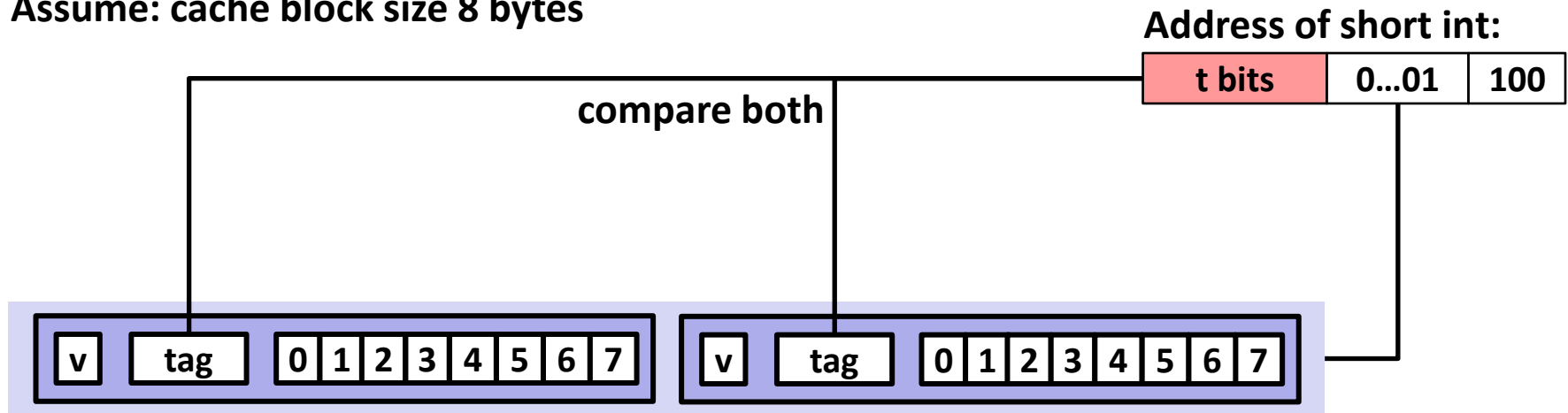**Assume: cache block size 8 bytes**

**Address of short int:**

| t bits | 0...01 | 100 |
|---|---|---|

**compare both**

**valid?  +  match: yes = hit**

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|

**block offset**

# E-way Set Associative Cache (Here: E = 2)

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

**Address of short int:**

| t bits | 0...01 | 100 |
|---|---|---|

**compare both**

**valid?  +  match: yes = hit**

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|

**block offset**

**short int (2 Bytes) is here**

# E-way Set Associative Cache (Here: E = 2)

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

**Address of short int:**

| t bits | 0…01 | 100 |
|--------|------|-----|

**compare both**

**valid?  +  match: yes = hit**

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

**block offset**

**short int (2 Bytes) is here**

## 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 | [00$\underline{0}$0$_2$], |
|---|---|
| 1 | [00$\underline{0}$1$_2$], |
| 7 | [01$\underline{1}$1$_2$], |
| 8 | [10$\underline{0}$0$_2$], |
| 0 | [00$\underline{0}$0$_2$] |

|  | v | Tag | Block |
|---|---|-----|-------|
| Set 0 | 0 | ? | ? |
|  | 0 | | |

|  | v | Tag | Block |
|---|---|-----|-------|
| Set 1 | 0 | | |
|  | 0 | | |

# 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 | $[00\underline{0}0_2]$, | miss |
| 1 | $[00\underline{0}1_2]$, | |
| 7 | $[01\underline{1}1_2]$, | |
| 8 | $[10\underline{0}0_2]$, | |
| 0 | $[00\underline{0}0_2]$ | |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 0 | 0 | ?   | ?     |
|       | 0 |     |       |
| Set 1 | 0 |     |       |
|       | 0 |     |       |

# 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 | $[00\underline{0}0_2]$, | miss |
| 1 | $[00\underline{0}1_2]$, | |
| 7 | $[01\underline{1}1_2]$, | |
| 8 | $[10\underline{0}0_2]$, | |
| 0 | $[00\underline{0}0_2]$ | |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 0 | 1 | 00  | M[0-1] |
|       | 0 |     |        |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 1 | 0 |     |        |
|       | 0 |     |        |

# 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 | $[00\underline{0}0_2]$, | miss |
|---|---|---|
| 1 | $[00\underline{0}1_2]$, | hit |
| 7 | $[01\underline{1}1_2]$, | |
| 8 | $[10\underline{0}0_2]$, | |
| 0 | $[00\underline{0}0_2]$ | |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 0 | 1 | 00  | M[0-1] |
|       | 0 |     |       |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 1 | 0 |     |       |
|       | 0 |     |       |

# 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 | [00$\underline{0}$0$_2$], | miss |
| 1 | [00$\underline{0}$1$_2$], | hit |
| 7 | [01$\underline{1}$1$_2$], | miss |
| 8 | [10$\underline{0}$0$_2$], | |
| 0 | [00$\underline{0}$0$_2$] | |

|  | v | Tag | Block |
|------|---|-----|-------|
| **Set 0** | 1 | 00 | M[0-1] |
|  | 0 | | |

|  | v | Tag | Block |
|------|---|-----|-------|
| **Set 1** | 0 | | |
|  | 0 | | |

# 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 | [$0000_2$], | miss |
|---|------------|------|
| 1 | [$0001_2$], | hit |
| 7 | [$0111_2$], | miss |
| 8 | [$1000_2$], | |
| 0 | [$0000_2$] | |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 0 | 1 | 00  | M[0-1] |
|       | 0 |     |       |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 1 | 1 | 01  | M[6-7] |
|       | 0 |     |       |

# 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 | [00$\underline{0}$0$_2$], | miss |
|---|---|---|
| 1 | [00$\underline{0}$1$_2$], | hit |
| 7 | [01$\underline{1}$1$_2$], | miss |
| 8 | [10$\underline{0}$0$_2$], | miss |
| 0 | [00$\underline{0}$0$_2$] | |

|  | v | Tag | Block |
|---|---|-----|-------|
| Set 0 | 1 | 00 | M[0-1] |
|  | 0 |  |  |

|  | v | Tag | Block |
|---|---|-----|-------|
| Set 1 | 1 | 01 | M[6-7] |
|  | 0 |  |  |

# 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 | [00$\underline{0}$0$_2$], | miss |
|---|---|---|
| 1 | [00$\underline{0}$1$_2$], | hit |
| 7 | [01$\underline{1}$1$_2$], | miss |
| 8 | [10$\underline{0}$0$_2$], | miss |
| 0 | [00$\underline{0}$0$_2$] | |

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

# 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$_2$], | miss |
|---|---|---|
| 1 | [0001$_2$], | hit |
| 7 | [0111$_2$], | miss |
| 8 | [1000$_2$], | miss |
| 0 | [0000$_2$] | 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, L3, Main Memory, Disk

# What about writes?

- **Multiple copies of data exist:**
  - L1, L2, L3, 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 bit (line different from memory or not)

Prof. Michaël Cadilhac (based on slides by Bryant and O'Hallaron, CMU)

45

# What about writes?

- **Multiple copies of data exist:**
  - L1, L2, L3, 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 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 straight to memory, does not load into cache)

# What about writes?

- **Multiple copies of data exist:**
  - L1, L2, L3, 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 bit (line different from memory or not)
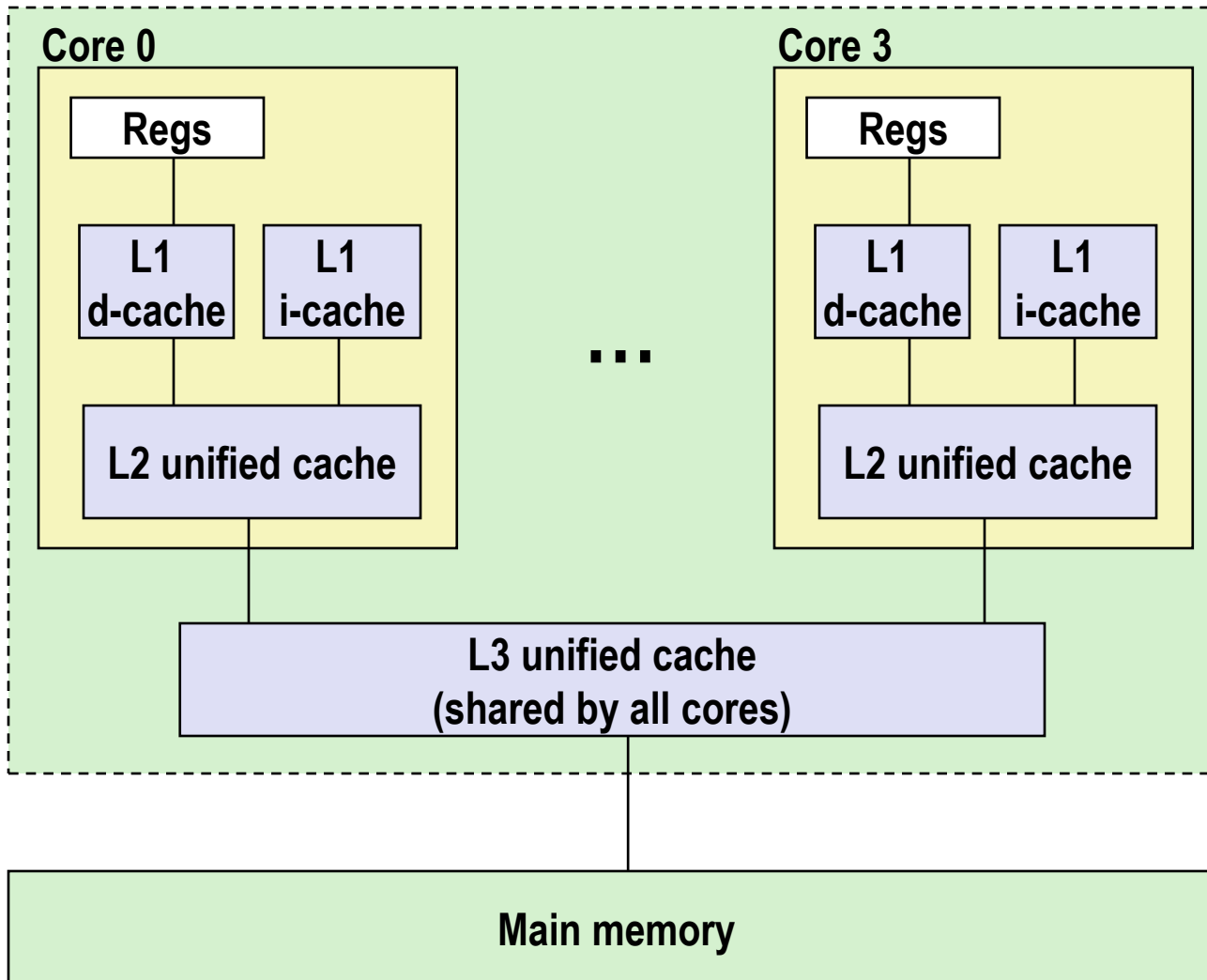- **What to do on a write-miss?**
  - Write-allocate (load into cache, update line in cache)
    - Good if more writes to the location follow
  - No-write-allocate (writes straight to memory, does not load into cache)
- **Typical**
  - Write-through + No-write-allocate
  - **Write-back + Write-allocate**

# Intel Core i7 Cache Hierarchy

**Processor package**



**L1 i-cache and d-cache:**
   32 KB,  8-way,
   Access: 4 cycles

**L2 unified cache:**
   256 KB, 8-way,
   Access: 10 cycles

**L3 unified cache:**
   8 MB, 16-way,
   Access: 40-75 cycles

**Block size**: 64 bytes for all caches.

# 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:
    - 4 clock cycle for L1
    - 10 clock cycles for L2

- **Miss Penalty**
  - Additional time required because of a miss
    - typically 50-200 cycles for main memory

# Let's think about those numbers

- **Huge difference between a hit and a miss**
  - Could be 100x, if just L1 and main memory

- **Would you believe 99% hits is twice as good as 97%?**
  - Consider:
  cache hit time of 1 cycle
  miss penalty of 100 cycles

# Let's think about those numbers

- **Huge difference between a hit and a miss**
  - Could be 100x, if just L1 and main memory

- **Would you believe 99% hits is twice as good as 97%?**
  - Consider:
    cache hit time of 1 cycle
    miss penalty of 100 cycles

  - Average access time:
    97% hits:  0.97 * 1 cycle + 0.03 * 100 cycles = **4 cycles**
    99% hits:  0.99 * 1 cycle + 0.01 * 100 cycles = **2 cycles**

# Let's think about those numbers

- **Huge difference between a hit and a miss**
  - Could be 100x, if just L1 and main memory

- **Would you believe 99% hits is twice as good as 97%?**
  - Consider:
    cache hit time of 1 cycle
    miss penalty of 100 cycles

  - Average access time:
    97% hits:  0.97 * 1 cycle + 0.03 * 100 cycles = **4 cycles**
    99% hits:  0.99 * 1 cycle + 0.01 * 100 cycles = **2 cycles**

- **This is why "miss rate" is used instead of "hit rate"**

# Writing Cache Friendly Code

- **Make the common case go fast**
  - Focus on the inner loops of the core functions

- **Minimize the misses in the inner loops**
  - Repeated references to variables are good (temporal locality)
  - Stride-1 reference patterns are good (spatial locality)

Prof. Michaël Cadilhac (based on slides by Bryant and O'Hallaron, CMU)

53

# Writing Cache Friendly Code

- **Make the common case go fast**
  - Focus on the inner loops of the core functions

- **Minimize the misses in the inner loops**
  - Repeated references to variables are good (temporal locality)
  - Stride-1 reference patterns are good (spatial locality)

**Key idea: Our qualitative notion of locality is quantified through our understanding of cache memories**