# CS61C

**Great Ideas
in
Computer Architecture**
(a.k.a. Machine Structures)

UC Berkeley
Teaching Professor
Lisa Yan

UC Berkeley
Lecturer
Justin Yokota

- EECS/CS **Undergraduate Town Hall**, Thursday April 11 4-5pm, Wozniak Lounge (tomorrow!)
- Data4All Kickoff event! **RSVP by today**. EdStem post #824 [link]

## Caches III: Direct Mapped

- Suppose the cache starts **cold**.

1. Load byte `0x43F`    `0x10F,0x3`
2. Load byte `0x5E2`    `0x178,0x2`
3. Load word `0x824`    `0x209,0x0`
4. Load word `0x5E0`    `0x178,0x0`

`0101 1110 0000`

a. **Cache hit**! Tag `0x178` is valid!

b. Read byte at `0x0` **offset** and return to processor.

| Val id | Tag | Data | | | |
|---|---|---|---|---|---|
| | | 11 | 10 | 01 | 00 |
| 1 | 0x10F | … | … | … | … |
| 1 | 0x178 | … | … | … | … |
| 1 | 0x209 | … | … | … | … |
| 0 | … | … | … | … | … |

Memory address:

| 11 | 2 | 1 0 |
|---|---|---|
| (10b) | | (2b) |

tag        offset

No access to main memory occurs on a **cache hit**.

Yan, Yokota

Berkeley
UNIVERSITY OF CALIFORNIA

- Cache line (i.e., cache) block
  - The smallest unit of memory that can be transferred between the main memory and the cache.
  - Each line has its own entry in the cache.
  - Line size/block size: The number of bytes in each cache line.

- When we bring data from the main memory into the cache, it is done in the granularity of a cache line (or cache block).
  - Typically **cache lines are 64 bytes**.
  - Cache blocks (cache lines) helps us take advantage of **spatial locality**.

- **Capacity**:
  - The total number of data bytes that can be stored in a cache.
  - For fully associative cache, **capacity = # lines * line size**.

# Hardware: Fully Associative Cache

- A Fully Associative Cache must compare tags in parallel:



Need hardware comparator for every single entry!

Sometimes infeasible; e.g., 1MB cache w/ 16B entries
→ $2^{20}/2^4 = 2^{16}$ = 16K comparators.

# Block Replacement Policies

- Block Replacement Policies
- Write Policies
- Direct Mapped Cache
- Types of Misses

# A Warmed up Cache Still Can Miss

- Suppose the cache starts **cold**.

1. Load byte `0x43F`    `0x10F`,`0x3`
2. Load byte `0x5E2`    `0x178`,`0x2`
3. Load word `0x824`    `0x209`,`0x0`
4. Load word `0x5E0`    `0x178`,`0x0`
5. Load word `0x524`    `0x149`,`0x0`

Fills up the cache

| Valid | Tag | Data 11 | Data 10 | Data 01 | Data 00 |
|-------|--------|-----|-----|-----|-----|
| 1 | 0x10F | … | … | … | … |
| 1 | 0x178 | … | … | … | … |
| 1 | 0x209 | … | … | … | … |
| 1 | 0x149 | … | … | … | … |

Yan, Yokota

# A Warmed up Cache Still Can Miss

- By the end of instr 5, cache is now **warm**.

1. Load byte `0x43F`    `0x10F,0x3`
2. Load byte `0x5E2`    `0x178,0x2`
3. Load word `0x824`    `0x209,0x0`
4. Load word `0x5E0`    `0x178,0x0`
5. Load word `0x524`    `0x149,0x0`

6. Load byte `0x972`    `1001 0111 0010`

| Valid | Tag | Data | | | |
|---|---|---|---|---|---|
| | | 11 | 10 | 01 | 00 |
| 1 | 0x10F | … | … | … | … |
| 1 | 0x178 | … | … | … | … |
| 1 | 0x209 | … | … | … | … |
| 1 | 0x149 | … | … | … | … |

**tag**
`0x25C`

**offset**
`0x2`

**Cache miss!**

When the FA cache **reaches capacity**, we can still miss.

Therefore we also need a policy for **block replacement**!

# Block Replacement Policies (i.e., Eviction)

- Least Recently Used (LRU)
  - Replace the entry that has not been used for the longest time, i.e., has the oldest previous access.
    - Pro: **temporal locality**!
      - recent past use implies likely future use
      - This is a very effective policy
    - Con: Complicated hardware to keep track of access history → performance hit

Yan, Yokota

- By the end of instr 5, cache is now **warm**.

1. Load byte `0x43F`    `0x10F`,`0x3`
2. Load byte `0x5E2`    `0x178`,`0x2`
3. Load word `0x824`    `0x209`,`0x0`
4. Load word `0x5E0`    `0x178`,`0x0`
5. Load word `0x524`    `0x149`,`0x0`

6. Load byte `0x972`    `1001 0111 0010`

tag `0x25C`   offset `0x2`

Suppose that LRU = 0 means **most** recently used, and 3 means **least** recently used. After the end of instruction 5, what are the LRU tags on each row?

| Valid | LRU | Tag | Data | | | |
|-------|-----|-----|------|------|------|------|
| | | | 11 | 10 | 01 | 00 |
| 1 | | 0x10F | … | … | … | … |
| 1 | | 0x178 | … | … | … | … |
| 1 | | 0x209 | … | … | … | … |
| 1 | | 0x149 | … | … | … | … |

A.
| LRU |
|-----|
| 0 |
| 1 |
| 2 |
| 3 |

B.
| LRU |
|-----|
| 3 |
| 2 |
| 1 |
| 0 |

C.
| LRU |
|-----|
| 3 |
| 1 |
| 2 |
| 0 |

D.
| LRU |
|-----|
| 0 |
| 2 |
| 3 |
| 1 |

E. Something else/don't know

Yan, Yokota

# At the end of instruction 5, what are the LRU tags on each row?



A

0%

B

0%

C

0%

D

0%

Something else/don't know

0%

# Fully Associative Cache with LRU policy

- By the end of instr 5, cache is now **warm**.

1. Load byte `0x43F`     `0x10F,0x3`
2. Load byte `0x5E2`     `0x178,0x2`
3. Load word `0x824`     `0x209,0x0`
4. Load word `0x5E0`     `0x178,0x0`
5. Load word `0x524`     `0x149,0x0`

6. Load byte `0x972`     `1001 0111 0010`

| Valid | LRU | Tag | Data | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | | | 11 | 10 | 01 | 00 |
| 1 | 3 | 0x10F | … | … | … | … |
| 1 | 1 | 0x178 | … | … | … | … |
| 1 | 2 | 0x209 | … | … | … | … |
| 1 | 0 | 0x149 | … | … | … | … |

**tag**        **offset**
0x25C         0x2

Suppose that LRU = 0 means **most** recently used, and 3 means **least** recently used. After the end of instruction 5, what are the LRU tags on each row?

Yan, Yokota

# Fully Associative Cache with LRU policy

- By the end of instr 5, cache is now **warm**.

1. Load byte `0x43F`    `0x10F`,`0x3`
2. Load byte `0x5E2`    `0x178`,`0x2`
3. Load word `0x824`    `0x209`,`0x0`
4. Load word `0x5E0`    `0x178`,`0x0`
5. Load word `0x524`    `0x149`,`0x0`

6. Load byte `0x972`    `0x25C`,`0x2`

| Valid | LRU | Tag | Data | | | |
|---|---|---|---|---|---|---|
| | | | 11 | 10 | 01 | 00 |
| 1 | 0 | 0x25C | … | … | … | … |
| 1 | 2 | 0x178 | … | … | … | … |
| 1 | 3 | 0x209 | … | … | … | … |
| 1 | 1 | 0x149 | … | … | … | … |

a. **Cache miss**! Tag `0x25C` is not valid.

Suppose our cache uses LRU (Least Recently Used).

b. **Replace according to block replacement policy**.
   Load into cache the **4-byte block** from `0x970` to `0x973`. Mark **valid bit**. Update block replacement policy fields.

c. Read byte at `0x2` **offset** and return to processor.

Yan, Yokota

# Block Replacement Policies (i.e., Eviction)

- Least Recently Used (LRU)
  - Replace the entry that has not been used for the longest time, i.e., has the oldest previous access.
    - Pro: **temporal locality**!
      - recent past use implies likely future use
      - This is a very effective policy
    - Con: Complicated hardware to keep track of access history → performance hit

- Most Recently Used (MRU)
  - Replace the entry that has the newest previous access.

- First In, First Out (FIFO)
  - Replace the oldest block in the set (queue).

- Last In, First Out (LIFO)
  - Replace the newest block in the set (stack).

- Random — Works surprisingly okay (when given a low temporal locality workload)

LRU is ideal for **temporal locality**, but the last three are used more commonly in practice.

Reasonable approximations to LRU and MRU without adding too much excess hardware.

Note both ignore order of accesses after (before) the first (last) access.

Yan, Yokota

# Write Policies

- Block Replacement Policies
- Write Policies
- Direct Mapped Cache
- Types of Misses

# What about Stores?

- By the end of instr 5, cache is now **warm**.

1. Load byte `0x43F`  `0x10F`,`0x3`
2. Load byte `0x5E2`  `0x178`,`0x2`
3. Load word `0x824`  `0x209`,`0x0`
4. Load word `0x5E0`  `0x178`,`0x0`
5. Load word `0x524`  `0x149`,`0x0`
6. Load byte `0x972`  `0x25C`,`0x2`

7. Store byte `0x524`  `0x149`,`0x0`

**Cache hit!...and then?**

| Valid | LRU | Tag | Data | | | |
|---|---|---|---|---|---|---|
| | | | 11 | 10 | 01 | 00 |
| 1 | 0 | 0x25C | … | … | … | … |
| 1 | 2 | 0x178 | … | … | … | … |
| 1 | 3 | 0x209 | … | … | … | … |
| 1 | 1 | 0x149 | … | … | … | … |

How do we handle stores?

Yan, Yokota

# Write-through vs. Write-back Policies

- **Store** instructions write to memory, which **changes values**.

- Hardware needs to ensure that cache and memory have consistent information.

- **Write-through**:
  - Write to the cache and memory at the same time.
  - (writes to memory take longer)

Yan, Yokota

Berkeley
UNIVERSITY OF CALIFORNIA

# Write-through vs. Write-back Policies

[LRU column fixed post-class]

- **Store** instructions write to memory, which **changes values**.

- Hardware needs to ensure that cache and memory have consistent information.

- **Write-through**:
  - Write to the cache and memory at the same time.
  - (writes to memory take longer)

- **Write-back**:
  - Write data in cache and set a **dirty bit** to 1.
  - When this block gets evicted from the cache (and "back" to memory), write to memory.

| Valid | Dirty | LRU | Tag | Data 11 | 10 | 01 | 00 |
|-------|-------|-----|-------|-----|-----|-----|-----|
| 1 | 0 | 1 | 0x25C | … | … | … | … |
| 1 | 0 | 2 | 0x178 | … | … | … | … |
| 1 | 0 | 3 | 0x209 | … | … | … | … |
| 1 | 1 | 0 | 0x149 | … | … | … | … |

Store byte 0x524      0x149,0x0

**Cache hit** w/**write-back**:
update data **within cache block**, and only write back to memory when this block is evicted.

Write-back policies require an additional **dirty bit** field.

Yan, Yokota

Berkeley
UNIVERSITY OF CALIFORNIA

# Fully Associative Cache w/ LRU, Write-Back

Note: LRU also updates for all blocks on writes as well!! Because you are "using" blocks. [LRU column fixed post-class]

| Valid | Dirty | LRU | Tag | Data | | | |
|---|---|---|---|---|---|---|---|
| | | | | 11 | 10 | 01 | 00 |
| 1 | 0 | 1 | 0x25C | … | … | … | … |
| 1 | 0 | 2 | 0x178 | … | … | … | … |
| 1 | 0 | 3 | 0x209 | … | … | … | … |
| 1 | 1 | 0 | 0x149 | … | … | … | … |

Store byte 0x524        0x149,0x0

**Cache hit** w/**write-back**:
update data **within cache block**, and only write back to memory when this block is evicted.

- **Write-back**:
  - Write data in cache and set a **dirty bit** to 1.
  - When this block gets evicted from the cache (and "back" to memory), write to memory.

Berkeley
UNIVERSITY OF CALIFORNIA

Yan, Yokota

# Write-through vs. Write-back Policies

- **Store** instructions write to memory, which **changes values**.

- Hardware needs to ensure that cache and memory have consistent information.

- **Write-through**:
  - Write to the cache and memory at the same time.
  - (writes to memory take longer)

  Very simple to implement

- **Write-back**:
  - Write data in cache and set a **dirty bit** to 1.
  - When this block gets evicted from the cache (and "back" to memory), write to memory.

  (typically) lower traffic to memory, because you likely write multiple times before evicting from cache

Yan, Yokota

Berkeley
UNIVERSITY OF CALIFORNIA

# Direct Mapped Cache

- Block Replacement Policies
- Write Policies
- Direct Mapped Cache
- Types of Misses

# Cache Design: Placement Policies

Fully Associative Cache

$\longleftrightarrow$

Set-Associative Cache

(out of scope)

**Direct Mapped Cache**

Put a new block anywhere

Put a new block in one specific place

(up next)

Fully associative caches need expensive hardware.

Yan, Yokota

# Direct Mapped Cache

- **Placement policy**: Each memory address is associated with exactly one possible block in the cache.
  - To check for existence in the cache, we only need to look in a single location in the cache.

| Valid | Dirty | Tag | Data | | | |
|---|---|---|---|---|---|---|
| | | | 11 | 10 | 01 | 00 |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |

How do we ensure this?

Yan, Yokota

| Val id | Dir ty | LRU | Tag | Data |
|--------|--------|-----|-----|------|
| … | … | … | … | … |
| … | … | … | … | … |
| … | … | … | … | … |
| … | … | … | … | … |

4B Direct Mapped Cache
(block size 1B)

0x                    0b

- Block starting @ 0 maps to index **0**.

0x | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 1A | 1B | 1C | 1D | 1E | 1F |

Memory

Yan, Yokota

Berkeley
UNIVERSITY OF CALIFORNIA

**4B Direct Mapped Cache**
(block size 1B)

| Valid | Dirty | LRU | Tag | Data |
|-------|-------|-----|-----|------|
| … | … | … | … | … |
| … | … | … | … | … |
| … | … | … | … | … |
| … | … | … | … | … |

0x　　　　　　　　0b

- Block starting @ 0 maps to index **0**.
- Block starting @ 4 maps to index **0**.

| 0x | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 1A | 1B | 1C | 1D | 1E | 1F |
|----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
|    |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |

**Memory**

Yan, Yokota

# (1/2) The Idea behind Direct Mapped Caches

| Valid | Dirty | LRU | Tag | Data |
|---|---|---|---|---|
| … | … | … | … | … |
| … | … | … | … | … |
| … | … | … | … | … |
| … | … | … | … | … |

**4B Direct Mapped Cache**
(block size 1B)

|  | 0x | 0b |
|---|---|---|

- Block starting @ 0 maps to index **0**.
- Block starting @ 4 maps to index **0**.
- Block starting @ 8 maps to index **0**.
- Block starting @ C maps to index **0**.
- etc.

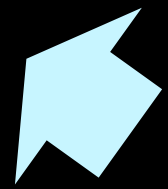Different addresses, get **same block index** but different **tags**.

| 0x | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 1A | 1B | 1C | 1D | 1E | 1F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

## Memory

**8B** Direct Mapped Cache
(block size **2B**)

| Valid | Dirty | LRU | Tag | Data | |
|---|---|---|---|---|---|
| | | | | 01 | 00 |
| … | … | … | … | … | … |
| … | … | … | … | … | … |
| … | … | … | … | … | … |
| … | … | … | … | … | … |

| 0x | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 1A | 1B | 1C | 1D | 1E | 1F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

Memory

| Valid | Dirty | LRU | Tag | Data 01 | Data 00 |
|---|---|---|---|---|---|
| … | … | … | … | … | … |
| … | … | … | … | | |
| … | … | … | … | | |
| … | … | … | … | | |

**8B** Direct Mapped Cache
(block size **2B**)

0x        0b

- Block starting @ 0 maps to index **0**.
- Block starting @ 8 maps to index **0**.
- Block starting @ 10 maps to index **0**.
- Block starting @ 18 maps to index **0**.
- etc.

Different addresses, get **same block index** but different **tags**.

| 0x | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 1A | 1B | 1C | 1D | 1E | 1F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Memory

Berkeley
UNIVERSITY OF CALIFORNIA

**8B** Direct Mapped Cache
(block size **2B**)

| Valid | Dirty | LRU | Tag | Data 01 | Data 00 |
|-------|-------|-----|-----|---------|---------|
| … | … | … | … | … | … |
| … | … | … | … | … | … |
| … | … | … | … | … | … |
| … | … | … | … | … | … |

| 0x | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 1A | 1B | 1C | 1D | 1E | 1F |
|----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|

Memory

Yan, Yokota

Berkeley
UNIVERSITY OF CALIFORNIA

# Direct Mapped Cache

- **Placement policy**: Each memory address is associated with exactly one possible block in the cache.
  - To check for existence in the cache, we only need to look in a single location in the cache.

| Valid | Dirty | Tag | Data | | | |
|---|---|---|---|---|---|---|
| | | | 11 | 10 | 01 | 00 |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |

Full address:

| 31 | 0 |
|---|---|
| | |

| tag | index | offset |
|---|---|---|

**tag** to check if have correct block

**index** to select block

**byte offset** within block

In direct mapped caches, split address to also get **index** of cache block to directly map to.

Yan, Yokota

Berkeley
UNIVERSITY OF CALIFORNIA

| Valid | Dirty | LRU | Tag | Data | |
|---|---|---|---|---|---|
| | | | | 01 | 00 |
| … | … | … | … | | |
| … | … | … | … | | |
| … | … | … | … | | |
| … | … | … | … | | |

**8B** Direct Mapped Cache
(block size **2B**)

0x                    0b      0b

- Block starting @ 0 maps to index **0**.  Tag 00
- Block starting @ 8 maps to index **0**.  Tag 01
- Block starting @ 10 maps to index **0**. Tag 10
- Block starting @ 18 maps to index **0**. Tag **11**

1 100x

Different addresses, get **same block index** but different **tags**.

| 0x | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 1A | 1B | 1C | 1D | 1E | 1F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

Memory

# Fill in the blank

| | Val id | Dir ty | Tag | Data 11 | 10 | 01 | 00 |
|---|---|---|---|---|---|---|---|
| | … | … | … | … | … | … | … |
| | … | … | … | … | … | … | … |
| | … | … | … | … | … | … | … |
| | … | … | … | … | … | … | … |

Suppose we have the following **direct mapped** cache, for **12b** addresses.

| | A. | B. | C. | D. | E. | F. |
|---|---|---|---|---|---|---|
| 1. What is the block size / line size, in bytes? | 2 | 4 | 8 | 10 | 16 | Other |
| 2. What is the capacity, in bytes? | 2 | 4 | 8 | 10 | 16 | Other |
| 3. For a 12b address, how many bits is the byte offset? | 2 | 4 | 8 | 10 | 16 | Other |
| 4. For a 12b address, how many bits is the index ? | 2 | 4 | 8 | 10 | 16 | Other |
| 5. For a 12b address, how many bits is the tag? | 2 | 4 | 8 | 10 | 16 | Other |

Yan, Yokota

# Fill in the blank

CS 61C

[activity slide]

Suppose we have the following **direct mapped** cache, for **12b** addresses.

| Valid | Dirty | LRU | Tag | Data 11 | Data 10 | Data 01 | Data 00 |
|---|---|---|---|---|---|---|---|
| ... | ... | | ... | ... | ... | ... | ... |
| ... | ... | | ... | ... | ... | ... | ... |
| ... | ... | | ... | ... | ... | ... | ... |
| ... | ... | | ... | ... | ... | ... | ... |

| | A. | B. | C. | D. | E. | F. |
|---|---|---|---|---|---|---|
| 1. What is the block size / line size, in bytes? | 2 | 4 | 8 | 10 | 16 | Other |
| 2. What is the capacity, in bytes? | 2 | 4 | 8 | 10 | 16 | Other |
| 3. For a 12b address, how many bits is the byte offset? | 2 | 4 | 8 | 10 | 16 | Other |
| 4. For a 12b address, how many bits is the index ? | 2 | 4 | 8 | 10 | 16 | Other |
| 5. For a 12b address, how many bits is the tag? | 2 | 4 | 8 | 10 | 16 | Other |

- Cache block/line: A single entry in the cache

1. Block size / line size: # bytes per cache block
   **4 bytes**

| Valid | Tag | Data | | | |
|---|---|---|---|---|---|
| | | 11 | 10 | 01 | 00 |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

2. Capacity      4 x 4 bytes = **16 bytes**
   - Total # data bytes that can be stored in a cache

3. Offset      $\log_2$(block size) = **2 bits**
   - Identifies byte offset of data stored within a given cache block

4. Index      $\log_2$(# lines) = **2 bits**
   - Selects block

5. Tag      # address bits - # offset bits - # index bits = **8 bits**
   - Identifies (checks) if given cache block is correct blocks

## Full Address

| Tag | Index | Byte Offset |
|---|---|---|

Read byte 0xFE2

T => 8 bits    I => 2 bits    O => 2 bits

0b 1111 1110 0010

T = 0xFE

I = 0x0

O = 0x2

| | Valid | Dirty | Tag | Data 11 | 10 | 01 | 00 |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0xFE | … | … | … | … |
| 1 | 0 | … | … | … | … | … | … |
| 2 | 0 | … | … | … | … | … | … |
| 3 | 0 | … | … | … | … | … | … |

4 Bytes

Cache Miss => bring in entire line

4B blocks, 4KB data

Much simpler than Fully Associative! Just check one tag.

# Types of Misses

- Block Replacement Policies

- Write Policies

- Direct Mapped Cache

- Types of Misses

# Direct Mapped: What is Possible?

- What policies can be implemented for a **direct-mapped cache**? Select all that apply.

1. Block Replacement (Eviction) Policy

   A. Least Recently Used
   B. Most Recently Used
   C. FIFO
   D. Random
   E. None of the above

2. Write-Back Policy

   A. Write-through (memory access per write)
   B. Write-back (dirty bit and write on eviction)

# What policies can be implemented for a direct-mapped cache? Select all that apply.

Replacement: LRU

0%

Replacement: MRU

0%

Replacement: FIFO

0%

Replacement: Random

0%

Replacement: None of the above

0%

Write: Write-through

0%

Write: Write-back

0%

# Direct Mapped: What is Possible?

- What policies can be implemented for a **direct-mapped cache**? Select all that apply.

1. Block Replacement (Eviction) Policy

   A. Least Recently Used
   B. Most Recently Used
   C. FIFO
   D. Random
   E. None of the above

2. Write Policy

   A. Write-through (memory access per write)
   B. Write-back (dirty bit and write to memory on eviction)

In direct mapped caches, there is only ever one block to evict—the existing block with matching index.

Yan, Yokota

# Types of Misses

- **Compulsory Miss**
    - Caused by the first access to a block that has never been in the cache

- **Capacity Miss**
    - Caused when the cache cannot contain all the blocks needed during the execution of a program
    - Occur when blocks were in the cache, replaced, and later retrieved

- **Conflict Miss**
    - Multiple blocks compete for the same location in the cache, even when the cache has not reached full capacity.

Yan, Yokota

# Types of Misses

- Compulsory Miss
  - Caused by the first access to a block that has never been in the cache

- Capacity Miss
  - Caused when the cache cannot contain all the blocks needed during the execution of a program
  - Occur when blocks were in the cache, replaced, and later retrieved

- Conflict Miss
  - Multiple blocks compete for the same location in the cache, even when the cache has not reached full capacity.

Which types of misses can occur for Fully Associative caches (FA)? For Direct Mapped caches (DM)? Select all that apply.

A. Compulsory
B. Capacity
C. Conflict
D. None of the above
E. All of the above

Yan, Yokota

# Which types of misses can occur for Fully Associative caches (FA)? For Direct Mapped caches (DM)? Select all that apply.

FA: Compulsory

0%

FA: Capacity

0%

FA: Conflict

0%

FA: None of the above

0%

DM: Compulsory

0%

DM: Capacity

0%

DM: Conflict

0%

DM: None of the above

0%

# Types of Misses

- Compulsory Miss
  - Caused by the first access to a block that has never been in the cache

- Capacity Miss
  - Caused when the cache cannot contain all the blocks needed during the execution of a program
  - Occur when blocks were in the cache, replaced, and later retrieved

- **Conflict Miss**
  - Multiple blocks compete for the same location in the cache, even when the cache has not reached full capacity.
  - Occurs in direct mapped caches, as well as in (out of scope) set associative caches.
  - Misses of this type would **not occur in a fully associative cache** with similar specs.

Yan, Yokota

# Summary: Cache Comparisons

- **Fully Associative**

  - A specific line of data can be stored in **any line** of the cache.

  - No index.

  - Need to choose replacement policy.

  - Need to choose write policy.

- **Direct Mapped**

  - A specific line of data can only be stored in **one index** of the cache.

  - Has index.

  - If the line you want to store the data in is occupied, you kick out that line.

  - Need to choose write policy.

Yan, Yokota