# Lecture 08

# Caches

이재진

서울대학교 컴퓨터공학부

http://aces.snu.ac.kr
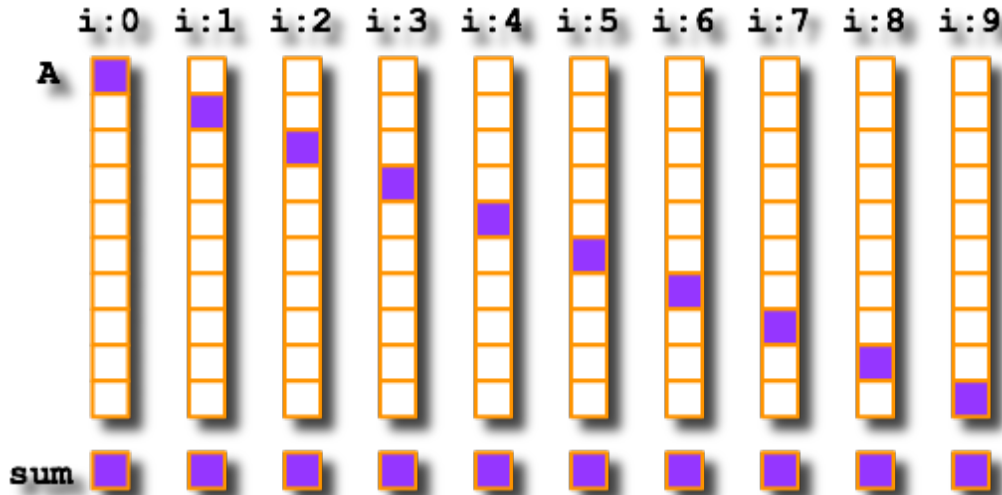
# Principle of Locality

- The reuse of data or instructions that were recently used, or near those that have been used recently
  - Predictable behavior
- Temporal locality
  - Recently referenced items are likely to be referenced in the near future
  - Within relatively small time durations
- Spatial locality
  - Items in nearby locations tend to be referenced close together in time
  - Within relatively close locations and relatively small time durations

# For Data

- Spatial locality
    - Reference array elements (A[i]) in succession (stride = 1)
- Temporal locality
    - Reference sum in each iteration

```
sum = 0;
for (i = 0; i < 10; i++)
  sum += A[i];
```

# For Instructions

- Spatial locality

    - Reference instructions in sequence

- Temporal locality

    - Cycle through loop repeatedly

```
sum = 0;
for (i = 0; i < 10; i++)
  sum += A[i];
```

```
    movl $0, -12(%ebp)
    movl $0, -16(%ebp)
    jmp  L2
L3:
    movl -16(%ebp), %eax
    movl -56(%ebp,%eax,4), %edx
    leal -12(%ebp), %eax
    addl %edx, (%eax)
    leal -16(%ebp), %eax
    incl (%eax)
L2:
    cmpl $9, -16(%ebp)
    jle  L3
```
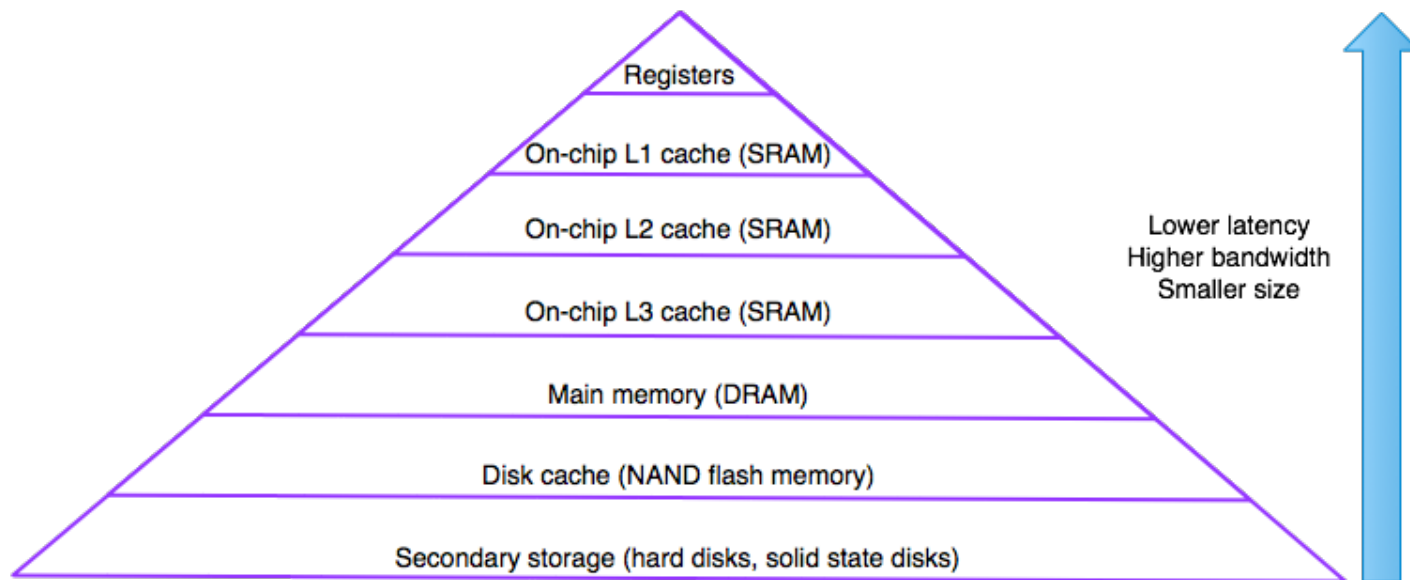
# Memory Hierarchies

- Hierarchical arrangement of storage

  - To exploit locality of reference

- Fast storage technologies cost more per byte and have less capacity

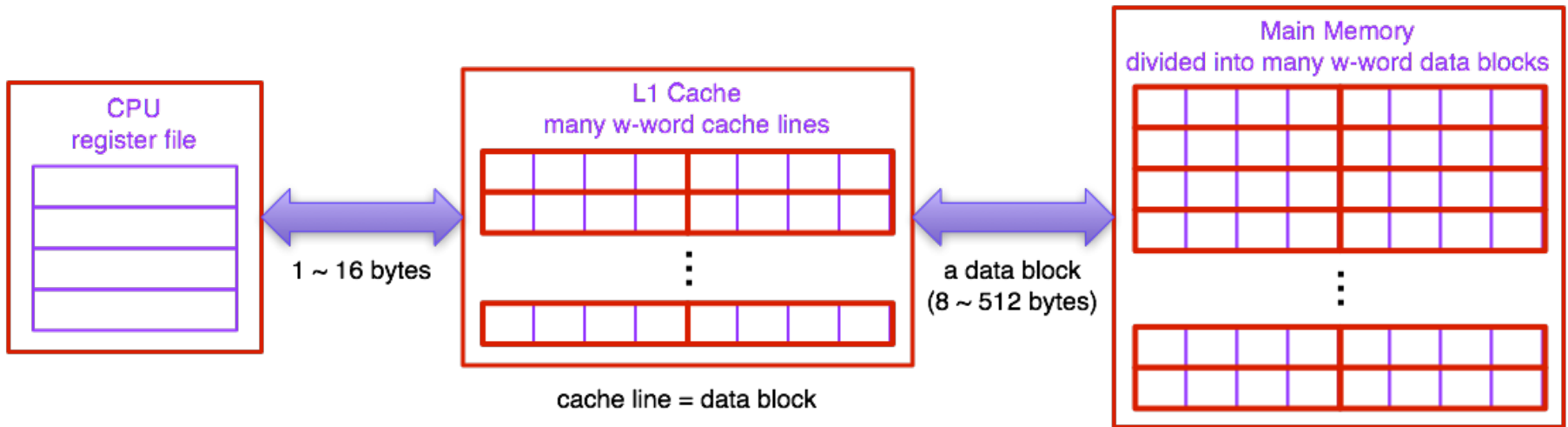- The gap between CPU and main memory speed is widening

# Caching

- Exploit temporal locality

    - Remember the contents of recently accessed locations

- Exploit spatial locality

    - Remember the blocks of recently accessed locations

- Cache block = cache line

    - The basic unit for cache storage

    - Multiple bytes or words

- Need an item d, which is stored in some block b

    - Cache hit

        - Find block b in the cache at level k

    - Cache miss

        - Block b is not in the cache at level k

        - The cache at level k must fetch b from level k+1

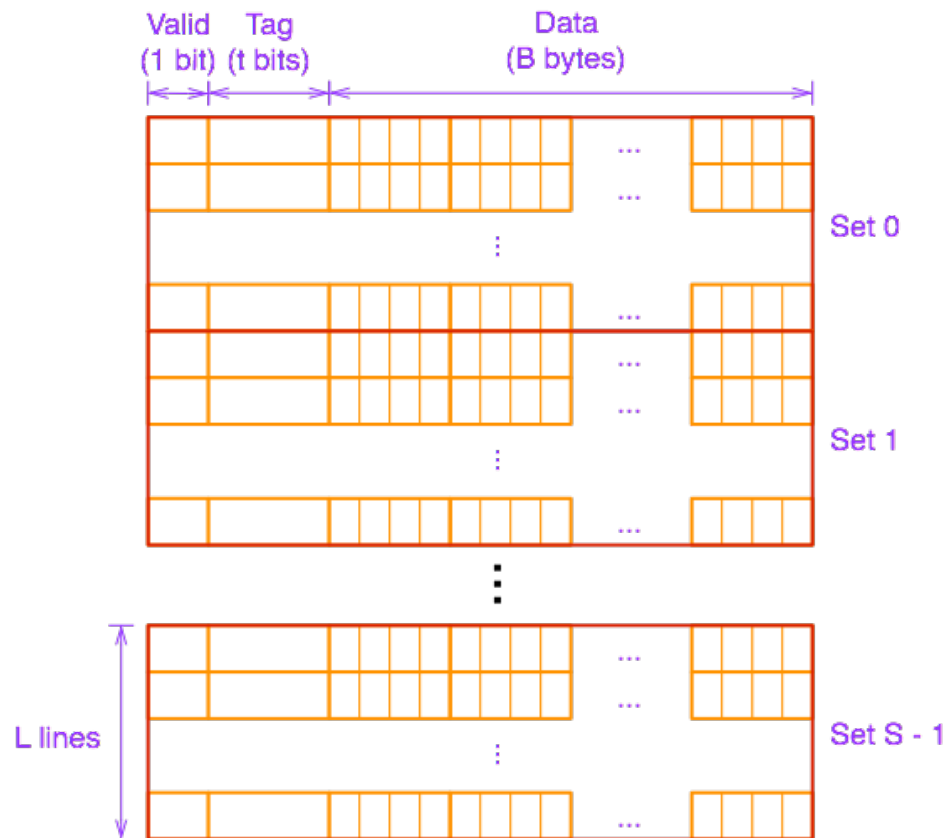            - If the cache at level k is full, then some block in the cache must be replaced
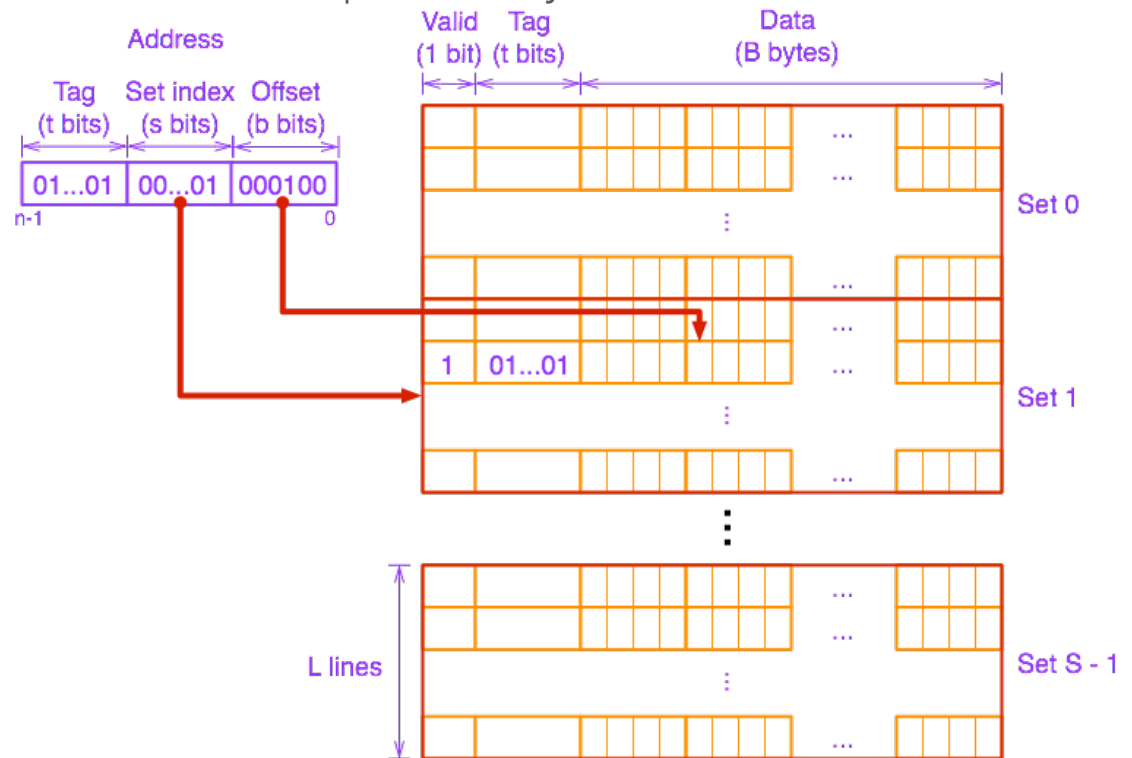
# L1 Cache between CPU and Main Memory



CPU
register file

1 ~ 16 bytes

L1 Cache
many w-word cache lines

cache line = data block

a data block
(8 ~ 512 bytes)

Main Memory
divided into many w-word data blocks

# Cache Organizations in General

- Cache size = L×S×B bytes

- A set is a collection of cache locations in which a given block may be placed
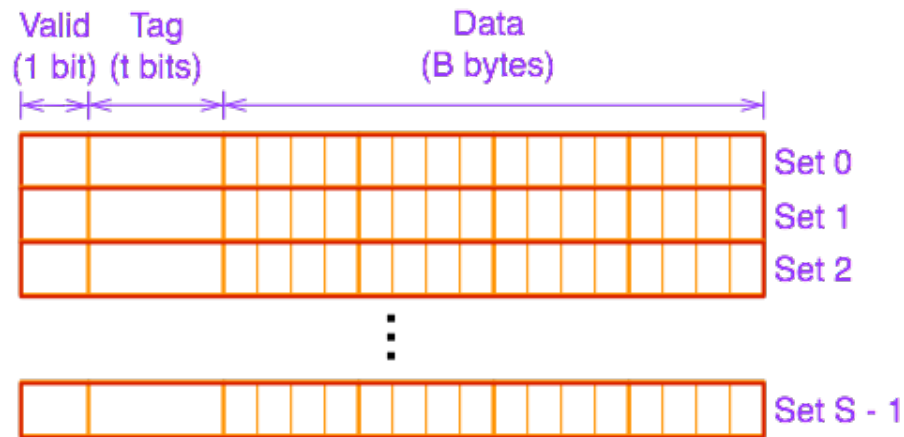
# Locating Data in the Cache

- The word at the requested address is in the cache if the tag bits in one of the valid lines in the specified set match the tag bits in the address
  - The set index is specified by the set index filed of the address
- The location of the word in the block is specified by the offset field in the address

$S = 2^s$
$B = 2^b$

THUNDER Research Group
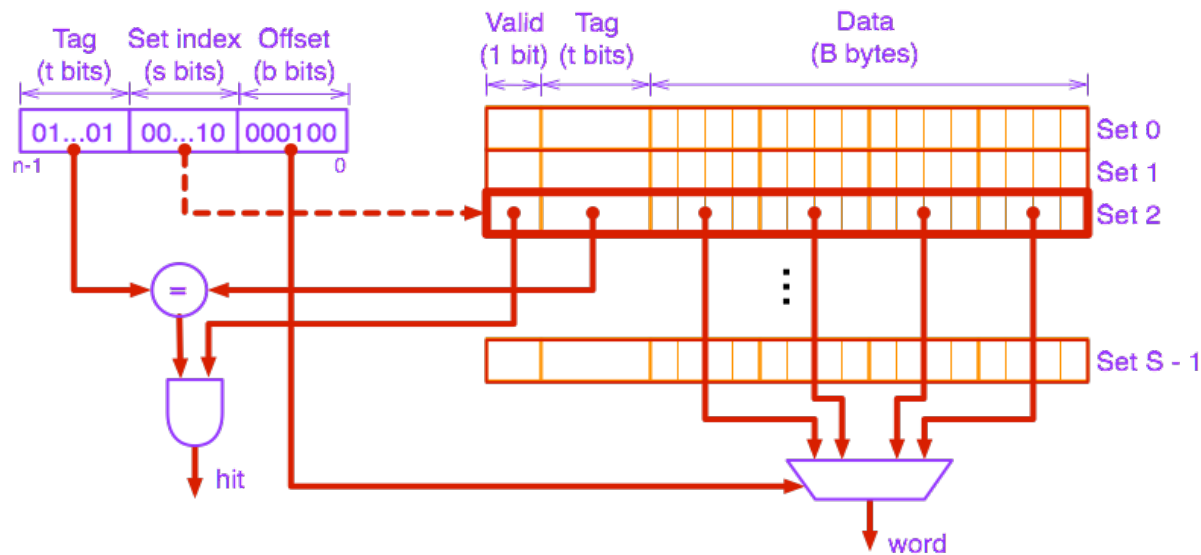Seoul National University
서울대학교 천둥 연구실

# Direct-Mapped Caches

- One cache line per set

- Simplest

- Data block can be only in one place in the cache

  - Replacement is straightforward

  - Collisions between data blocks for the same cache line can occur

**THUNDER Research Group**
Seoul National University
서울대학교 천둥 연구실

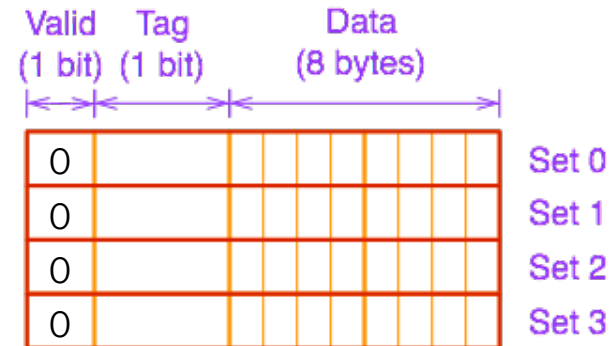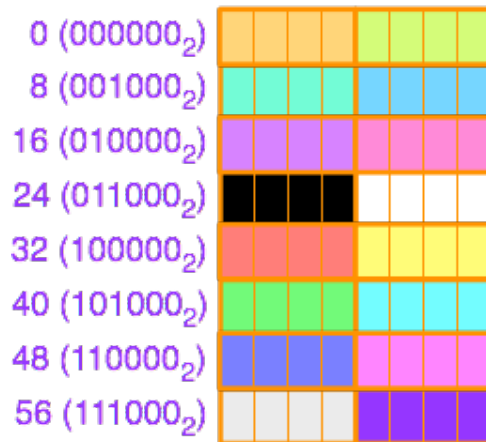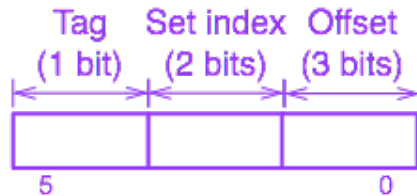# Addressing Direct-Mapped Caches

- Find a valid line in the selected set with a matching tag

- If there is one such line, extract the word with the offset field

- Otherwise, fetch the line from the lower level memory, place it in the selected set, and update the valid bit

# Addressing Direct-Mapped Caches (cont'd)

- Lower level memory size = 64 bytes

- B = 8 bytes/block, S = 4 sets, L = 1 line/set

- Address size = 6 bits

0 ($000000_2$)

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Addressing Direct-Mapped Caches (cont'd)

- Lower level memory size = 64 bytes

- B = 8 bytes/block, S = 4 sets, L = 1 line/set

- Address size = 6 bits

0 ($000000_2$)

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Addressing Direct-Mapped Caches (cont'd)

- Lower level memory size = 64 bytes

- B = 8 bytes/block, S = 4 sets, L = 1 line/set
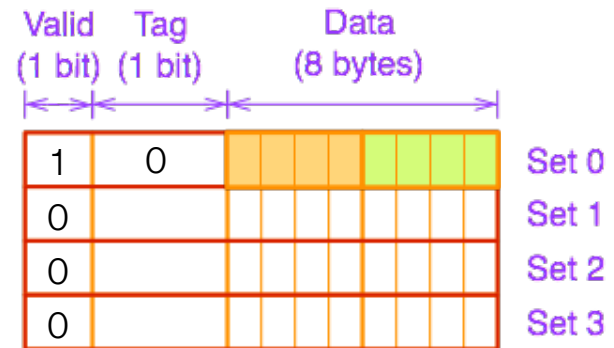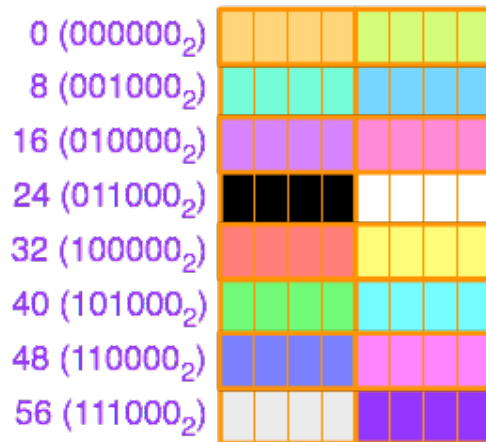
- Address size = 6 bits

0 (000000$_2$)          4 (000100$_2$)

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Addressing Direct-Mapped Caches (cont'd)

- Lower level memory size = 64 bytes

- B = 8 bytes/block, S = 4 sets, L = 1 line/set
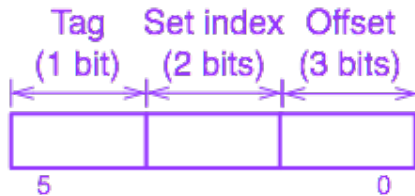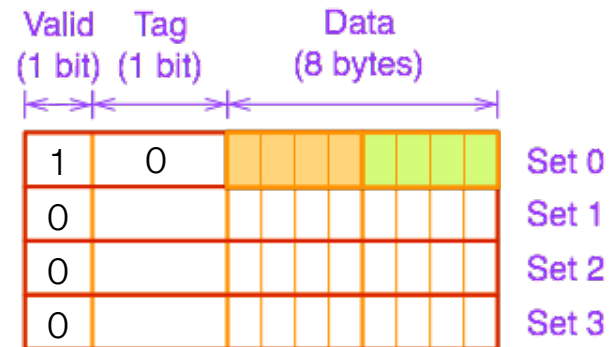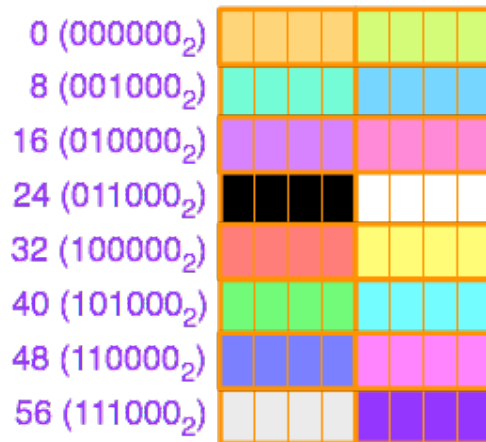
- Address size = 6 bits

0 ($000000_2$)          4 ($000100_2$)          20 ($010100_2$)

Tag          Set index          Offset
(1 bit)      (2 bits)           (3 bits)

5                                0

0 ($000000_2$)
8 ($001000_2$)
16 ($010000_2$)
24 ($011000_2$)
32 ($100000_2$)
40 ($101000_2$)
48 ($110000_2$)
56 ($111000_2$)

Valid      Tag          Data
(1 bit)    (1 bit)      (8 bytes)

| 1 | 0 | | | Set 0 |
| 0 | | | | Set 1 |
| 0 | | | | Set 2 |
| 0 | | | | Set 3 |

# Addressing Direct-Mapped Caches (cont'd)

- Lower level memory size = 64 bytes

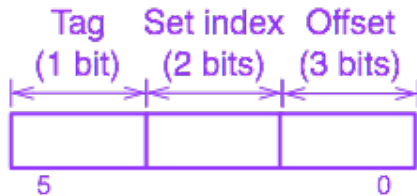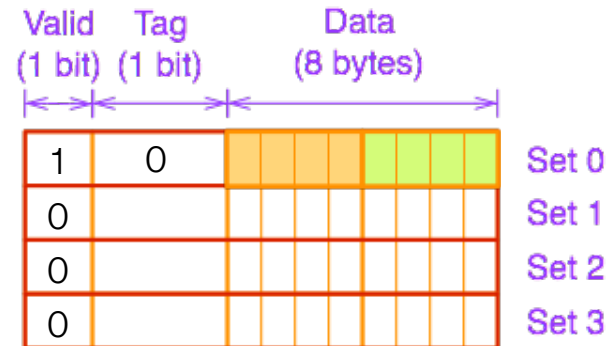- B = 8 bytes/block, S = 4 sets, L = 1 line/set

- Address size = 6 bits

0 ($000000_2$)        4 ($000100_2$)        20 ($010100_2$)

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Addressing Direct-Mapped Caches (cont'd)

- Lower level memory size = 64 bytes

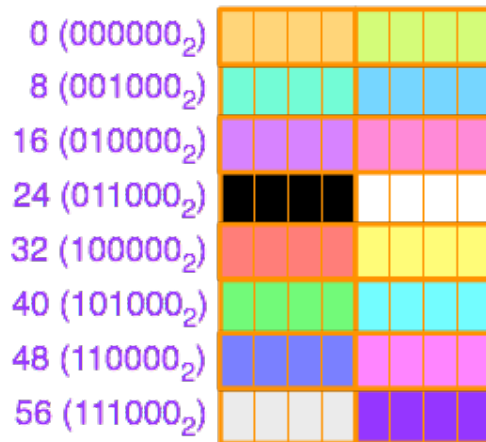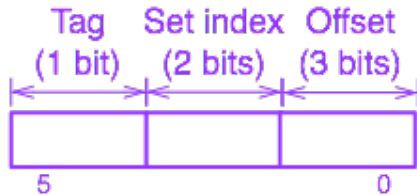- B = 8 bytes/block, S = 4 sets, L = 1 line/set

- Address size = 6 bits

0 ($000000_2$)     4 ($000100_2$)     20 ($010100_2$)     48 ($110000_2$)

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Addressing Direct-Mapped Caches (cont'd)

- Lower level memory size = 64 bytes

- B = 8 bytes/block, S = 4 sets, L = 1 line/set
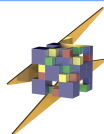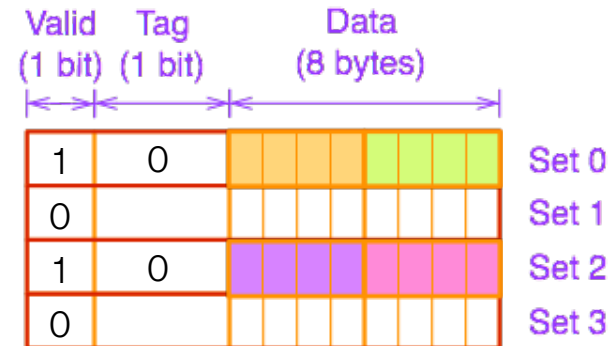
- Address size = 6 bits

0 ($000000_2$)        4 ($000100_2$)        20 ($010100_2$)        48 ($110000_2$)

**THUNDER Research Group**
Seoul National University
서울대학교 천둥 연구실

# Addressing Direct-Mapped Caches (cont'd)

- Lower level memory size = 64 bytes

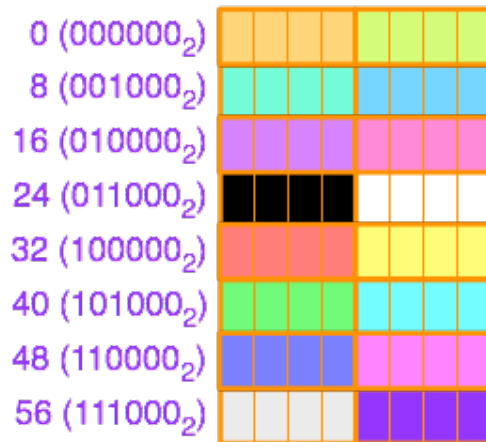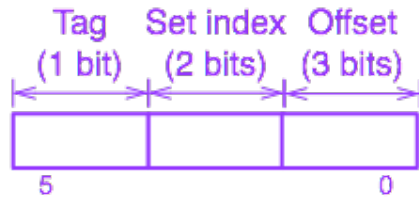- B = 8 bytes/block, S = 4 sets, L = 1 line/set

- Address size = 6 bits

$0 \ (000000_2)$      $4 \ (000100_2)$      $20 \ (010100_2)$      $48 \ (110000_2)$      $36 \ (100100_2)$

# Addressing Direct-Mapped Caches (cont'd)

- Lower level memory size = 64 bytes

- B = 8 bytes/block, S = 4 sets, L = 1 line/set
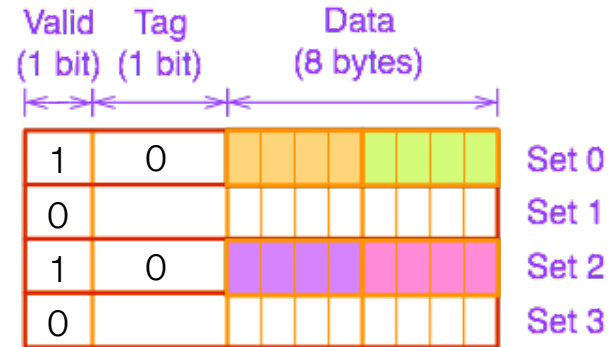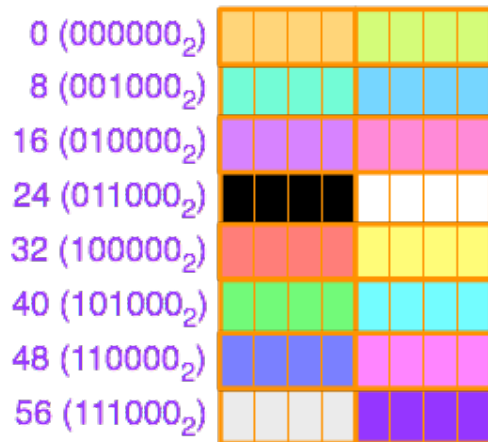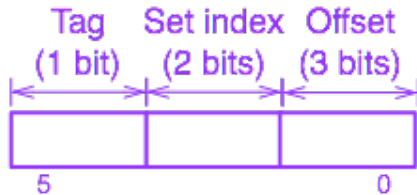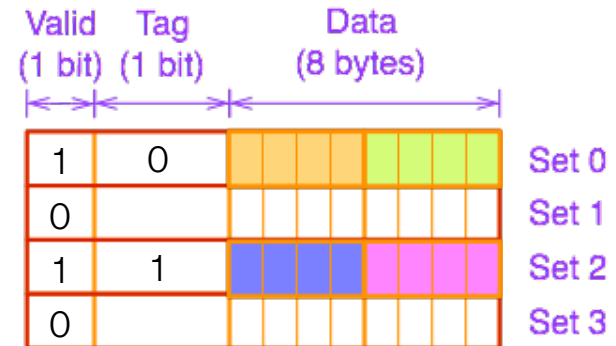
- Address size = 6 bits

0 ($000000_2$)　　　4 ($000100_2$)　　　20 ($010100_2$)　　　48 ($110000_2$)　　　36 ($100100_2$)

# Set Associative Caches

- Data block can be in a few places in the cache
    - Need a good replacement policy
    - Less collisions between data blocks for the same cache line than the direct-mapped cache
- Complex tag comparison hardware on the lines in a set

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Addressing Set Associative Caches

- Find a valid line in the selected set with a matching tag

- If there is one such line, extract the word with the offset field

- Otherwise, fetch the line from the lower level memory, place it in the selected set by deciding which line should be used, and update the valid bit

  - Need a sophisticated replacement policy

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Fully Associative Caches

- Only one set

- Data block can be any place in the cache
    - Less collisions between data blocks for the same cache line than the set associative cache

- Complex tag comparison hardware on the lines in the cache

# Addressing Fully Associative Caches

- Find a valid line with a matching tag

- If there is one such line, extract the word with the offset field

- Otherwise, fetch the line from the lower level memory, place it in the cache by deciding which line should be used, and update the valid bit
  - Need a sophisticated replacement policy

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Types of Cache Misses

- Cold (compulsory) miss

  - When the cache is empty

- Conflict miss

  - When the cache is large enough, but multiple data items map to the same cache line

- Capacity miss

  - When the set of active cache lines (working set) is larger than the cache

  - Working set

    - The set of referenced blocks that are active during a given period of time

# Replacement Policies

- After a miss, what cache block should be replaced with the block read from memory?
  - Which way in a multiway (i.e., set associative or fully associative) cache should be replaced?
  - Ideally, any cached data which is no longer needed  would be chosen to be replaced
- LRU (Least Recently Used)
- Pseudo LRU
- FIFO (First In, First Out)
  - Select a block that has been in the set for the longest time
- Random

# Least Recently Used (LRU)

- Select a block that has not been used for the longest time
  - Need to maintain LRU statistics for each cache line in a set
    - 2-way set associative cache: 1 bit to encode 2 states in a set
    - 4-way set associative cache: 5 bits to encode 4! = 24 states in a set
    - 8-way set associative cache: 16 bits to encode 8! = 40320 states in a set
    - ...
- A time consuming read/modify/write cycle is needed to maintain the set state on a cache access
  - Too costly
- Instead, use pseudo LRU

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Least Recently Used (cont'd)

- Select a block that has not been used for the longest time
  - Need to maintain LRU statistics for each cache line in a set
    - 2-way set associative cache: 1 bit to encode 2 states in a set
    - 4-way set associative cache: 5 bits to encode 4! = 24 states in a set
    - 8-way set associative cache: 16 bits to encode 8! = 40320 states in a set
    - ...
- A time consuming read/modify/write cycle is needed to maintain the set state on a cache access
  - Too costly
- Instead, use pseudo LRU

... A

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Least Recently Used (cont'd)

- Select a block that has not been used for the longest time
  - Need to maintain LRU statistics for each cache line in a set
    - 2-way set associative cache: 1 bit to encode 2 states in a set
    - 4-way set associative cache: 5 bits to encode $4! = 24$ states in a set
    - 8-way set associative cache: 16 bits to encode $8! = 40320$ states in a set
    - ...
- A time consuming read/modify/write cycle is needed to maintain the set state on a cache access
  - Too costly
- Instead, use pseudo LRU

## ... A B

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Pseudo LRU

- A binary decision tree
  - 2-way set associative cache: 1 bit
  - 4-way set associative cache: $(2^3 - 1) - 4 = 3$ bits
  - N-way set associative cache: $(2^{\log_2 N+1} - 1) - N$ bits
- The difference between pseudo LRU and true LRU is statistically small
- Each bit represents the left or right child in the binary decision tree
  - 1: the left side has been referenced more recently than the right side
  - 0: vice versa
- A write cycle to update the pseudo-LRU bits on a hit
- A read cycle for the pseudo-LRU bits during a line replacement

| access | next state |
|--------|-----------|
| A | 11_ |
| B | 10_ |
| C | 0_1 |
| D | 0_0 |

| state | replace |
|-------|---------|
| 00X | A |
| 01X | B |
| 1X0 | C |
| 1X1 | D |

| AB/CD | A/B | C/D |
|-------|-----|-----|
| 0 | 1 | 0 |

THUNDER Research Group
Seoul National University
서울대학교 천둥 연구실

# Pseudo LRU (cont'd)
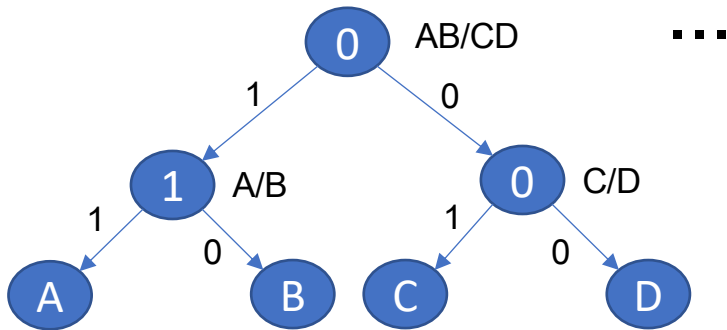
- A binary decision tree
  - 2-way set associative cache: 1 bit
  - 4-way set associative cache: $(2^3 - 1) - 4 = 3$ bits
  - N-way set associative cache: $(2^{\log_2 N + 1} - 1) - N$ bits
- The difference between pseudo LRU and true LRU is statistically small
- Each bit represents the left or right child in the binary decision tree
  - 1: the left side has been referenced more recently than the right side
  - 0: vice versa
- A write cycle to update the pseudo-LRU bits on a hit
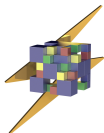- A read cycle for the pseudo-LRU bits during a line replacement

| access | next state |
|--------|-----------|
| A | 11_ |
| B | 10_ |
| C | 0_1 |
| D | 0_0 |

| state | replace |
|-------|---------|
| 00X | A |
| 01X | B |
| 1X0 | C |
| 1X1 | D |

... A

| AB/CD | A/B | C/D |
|-------|-----|-----|
| 0 | 1 | 0 |

THUNDER Research Group
Seoul National University
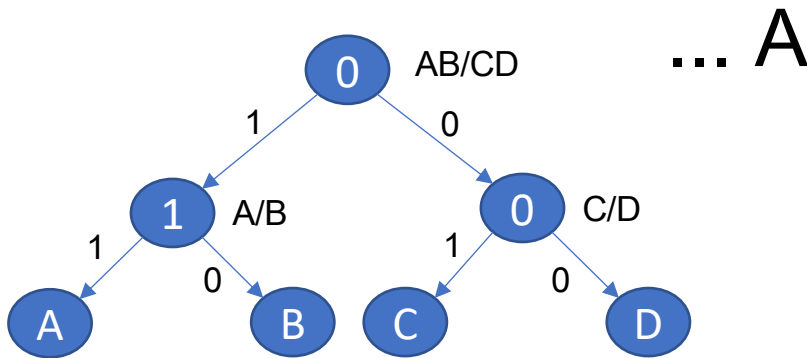서울대학교 천둥 연구실

# Pseudo LRU (cont'd)

- A binary decision tree
  - 2-way set associative cache: 1 bit
  - 4-way set associative cache: $(2^3 - 1) - 4 = 3$ bits
  - N-way set associative cache: $(2^{\log_2 N+1} - 1) - N$ bits
- The difference between pseudo LRU and true LRU is statistically small
- Each bit represents the left or right child in the binary decision tree
  - 1: the left side has been referenced more recently than the right side
  - 0: vice versa
- A write cycle to update the pseudo-LRU bits on a hit
- A read cycle for the pseudo-LRU bits during a line replacement

| access | next state |
|--------|-----------|
| A | 11_ |
| B | 10_ |
| C | 0_1 |
| D | 0_0 |

| state | replace |
|-------|---------|
| 00X | A |
| 01X | B |
| 1X0 | C |
| 1X1 | D |

... A

| AB/CD | A/B | C/D |
|-------|-----|-----|
| 1 | 1 | 0 |

THUNDER Research Group
Seoul National University
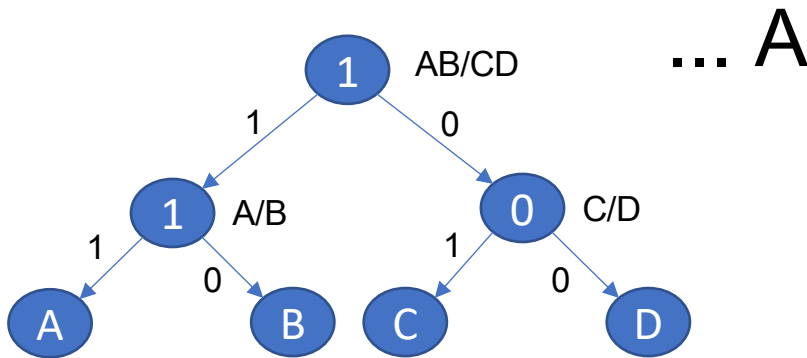서울대학교 천둥 연구실

# Pseudo LRU (cont'd)

- A binary decision tree
  - 2-way set associative cache: 1 bit
  - 4-way set associative cache: $(2^3 - 1) - 4 = 3$ bits
  - N-way set associative cache: $(2^{\log_2 N + 1} - 1) - N$ bits
- The difference between pseudo LRU and true LRU is statistically small
- Each bit represents the left or right child in the binary decision tree
  - 1: the left side has been referenced more recently than the right side
  - 0: vice versa
- A write cycle to update the pseudo-LRU bits on a hit
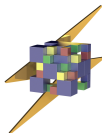- A read cycle for the pseudo-LRU bits during a line replacement

| access | next state |
|--------|-----------|
| A | 11_ |
| B | 10_ |
| C | 0_1 |
| D | 0_0 |

| state | replace |
|-------|---------|
| 00X | A |
| 01X | B |
| 1X0 | C |
| 1X1 | D |

... A C

| AB/CD | A/B | C/D |
|-------|-----|-----|
| 1 | 1 | 0 |

THUNDER Research Group
Seoul National University
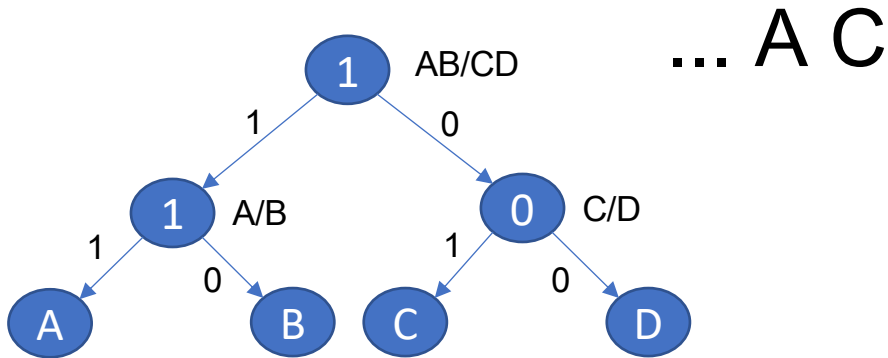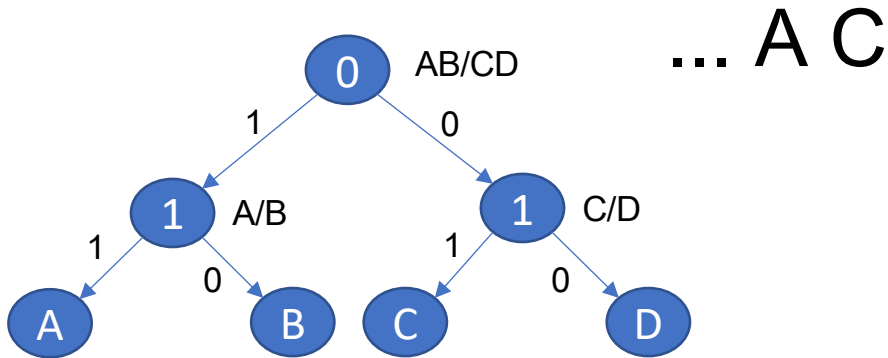서울대학교 천둥 연구실

# Pseudo LRU (cont'd)

- A binary decision tree
  - 2-way set associative cache: 1 bit
  - 4-way set associative cache: $(2^3 - 1) - 4 = 3$ bits
  - N-way set associative cache: $(2^{\log_2 N + 1} - 1) - N$ bits
- The difference between pseudo LRU and true LRU is statistically small
- Each bit represents the left or right child in the binary decision tree
  - 1: the left side has been referenced more recently than the right side
  - 0: vice versa
- A write cycle to update the pseudo-LRU bits on a hit
- A read cycle for the pseudo-LRU bits during a line replacement

| access | next state |
|--------|-----------|
| A | 11_ |
| B | 10_ |
| C | 0_1 |
| D | 0_0 |

| state | replace |
|-------|---------|
| 00X | A |
| 01X | B |
| 1X0 | C |
| 1X1 | D |

... A C

| AB/CD | A/B | C/D |
|-------|-----|-----|
| 0 | 1 | 1 |

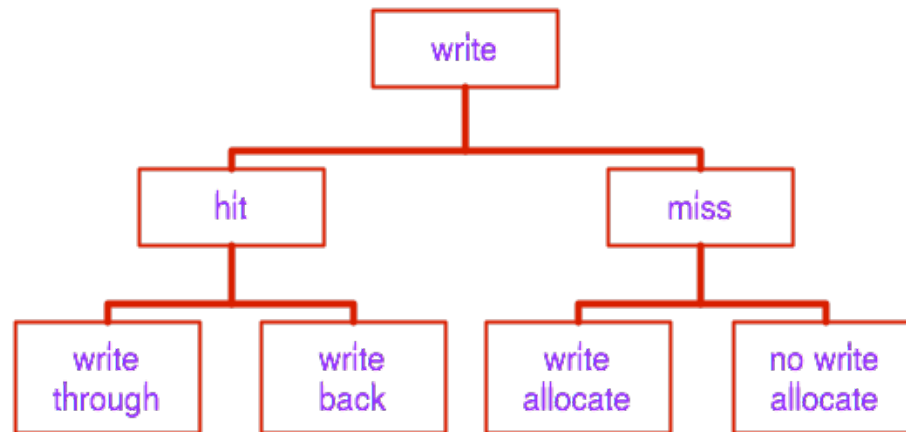THUNDER Research Group
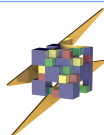Seoul National University
서울대학교 천둥 연구실

# Write Policies

- For reads, the block can be read at the same time that the tag is compared
    - If a miss, just ignore the value read
- For writes, modifying the block cannot begin until the tag is compared
    - Only some part of the entire block is modified
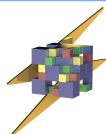
# Write Policies When a Hit

| Write through | Write back |
|---|---|
| • Both the block in the cache and the block in the lower level memory are modified | • Only the block in the cache is modified<br>  ▪ The block is written back to the lower level memory when it is replaced<br>  ▪ A dirty bit is used to reduce the frequency of writing back blocks on replacement |
| • Simpler to implement<br>• Writes are slower than reads<br>• The lower level memory is always consistent with the cache<br>• Every write requires the lower level memory access (need more memory bandwidth)<br>• Read misses never result in writes to the lower level memory | • Harder to implement<br>• Writes and reads are preformed at the same speed<br>• The lower level memory is not always consistent with the cache<br>• Multiple writes within a block require only one write to the lower level memory (need less memory bandwidth)<br>• Read misses may cause writes of dirty blocks to the lower level memory due to replacement |

# Write Policies when a Miss

- Write allocate

  - The block is loaded into the cache on a write miss

- No write allocate

  - The block is modified in the lower level memory and not loaded into the cache

| Write through and write allocate | Write back and write allocate |
|---|---|
| - Subsequent writes to the same block will generate a write to the lower level memory anyway<br>- Bringing the block in the cache is a waste of time | - On a miss it updates the block in the lower level memory and brings the block to the cache<br>- Subsequent writes to the same block will hit in the cache |
| Write through and no write allocate | Write back and no write allocate |
| - Not bringing the block in the cache on a miss saves time | - Subsequent writes to the same block will generate misses |

# Non-Blocking/Lockup-Free Caches

- Most caches can handle only one outstanding request at a time

  - On a miss, the cache must wait for the lower level memory to supply the requested data and until then it is blocked

- A non-blocking cache continues to supply cache hits during a miss

  - Reduce effective miss penalty

  - Another option: supporting multiple outstanding misses

    - A special state need to be maintained for each outstanding miss

      - Miss Status/Information Holding Registers (MSHRs)

# Cache Performance Metrics

- Miss Rate

    - Fraction of memory references not found in the cache

        (misses/references)

- Hit Time

    - Time to deliver a line in the cache to the processor (includes time to

        determine whether the line is in the cache)

- Miss Penalty

    - Additional time required due to the miss