# CSE 321a
# Computer Organization (1)
# تنظيم الحاسبات (1)

3rd year, Computer Engineering
Fall 2016

## Lecture #5

Dr. Hazem Ibrahim Shehata

Dept. of Computer & Systems Engineering

Credits to Dr. Ahmed Abdul-Monem Ahmed for the slides

# Administrivia

- Lecture:
  - Will get back to regular schedule starting next week!
  - Day/Time: Tuesday,  10:15am – 12:45pm
- Assignment #1:
  - Solution will posted on Thursday.
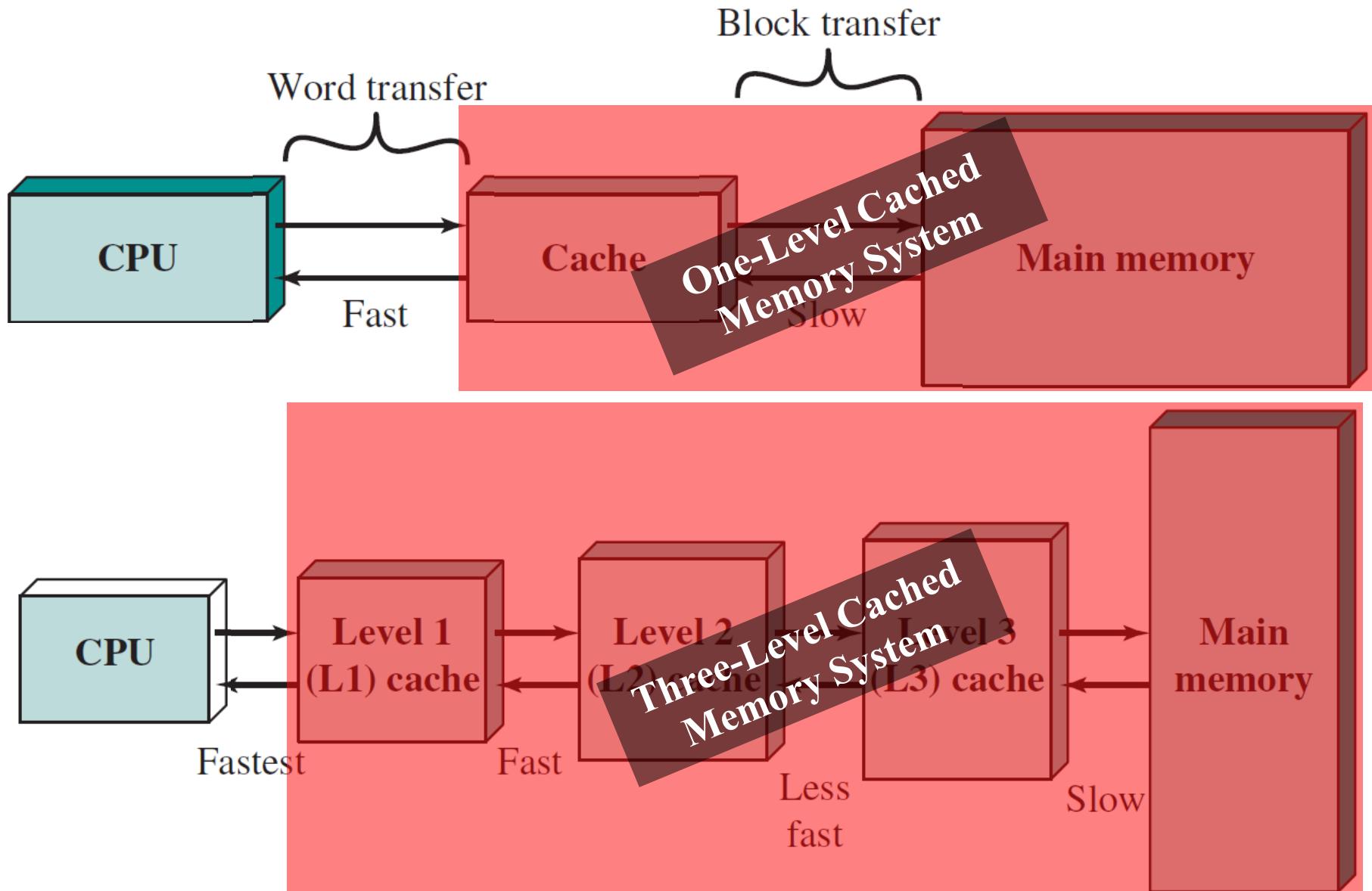
Website: http://hshehata.github.io/courses/zu/cse321a
Office hours: Sunday 12:00pm-1:00pm

# Chapter 4. Cache Memory (*cont.*)

# Cache Memory - Concept

Block transfer

Word transfer

CPU

Cache

**One-Level Cached Memory System**

Main memory

Fast

Slow

CPU

Level 1 (L1) cache

Level 2 (L2) cache

**Three-Level Cached Memory System**

Level 3 (L3) cache

Main memory

Fastest

Fast

Less fast

Slow

# Cache Memory – Design

1. Mapping function
2. Replacement algorithm
3. Read/Write policies
4. Number of caches
5. Addresses
6. Size
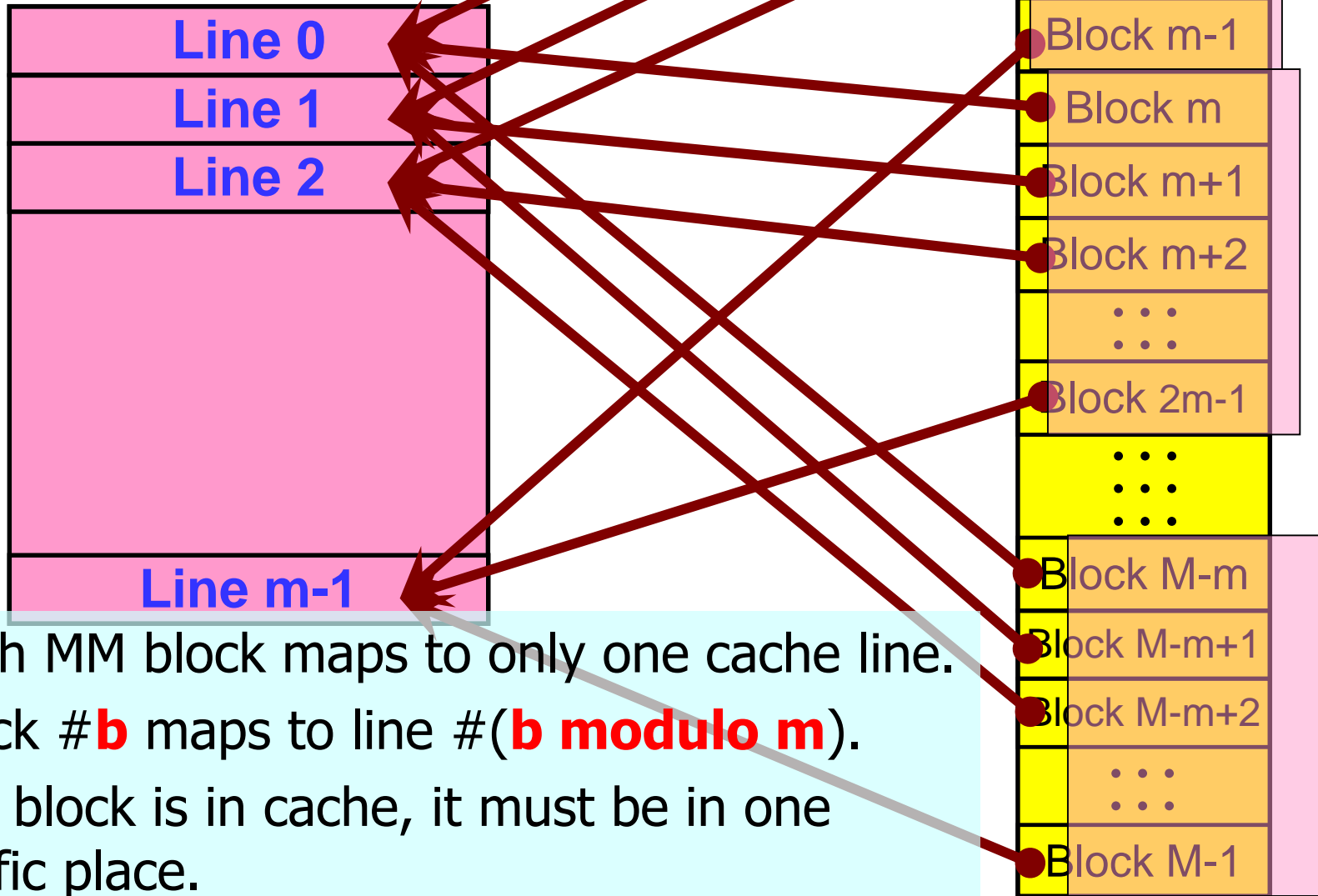7. Block/line size

# Cache Memory – Design

1. Mapping function
2. Replacement algorithm
3. Read/Write policies
4. Number of caches
5. Addresses
6. Size
7. Block/line size

# (I) Direct Mapping

**Main Memory**

**Cache Memory**

| Line 0 |
|---|
| Line 1 |
| Line 2 |
| Line m-1 |

| Block 0 |
|---|
| Block 1 |
| Block 2 |
| . . . |
| Block m-1 |
| Block m |
| Block m+1 |
| Block m+2 |
| . . . |
| Block 2m-1 |
| . . . |
| Block M-m |
| Block M-m+1 |
| Block M-m+2 |
| . . . |
| Block M-1 |

- Each MM block maps to only one cache line.
- Block #**b** maps to line #(**b modulo m**).
- If a block is in cache, it must be in one specific place.

# (I) Direct Mapping
# Cache Line Table

- Number of MM blocks = M
- Number of cache lines = m

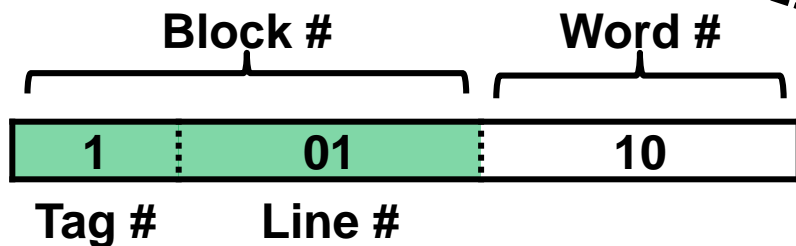| Cache line # (b modulo m) | MM block # (b) |
|---|---|
| 0 | ← 0, m, 2m, 3m, …, M-m |
| 1 | ← 1, m+1, 2m+1, …, M-m+1 |
| 2 | ← 2, m+2, 2m+2, …, M-m+2 |
| …… | ← ………………………………… |
| m-1 | ← m-1, 2m-1, 3m-1, …, M-1 |

# (I) Direct Mapping Tiny Example
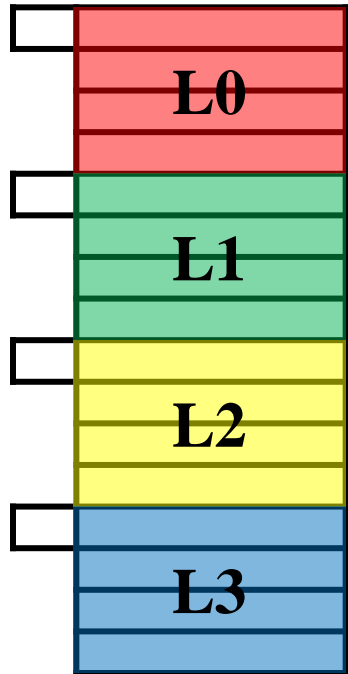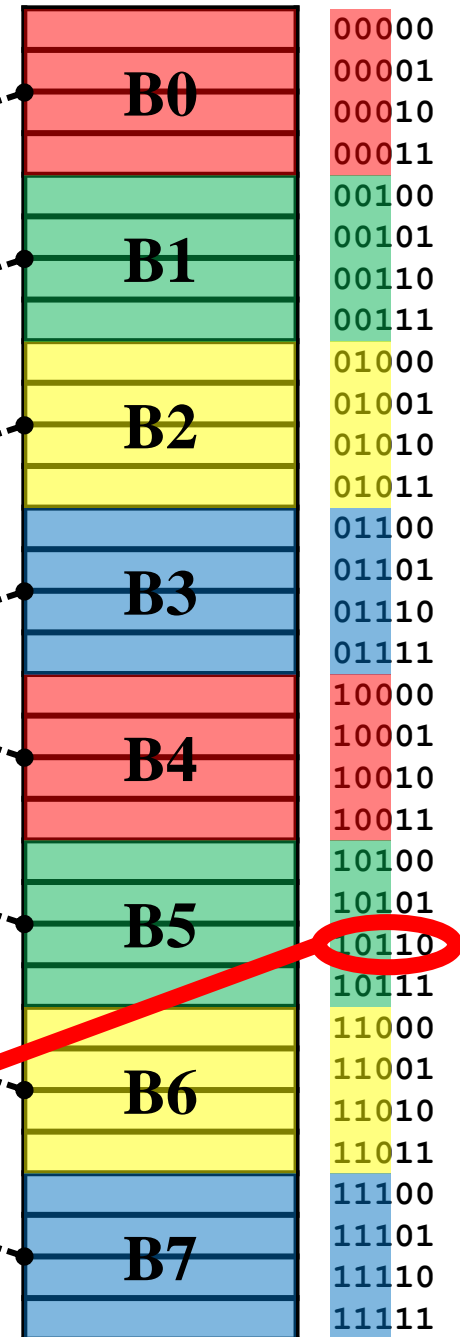
- Main memory:
  - 1 Location ≡ 1 word ≡ 1 byte
  - 32 byte
    - Length of address = $\log_2 32 = 5$
  - 4-byte **blocks**
    - # of blocks (M) = 32 / 4 = 8
- Cache:
  - 16 byte
  - 4-bytes **lines**
    - # of lines (m) = 16 / 4 = 4

# (I) Direct Mapping
## Tiny Example

**Cache Memory**

L0

L1

L2

L3

**Main Memory**

B0

B1

B2

B3

B4

B5

B6

B7

001 → tag=0, line=01

010 → tag=0, line=10

101 → tag=1, line=01

110: tag=1, line=10

00000
00001
00010
00011
00100
00101
00110
00111
01000
01001
01010
01011
01100
01101
01110
01111
10000
10001
10010
10011
10100
10101
10110
10111
11000
11001
11010
11011
11100
11101
11110
11111

**Block #**   **Word #**

| 1 | 01 | 10 |

**Tag #**   **Line #**

# (I) Direct Mapping Address Format

- Memory address is split (based on block size) into:
  1. Least significant **w** bits ➜ identify **word** in block.
     - $w = \log_2$ (# of words in block).
  2. Most significant **s** bits ➜ identify **block** in MM.
     - $s = \log_2$ (# of blocks in MM).
     - The **s** bits are split (based on cache size) into:
       1. Least significant **r** bits ➜ identify cache **line**.
          - $r = \log_2$ (# of lines in cache).
       2. Most significant **s − r** bits ➜ **tag** (identify group).
- In the "tiny example": w=2, s=3, r=2, s-r=1.

| Block # (**s** bits) | | Word # |
|---|---|---|
| Tag # (**s-r** bits) | Line # (**r** bits) | (**w** bits) |

# (I) Direct Mapping
# Big Example

- Main memory:
  - 1 Location ≡ 1 word ≡ 1 byte
  - 16M byte
    - Length of address = $\log_2 16M = 24$
  - 4-byte **blocks**
    - # of blocks (M) = 16M / 4 = 4M
- Cache:
  - 64K byte
  - 4-bytes **lines**
    - # of lines (m) = 64K / 4 = 16K

# (I) Direct Mapping
# Big Example – Address Format

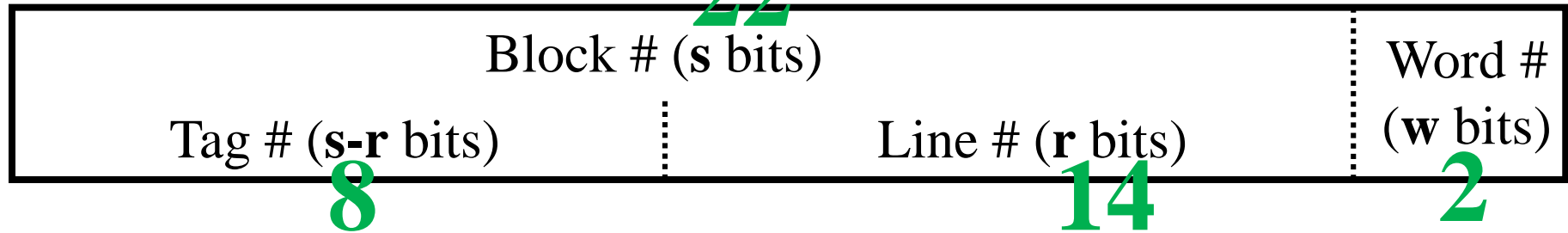| Block # (**s** bits) | | Word # |
|---|---|---|
| **22** | | |
| Tag # (**s-r** bits) | Line # (**r** bits) | (**w** bits) |
| **8** | **14** | **2** |

- 24-bit address
  - ➢ s + w = (s-r) + r + w = 24
- 4-byte block ➔ 4M-block MM
  - ➢ w = $\log_2 4$ = 2 bits
  - ➢ s = 24 − w = 22 bits (Another way: s = $\log_2$ 4M = 22 bits)
- 64K-byte Cache ➔ 16K-line Cache
  - — r = $\log_2$ 16K = 14 bits
  - — s-r = 22 − 14 = 8 bits
- Notes:
  - — No 2 blocks that map to the same line have the same tag field
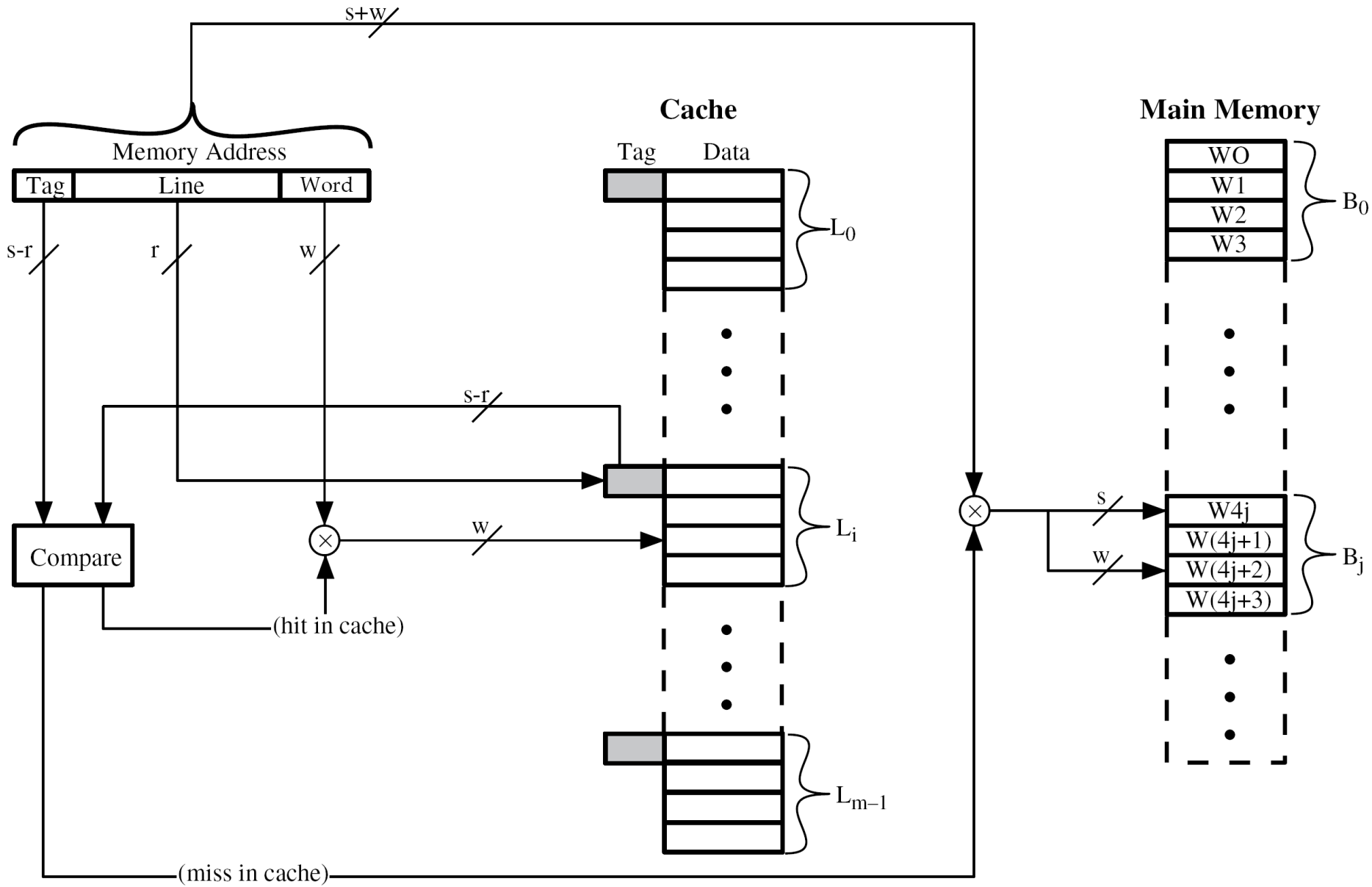  - — Check contents of cache by finding line and checking tag
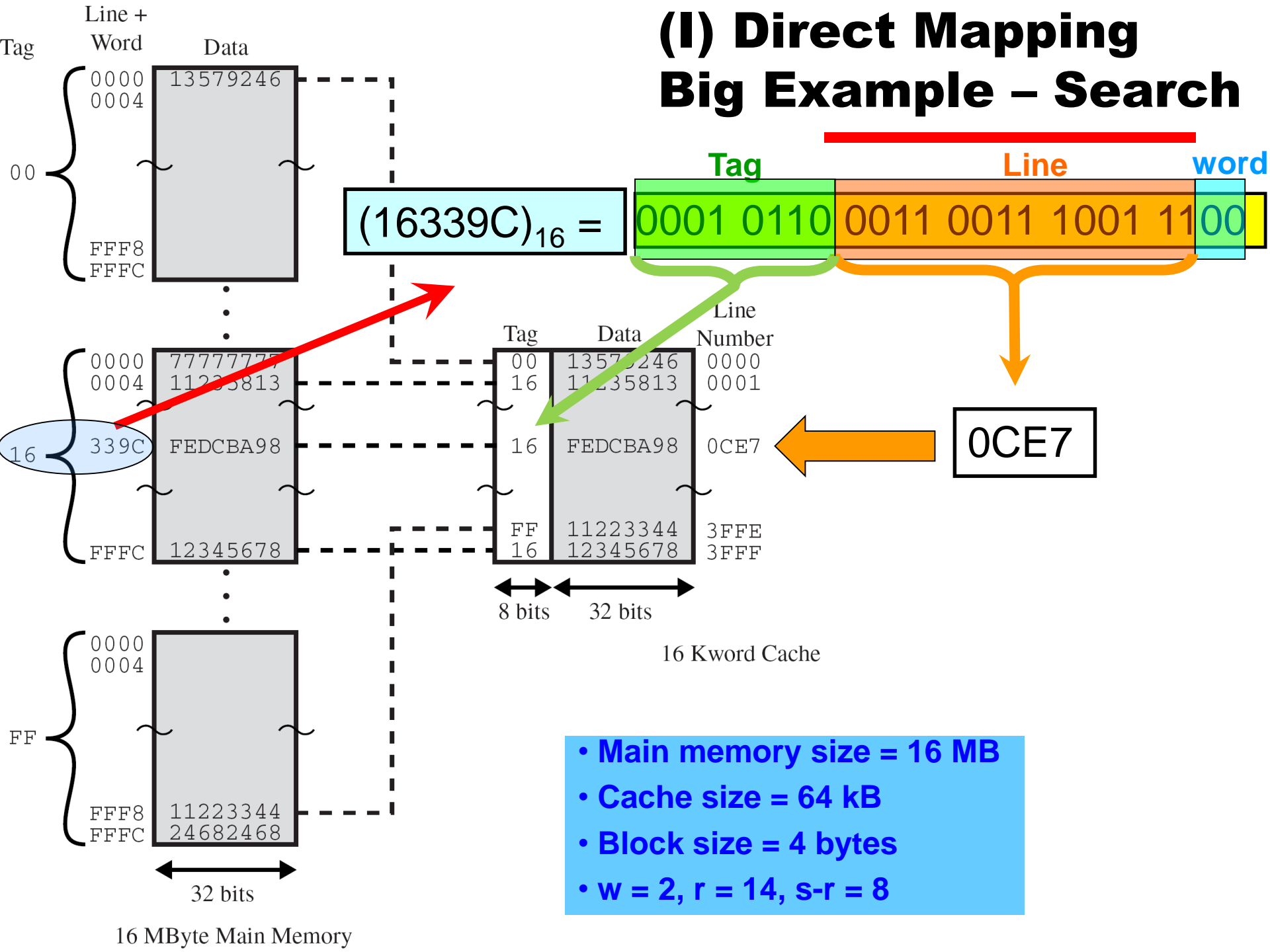
# (I) Direct Mapping
# Address Format Summary

- Address length = (s + w) bits.
- Number of addressable units = $2^{s+w}$ words.
- Block size = line size = $2^w$ words.
  - ➢ w = $\log_2$ (# of words in block).
- Number of blocks in MM = M = $2^{s+w}/2^w = 2^s$.
  - ➢ s = $\log_2$ (# of words in MM / # of words in block)
- Number of lines in cache = m = $2^r$.
  - ➢ r = $\log_2$ (# of words in cache / # of words in line).
- Size of tag = (s − r) bits.
  - ➢ (s-r) = $\log_2$ (# of words in MM / # of words in cache)

# (I) Direct Mapping Cache Organization

# (I) Direct Mapping
# Big Example – Search

Tag | Line + Word | Data

| | |
|---|---|
| 00 | 0000 | 13579246 |
| | 0004 | |
| | FFF8 | |
| | FFFC | |

$(16339C)_{16}$ =

**Tag** | **Line** | **word**

0001 0110 | 0011 0011 1001 11 | 00

Line Number

| | | | |
|---|---|---|---|
| 0000 | 77777777 | | |
| 0004 | 11235813 | | |
| 16 | 339C | FEDCBA98 | |
| FFFC | 12345678 | | |

Tag | Data | Line Number
| | | |
|---|---|---|
| 00 | 13579246 | 0000 |
| 16 | 11235813 | 0001 |
| 16 | FEDCBA98 | 0CE7 |
| FF | 11223344 | 3FFE |
| 16 | 12345678 | 3FFF |

0CE7

8 bits    32 bits

16 Kword Cache

| | | |
|---|---|---|
| FF | 0000 | |
| | 0004 | |
| | FFF8 | 11223344 |
| | FFFC | 24682468 |

32 bits

16 MByte Main Memory

- **Main memory size = 16 MB**
- **Cache size = 64 kB**
- **Block size = 4 bytes**
- **w = 2, r = 14, s-r = 8**

# (I) Direct Mapping
# Pros & Cons

- Simple.

- Inexpensive.

- Fixed location for given block.
  - If a program accesses 2 blocks that map to the same line repeatedly, cache misses are very high.

# (II) (Fully-)Associative Mapping

- A main memory block can load into any line of cache.

- Memory address is interpreted as tag and word.

- Tag uniquely identifies block of memory.

- Every line's tag is examined for a match.

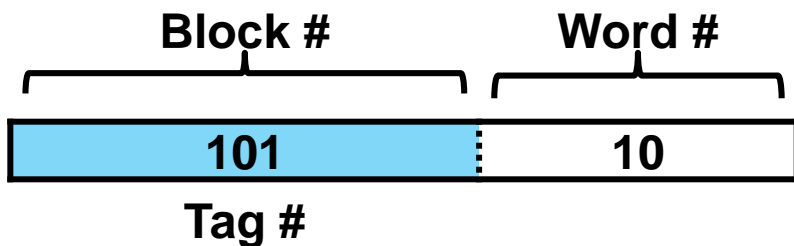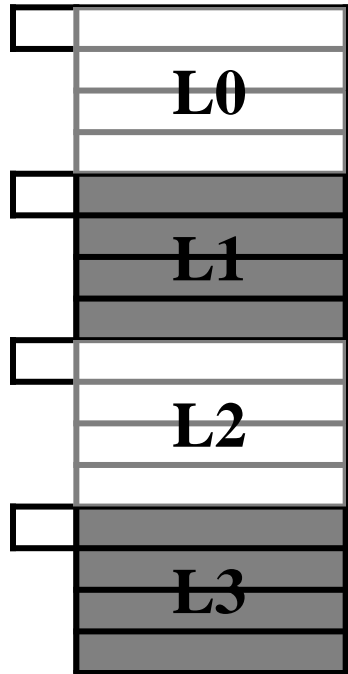- Cache searching gets expensive.

# (II) Associative Mapping Tiny Example

- Main memory:
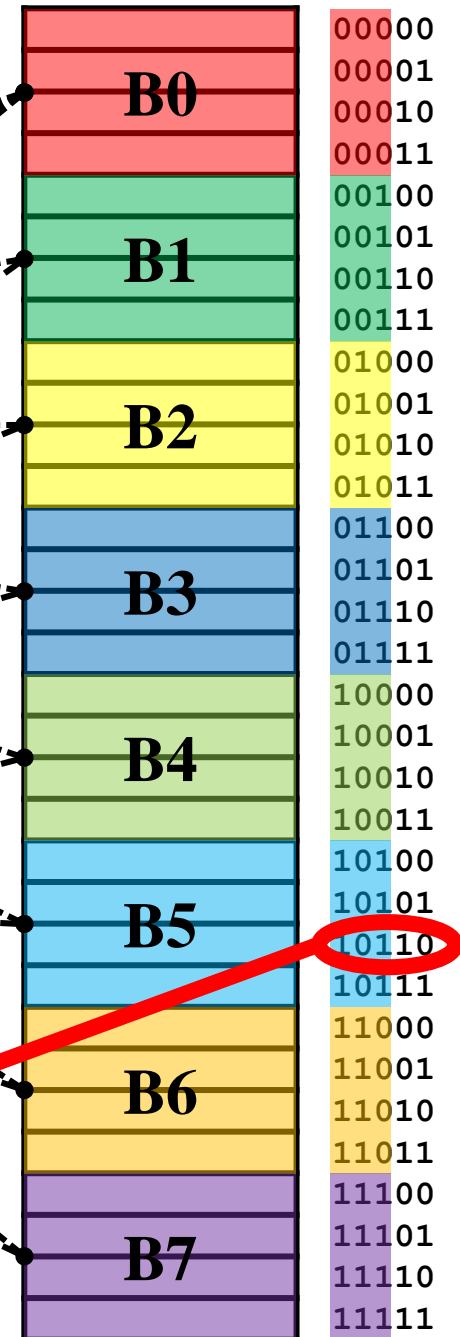    - —1 Location ≡ 1 word ≡ 1 byte
    - —32 byte
        - ➢ Length of address = $\log_2 32 = 5$
    - —4-byte **blocks**
        - ➢ # of blocks (M) = 32 / 4 = 8
- Cache:
    - —16 byte
    - —4-bytes **lines**
        - ➢ # of lines (m) = 16 / 4 = 4

# (II) Associative Mapping
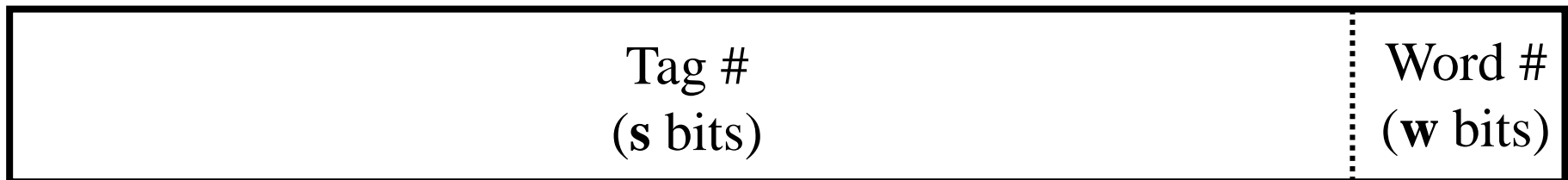# Tiny Example

**Cache Memory**

**Main Memory**



Cache lines: L0, L1, L2, L3

Main Memory blocks: B0, B1, B2, B3, B4, B5, B6, B7

000 → tag=000

101 → tag=101

111 → tag=111

Memory addresses:
00000, 00001, 00010, 00011
00100, 00101, 00110, 00111
01000, 01001, 01010, 01011
01100, 01101, 01110, 01111
10000, 10001, 10010, 10011
10100, 10101, 10110, 10111
11000, 11001, 11010, 11011
11100, 11101, 11110, 11111

**Block #**     **Word #**

| Tag # | |
|---|---|
| 101 | 10 |

# (II) Associative Mapping Address Format

- Memory address is split (based on block size) into:
  1. Least significant **w** bits ➜ identify **word** in block.
     - w = $\log_2$ (# of words in block).
  2. Most significant **s** bits ➜ identify **block** in MM, and used as a **tag**.
     - s = $\log_2$ (# of blocks in MM).
- In the "tiny example": w=2, s=3.

| Tag # (s bits) | Word # (w bits) |
|:---:|:---:|

# (II) Associative Mapping Big Example

- Main memory:
  - 1 Location ≡ 1 word ≡ 1 byte
  - 16M byte
    - Length of address = $\log_2 16M = 24$
  - 4-byte **blocks**
    - # of blocks (M) = 16M / 4 = 4M
- Cache:
  - 64K byte
  - 4-bytes **lines**
    - # of lines (m) = 64K / 4 = 16K

# (II) Associative Mapping
# Big Example – Address Format

| Tag #<br>(**s** bits)<br>22 | Word #<br>(**w** bits)<br>2 |
|---|---|

- 24-bit address
  - ➤ s + w = (s-r) + r + w = 24
- 4-byte block ➤ 4M-block cache
  - ➤ w = $\log_2 4$ = 2 bits
  - ➤ s = 24 − w = 22 bits (Another way: s = $\log_2$ 4M = 22 bits)
- Note:
  - — Each cache line can store a 4-byte block and a 22-bit tag (identifying that block).
  - — To find a block in cache, compare the address tag field (s) against all cache line tags until hit!!
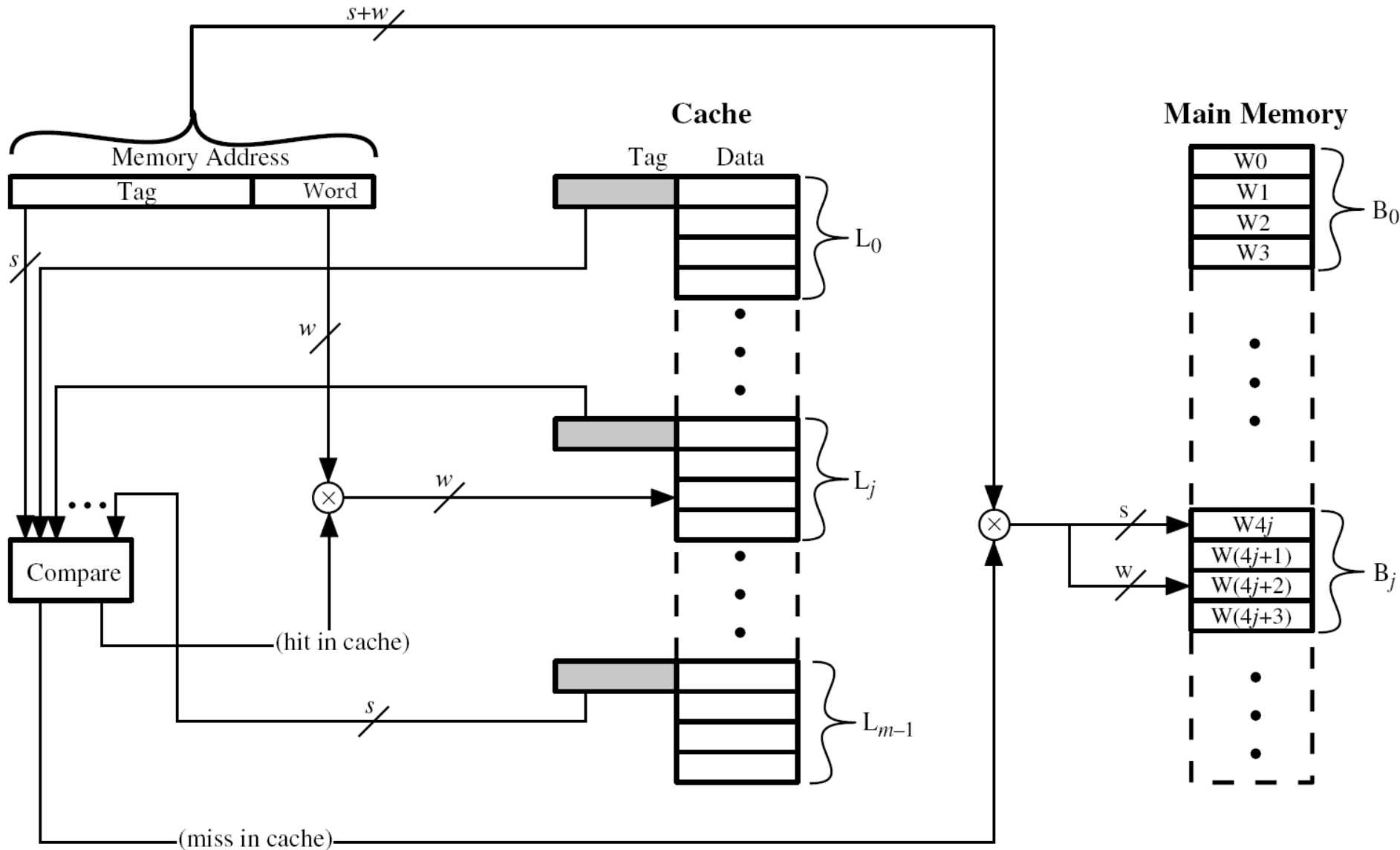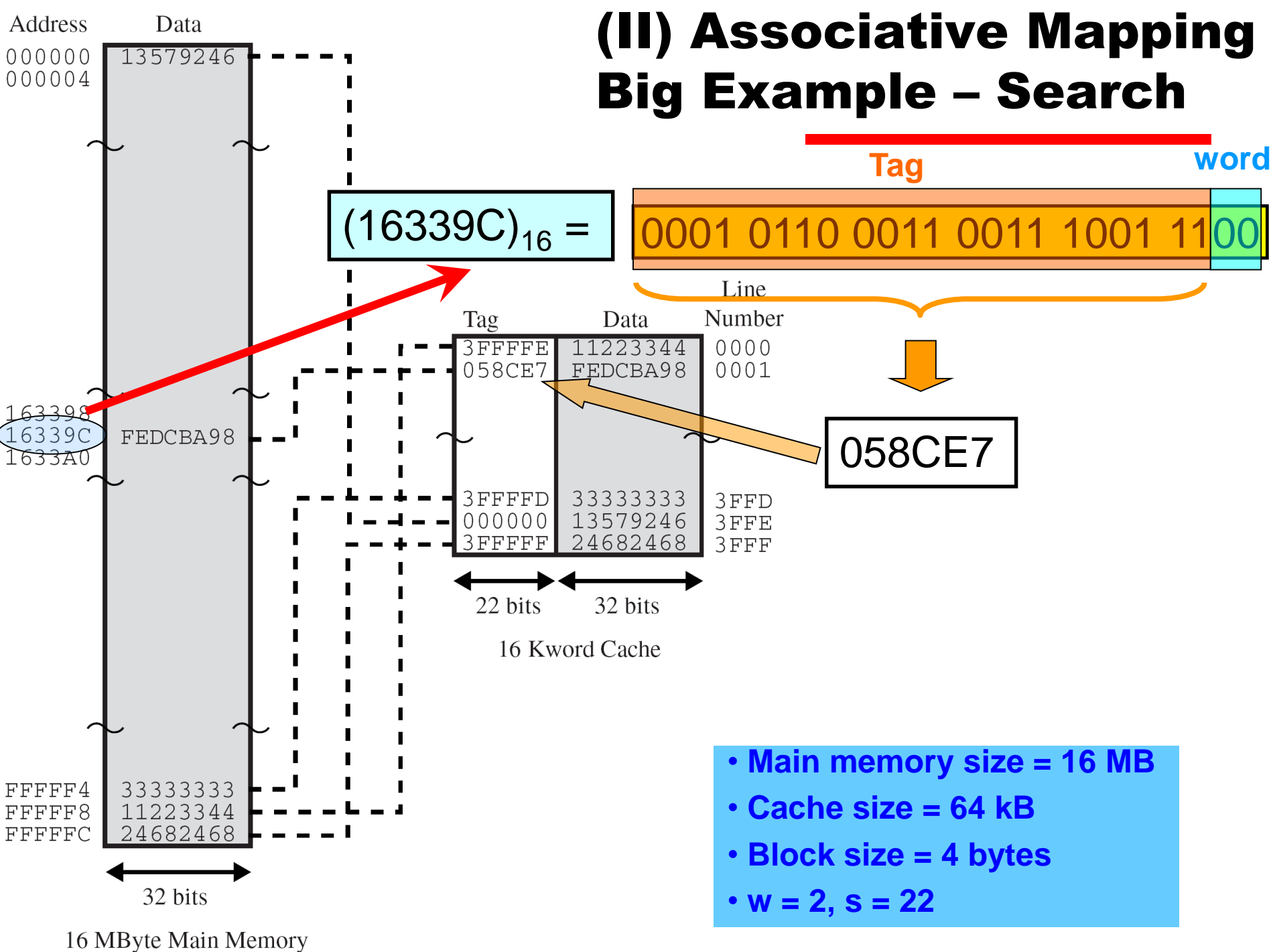
# (II) Associative Mapping
## Address Format Summary

- Address length = $(s + w)$ bits.
- Number of addressable units = $2^{s+w}$ words.
- Block size = line size = $2^w$ words.
  - ➤ $w = \log_2$ (# of words in block).
- Number of blocks in MM = M = $2^{s+w}/2^w = 2^s$.
  - ➤ $s = \log_2$ (# of words in MM / # of words in block)
- Number of lines in cache = m ➔ unknown!
- Size of tag = s.

# (II) Associative Mapping Cache Organization

# (II) Associative Mapping Big Example – Search

**Address** | **Data**

| Address | Data |
|---|---|
| 000000 | 13579246 |
| 000004 | |

163398
16339C | FEDCBA98
1633A0

| | |
|---|---|
| FFFFF4 | 33333333 |
| FFFFF8 | 11223344 |
| FFFFFC | 24682468 |

32 bits

16 MByte Main Memory

$(16339C)_{16} =$

**Tag** | **word**

0001 0110 0011 0011 1001 11|00

Line Number

| Tag | Data | Line Number |
|---|---|---|
| 3FFFFE | 11223344 | 0000 |
| 058CE7 | FEDCBA98 | 0001 |
| 3FFFFD | 33333333 | 3FFD |
| 000000 | 13579246 | 3FFE |
| 3FFFFF | 24682468 | 3FFF |

22 bits     32 bits

16 Kword Cache

058CE7

- **Main memory size = 16 MB**
- **Cache size = 64 kB**
- **Block size = 4 bytes**
- **w = 2, s = 22**

# (III) Set-Associative Mapping
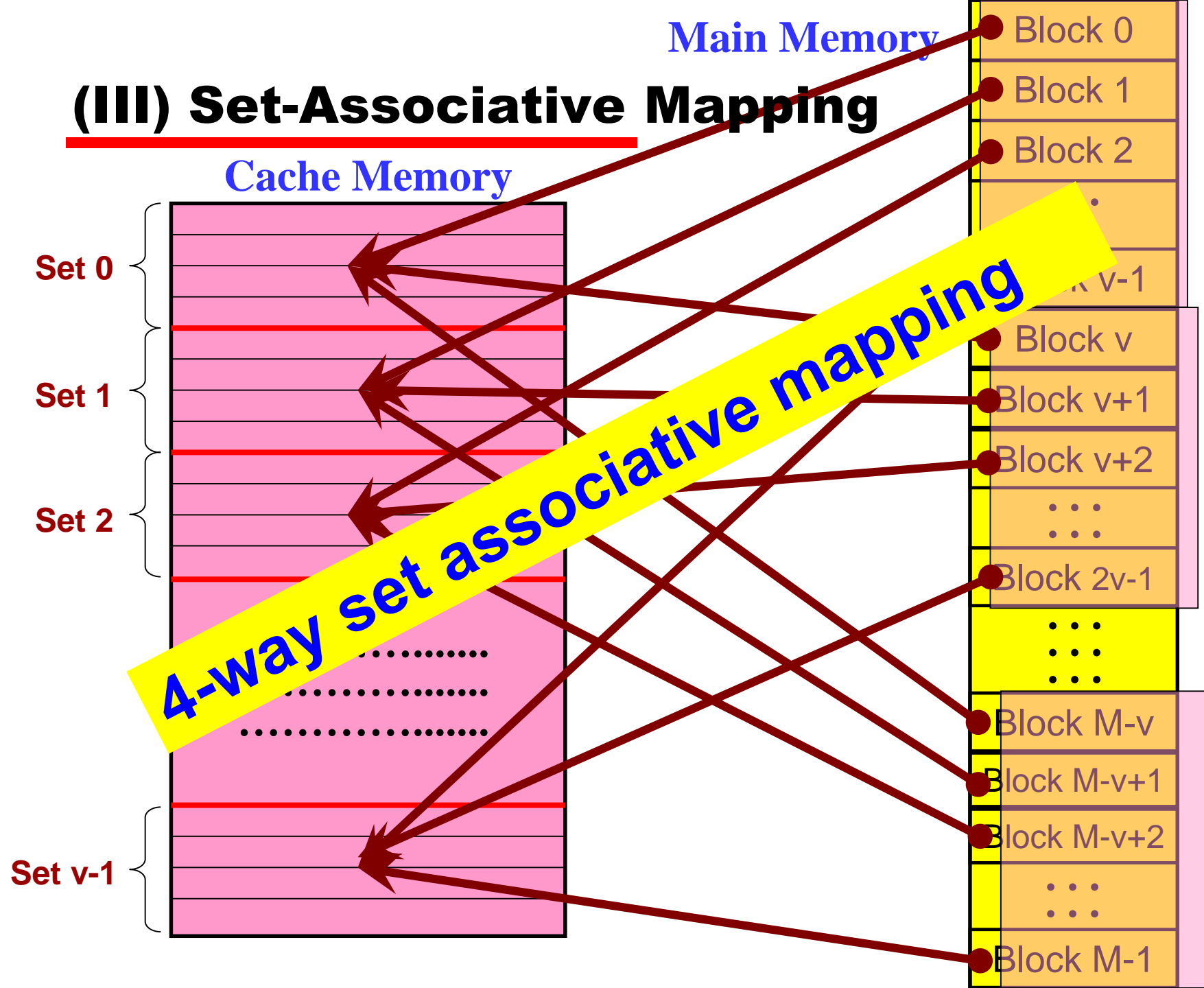
- Cache is divided into a number of **sets** (v) of equal size.

- Each set contains a number of lines (k).
  - ➤ **k-way** set-associative mapping!

- A block b could map to **any line** in a set i if and only if i = b modulo v.

# (III) Set-Associative Mapping

## Cache Memory

**Main Memory**

Set 0

Set 1

Set 2

Set v-1

**4-way set associative mapping**

Block 0

Block 1

Block 2

·

k v-1

Block v

Block v+1

Block v+2

· · ·
· · ·

Block 2v-1

· · ·
· · ·
· · ·

Block M-v

Block M-v+1

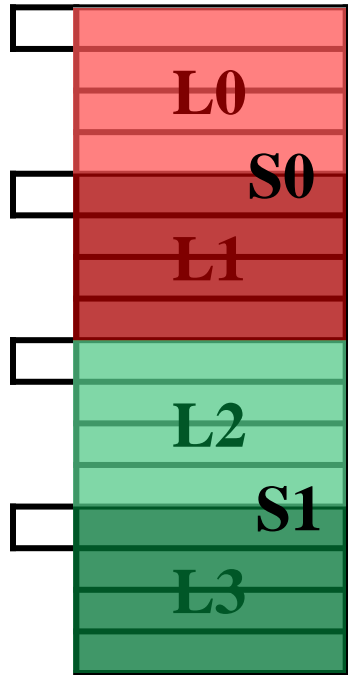Block M-v+2

· · ·
· · ·

Block M-1

# (III) Set-Associative Mapping Tiny Example
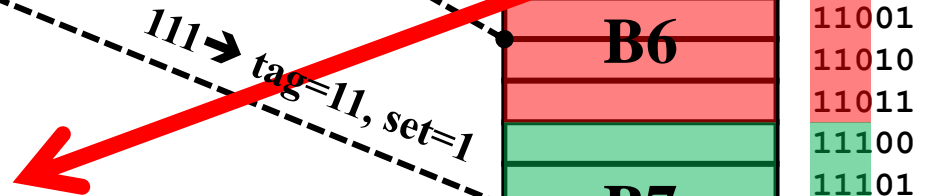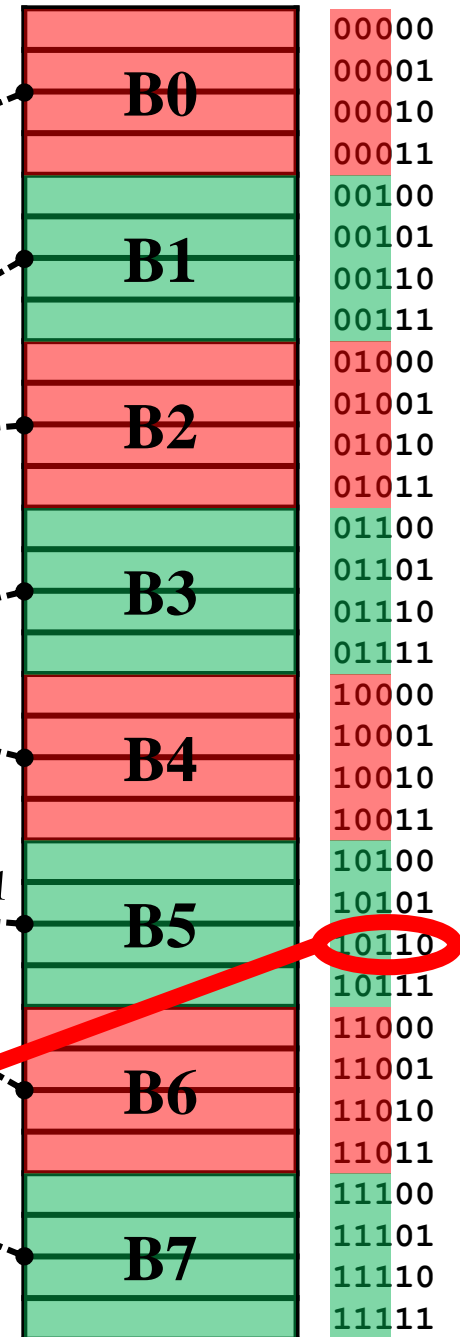
- Main memory:
  - 1 Location ≡ 1 word ≡ 1 byte
  - 32 byte
    - Length of address = $\log_2 32 = 5$
  - 4-byte **blocks**
    - # of blocks (M) = 32 / 4 = 8
- Cache:
  - 16 byte
  - 4-bytes **lines**
    - # of lines (m) = 16 / 4 = 4
  - 2-line **sets** ➔ k=2 ➔ 2-way set-associative
    - # of sets (v) = 4 / 2 = 2

# (III) Set-Associative Mapping Tiny Example

**Cache Memory**

| L0 | S0 |
| L1 | |
| L2 | S1 |
| L3 | |

**Main Memory**

| Block | Address |
|---|---|
| B0 | 00000, 00001, 00010, 00011 |
| B1 | 00100, 00101, 00110, 00111 |
| B2 | 01000, 01001, 01010, 01011 |
| B3 | 01100, 01101, 01110, 01111 |
| B4 | 10000, 10001, 10010, 10011 |
| B5 | 10100, 10101, 10110, 10111 |
| B6 | 11000, 11001, 11010, 11011 |
| B7 | 11100, 11101, 11110, 11111 |

001 → *tag= 00, set=1*

011 → *tag= 01, set=1*

101 → *tag= 10, set= 1*

111 → *tag=11, set=1*

**Block #** **Word #**

| 10 | 1 | 10 |

**Tag #** **Set #**

# (III) Set-Associative Mapping Address Format

- Memory address is split (based on block size) into:
  1. Least significant **w** bits ➜ identify **word** in block.
     - w = $\log_2$ (# of words in block).
  2. Most significant **s** bits ➜ identify **block** in MM.
     - s = $\log_2$ (# of blocks in MM).
     - The **s** bits are split (based on cache size) into:
       1. Least significant **d** bits ➜ identify cache **set**.
          - d= $\log_2$ (# of sets in cache).
       2. Most significant **s − d** bits ➜ **tag** (identify group).
- In the "tiny example": w=2, s=3, d=1, s-d=2.

| Block # (**s** bits) | | Word # |
|---|---|---|
| Tag # (**s-d** bits) | Set # (**d** bits) | (**w** bits) |

# (III) Set-Associative Mapping Big Example
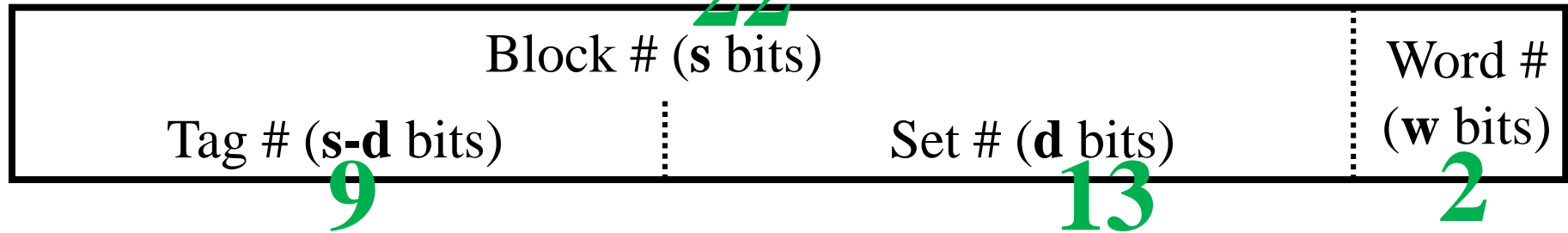
- ## Main memory:
  - 1 Location ≡ 1 word ≡ 1 byte
  - 16M byte
    - Length of address = $\log_2 16M = 24$
  - 4-byte **blocks**
    - # of blocks (M) = 16M / 4 = 4M
- ## Cache:
  - 64K byte
  - 4-bytes **lines**
    - # of lines (m) = 64K / 4 = 16K
  - 2-line **sets** ➔ k=2 ➔ 2-way set-associative
    - # of sets (v) = 16K / 2 = 8K

# (III) Set-Associative Mapping
# Big Example – Address Format

| Block # ($s$ bits) | | Word #<br>($w$ bits) |
|:---:|:---:|:---:|
| **22** | | **2** |
| Tag # ($s$-$d$ bits) | Set # ($d$ bits) | |
| **9** | **13** | |

- 24-bit address
  - ➢ s + w = (s-r) + r + w = 24
- 4-byte block
  - ➢ w = log$_2$ 4 = 2 bits
  - ➢ s = 24 − w = 22 bits
- 64K-byte Cache ➔ 16K-line Cache ➔ 8K-set Cache
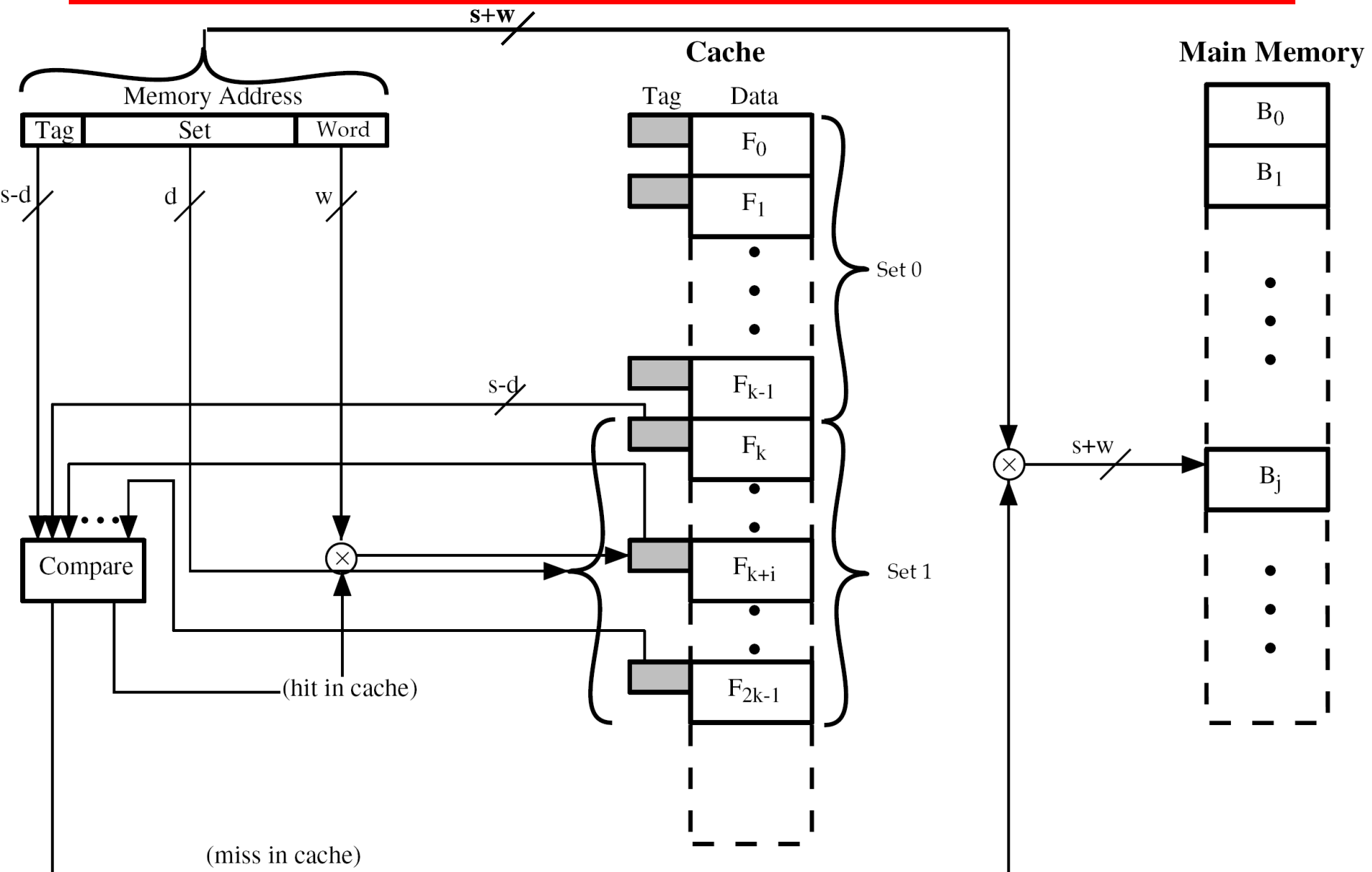  - — d = log$_2$ 8K = 13 bits
  - — s-r = 22 − 13 = 9 bits
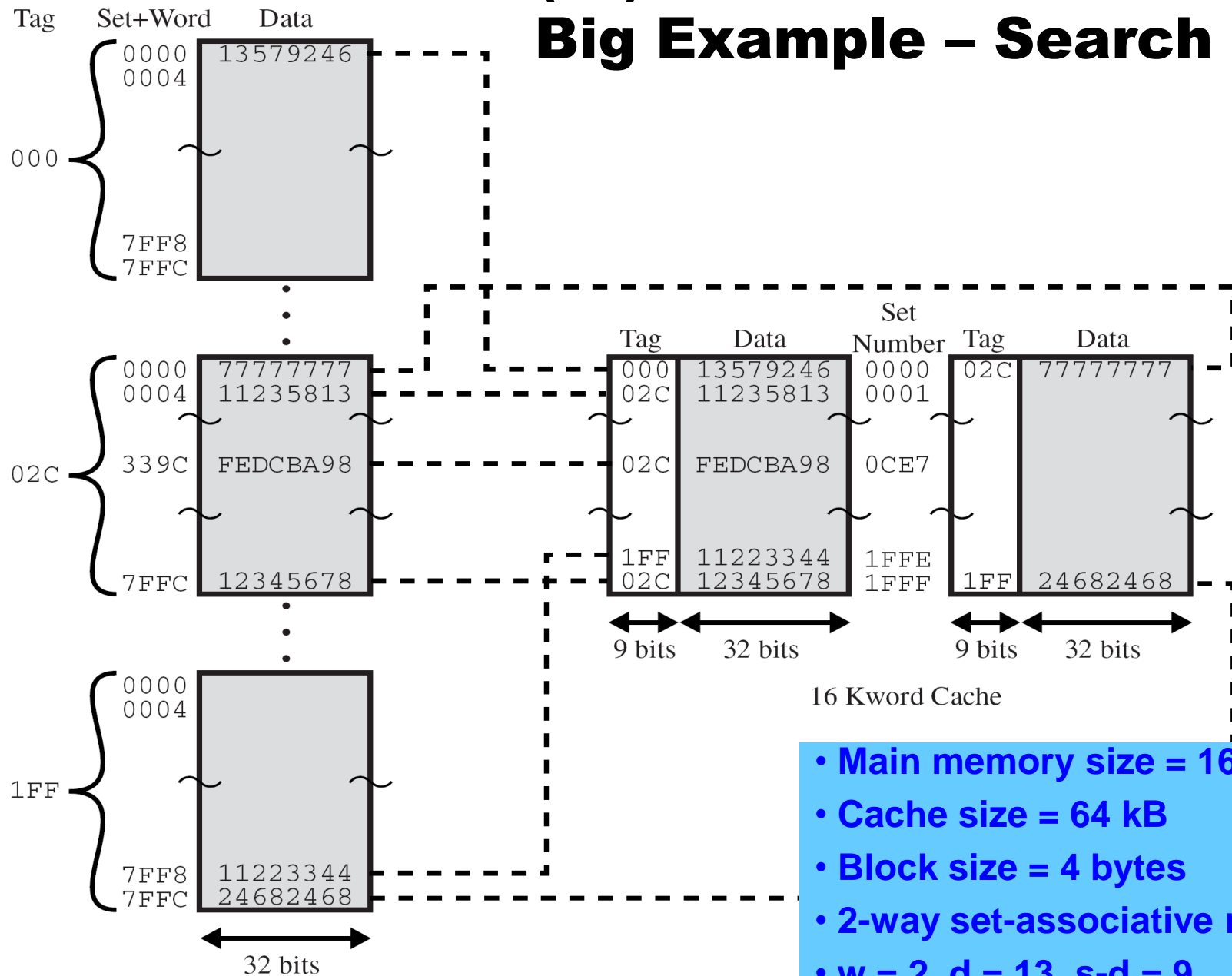
# (III) Set-Associative Mapping Address Format Summary

- Address length = $(s + w)$ bits.
- Number of addressable units = $2^{s+w}$ words.
- Block size = line size = $2^w$ words.
  - ➢ $w = \log_2$ (# of words in block).
- Number of blocks in MM = M = $2^{s+w}/2^w = 2^s$.
  - ➢ $s = \log_2$ (# of words in MM / # of words in block)
- Number of lines in set = k ➔ **k-way set-associative**
- Number of sets = $v = 2^d$.
- Number of lines in cache = m = $k * v = k * 2^d$.
  - ➢ $d = \log_2$ (# of words in cache / # of words in set).
- Size of tag = $(s - d)$ bits.
  - ➢ $(s-d) = \log_2$ (# of blocks in MM / # of sets in cache).

# (III) Set-Associative Mapping Cache Organization (k-way)

# (III) Set-Associative Mapping Big Example – Search

Tag | Set+Word | Data

000

```
0000   13579246
0004
7FF8
7FFC
```

02C

```
0000   77777777
0004   11235813
339C   FEDCBA98
7FFC   12345678
```

1FF

```
0000
0004
7FF8   11223344
7FFC   24682468
```

32 bits

16 MByte Main Memory

|  | Tag | Data | Set Number | Tag | Data |
|---|---|---|---|---|---|
|  | 000 | 13579246 | 0000 | 02C | 77777777 |
|  | 02C | 11235813 | 0001 |  |  |
|  | 02C | FEDCBA98 | 0CE7 |  |  |
|  | 1FF | 11223344 | 1FFE |  |  |
|  | 02C | 12345678 | 1FFF | 1FF | 24682468 |

9 bits   32 bits          9 bits   32 bits

16 Kword Cache

- **Main memory size = 16 MB**
- **Cache size = 64 kB**
- **Block size = 4 bytes**
- **2-way set-associative mapping**
- **w = 2, d = 13, s-d = 9**

# Reading Material

- Stallings, Chapter 4:
    - Pages 123 – 136