# Cache Memory

CPU

1-IF
2-ID
3-EX
4-Mem
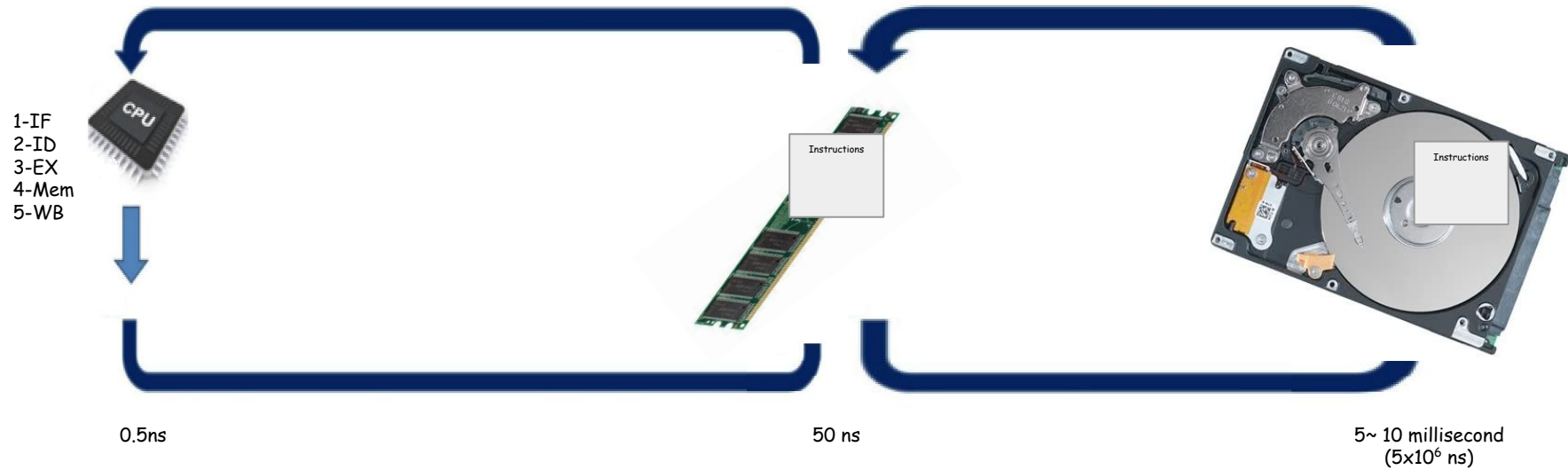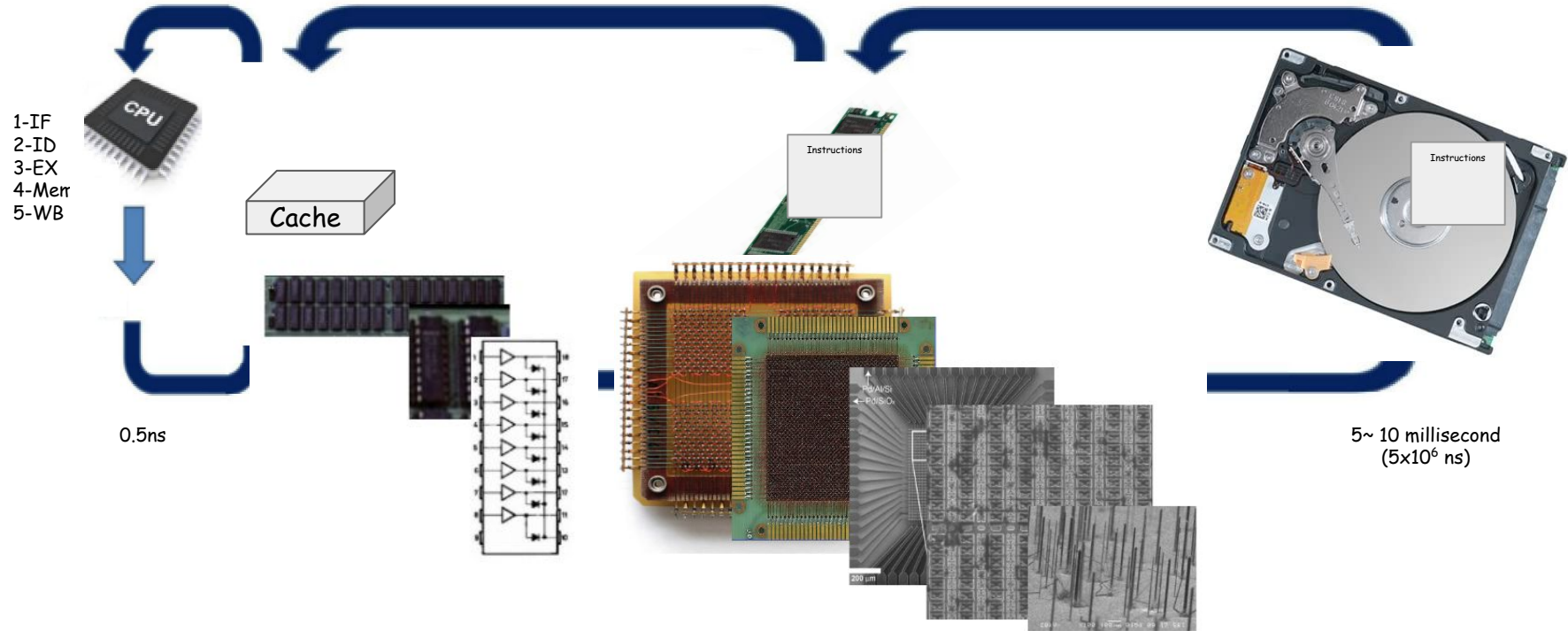5-WB

Instructions
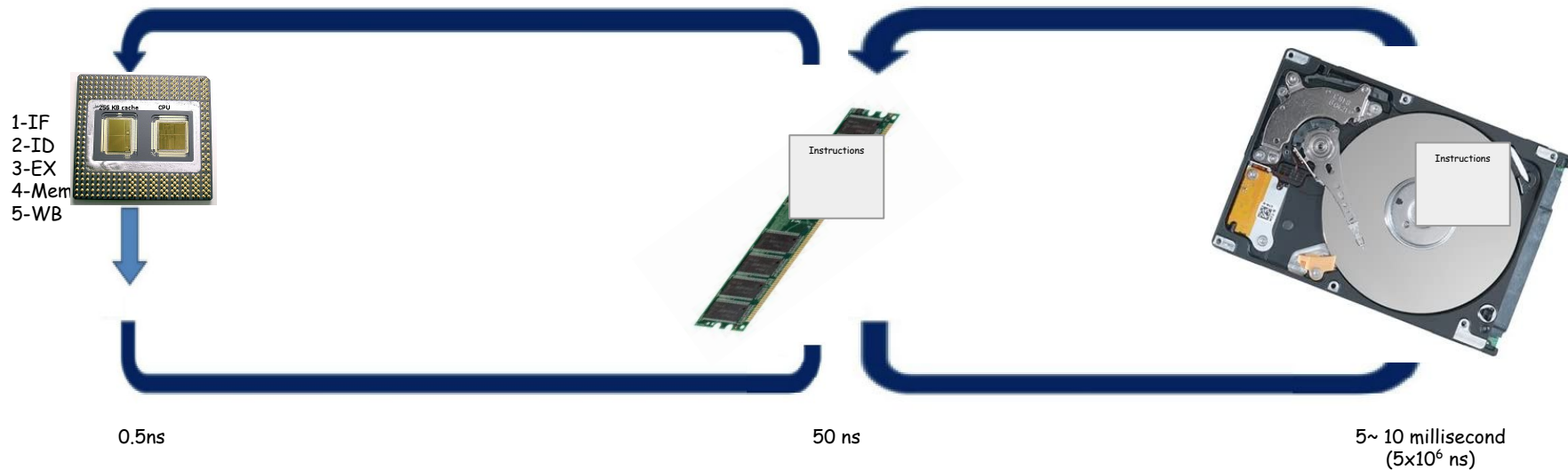
1-IF
2-ID
3-EX
4-Mem
5-WB

0.5ns

Instructions

5~ 10 millisecond
($5 \times 10^6$ ns)

The Von Neumann Bottleneck

1-IF
2-ID
3-EX
4-Mem
5-WB

Instructions

Instructions

0.5ns

50 ns

5~ 10 millisecond
$(5 \times 10^6 \text{ ns})$

1-IF
2-ID
3-EX
4-Mem
5-WB

Cache

Instructions

Instructions

0.5ns

5~ 10 millisecond
$(5 \times 10^6$ ns$)$

1-IF
2-ID
3-EX
4-Mem
5-WB

Instructions

Instructions

0.5ns

50 ns

5~ 10 millisecond
($5 \times 10^6$ ns)

**Spatial Locality:** (Need for adjacent Data)

**Temporal Locality:** (Need to be close to data for some time)
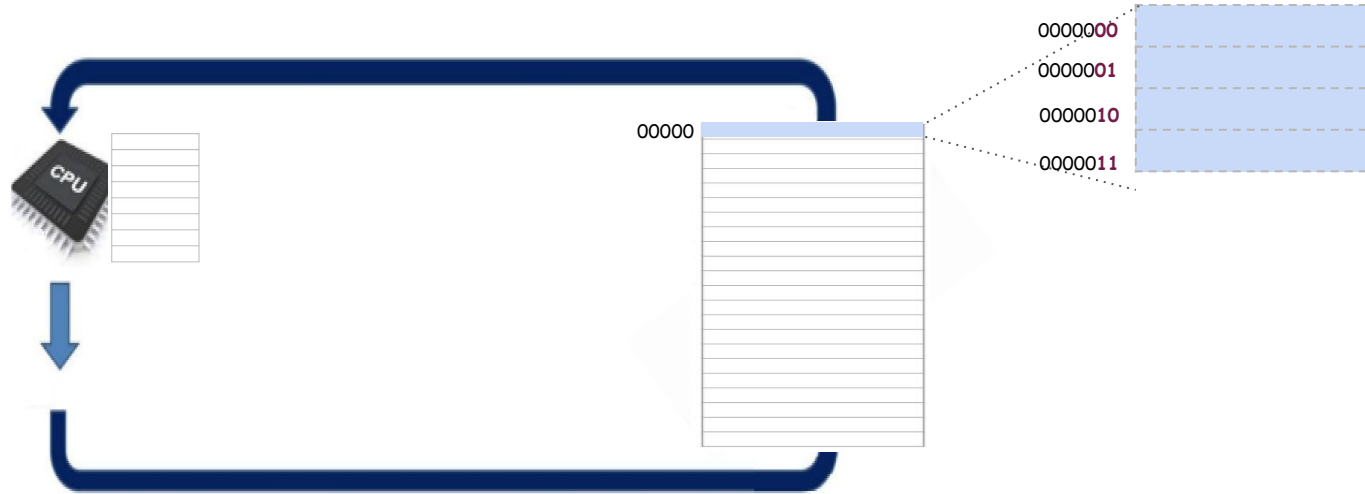
```
int a={3,4,5,6,7,1,2,3,8,3}

total=0

for(int i=0; i<10;i++){

        total+=a[i]
```
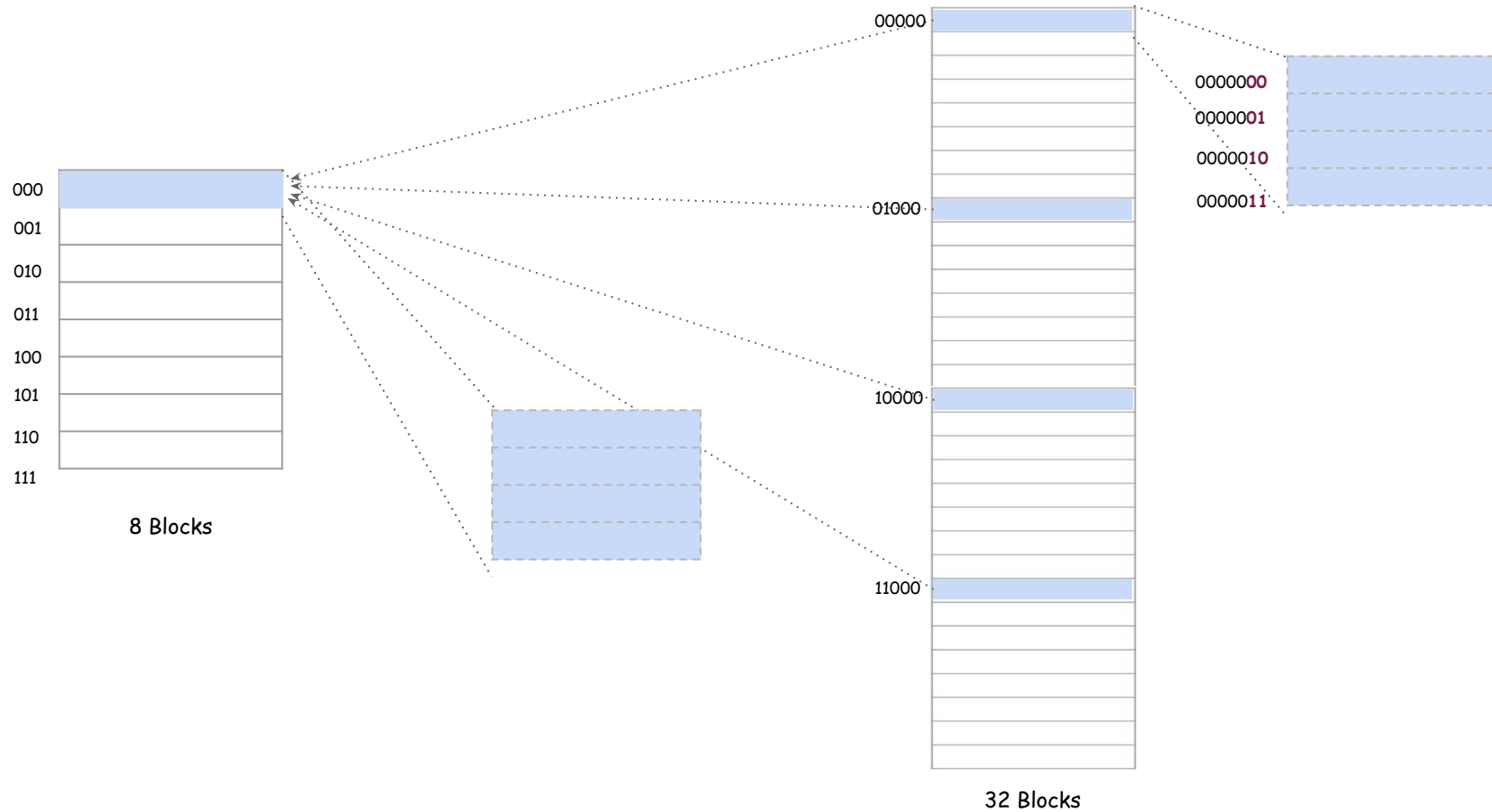
**Cache** is a small amount of memory that's on the CPU itself or right next to it. It can provide the cpu with same of its speed.

- It stores a copy of info. From the main memory.

- CPU asks cache if yes (cache hit) if not (cache miss)

- The greater the cache hits ⇒ the greater the performance

- The greater the cache misses ⇒ the lower the performance

00000

0000000
0000001
0000010
0000011

8 Blocks

32 Blocks

000
001
010
011
100
101
110
111

00000
01000
10000
11000

0000000
0000001
0000010
0000011

8 Blocks

32 Blocks

8 Blocks

32 Blocks

index

tag

| | |
|---|---|
| 000 | _ _ |
| 001 | _ _ |
| 010 | _ _ |
| 011 | _ _ |
| 100 | _ _ |
| 101 | _ _ |
| 110 | _ _ |
| 111 | _ _ |

8 Blocks

00000
00011
01000
01011
10000
10011
11000
11011

32 Blocks

index

tag

validation

| | | |
|---|---|---|
| 000 | _ _ | _ |
| 001 | _ _ | _ |
| 010 | _ _ | _ |
| 011 | _ _ | _ |
| 100 | _ _ | _ |
| 101 | _ _ | _ |
| 110 | _ _ | _ |
| 111 | _ _ | _ |

8 Blocks

00000

00011

01000

01011

10000

10011

11000

11011

32 Blocks

CPU

Request
Address

000
001
010
011
100
101
110
111

8 Blocks

00000
00011
01000
01011
10000
10011
11000
11011

32 Blocks

| Block Address | | Offset |
|---|---|---|
| Tag | index | |

CPU

Request
Address

| 000 | _ _ | _ | |
| 001 | _ _ | _ | |
| 010 | _ _ | _ | |
| 011 | _ _ | _ | |
| 100 | _ _ | _ | |
| 101 | _ _ | _ | |
| 110 | _ _ | _ | |
| 111 | _ _ | _ | |

8 Blocks

| Block Address | | Offset |
| Tag | index | |

Given ByteAddress (0x59), find its location in cache (which block address)?
(Each block has 4 bytes, and the cache has 8 Blocks in total)

00000
00011
01000
01011
10000
10011
11000
11011

32 Blocks

**Exercise:**

Consider a 64 Blocks cache and a block size of 16 bytes and address length is 16 bits. To which Block number does byte address 0x4B0 map?

**Exercise Solution:**

CacheSize is 64 blocks = $2^6$  $\Rightarrow$  so the INDEX is 6 bits

BlockSize is 16 bytes = $2^4$  $\Rightarrow$  so the OFFSET is 4 bits

Address is $(4B0)_{hex}$ = $(0000010010110000)_2$ $\Rightarrow$ the length of the address is 16 $\Rightarrow$ tag is 6 bits

OFFSET

TAG

INDEX

index   tag   validation

| | | |
|---|---|---|
| 000 | | |
| 001 | | |
| 010 | | |
| 011 | | |
| 100 | | |
| 101 | | |
| 110 | | |
| 111 | | |

Storage Area

Cache Memory

The offset is used to determine which byte was requested exactly

00
01
10
11

Bytes within the block

**Exercise:**

How many total bits are required for a direct-mapped cache with 16KB of data and 4-word

blocks, assuming 32-bit address (word size=32 bits)?

**Exercise Solution:**

BlockSize = 4 words = 4 x (32 bits) = 16 bytes = $2^4$ ⇒      the OFFSET is 4

CacheSize = 16kB, Blocksize = 16 bytes ⇒ No. of blocks in cache = 16kB/16 = $2^{10}$ ⇒ the INDEX is 10 bits

Address = TAG + INDEX + OFFSET ⇒          TAG = 32 - (10 + 4) = 18 bits

⇒     Total bits in the cache    = (BlockSizeinBits + validationBit + TAG) x (No. of Blocks)

                               = (128 + 1 + 18 ) x $2^{10}$

                               = 147 kbits

Request Byte Address

| tag | index | offset |

index    tag

validation

000
001
010
011
100
101
110
111

=

Hit        Data

00
01
10
11

Bytes within the block

1-Initial State



Request Byte
Address

| | | |
|---|---|---|
| 000 | __ _ | |
| 001 | __ _ | |
| 010 | __ _ | |
| 011 | __ _ | |
| 100 | __ _ | |
| 101 | __ _ | |
| 110 | __ _ | |
| 111 | __ _ | |

8 Blocks

00000

00011

01000

01011

10000

10011

11000

11011

32 Blocks

2.



Request Byte
Address

0101110

Cache Miss

| | | |
|---|---|---|
| 000 | _ _ | 0 |
| 001 | _ _ | 0 |
| 010 | _ _ | 0 |
| 011 | _ _ | 0 |
| 100 | _ _ | 0 |
| 101 | _ _ | 0 |
| 110 | _ _ | 0 |
| 111 | _ _ | 0 |

8 Blocks

00
01
10
11

00000

00011

01000

01011

10000

10011

11000

11011

32 Blocks

0101100
0101101
0101110
0101111

Some
Data
Here

3.

CPU

Request Byte
Address

0101110

| | | | |
|---|---|---|---|
| 000 | _ _ | 0 | |
| 001 | _ _ | 0 | |
| 010 | _ _ | 0 | |
| 011 | 0 1 | 1 | Mem(01011) |
| 100 | _ _ | 0 | |
| 101 | _ _ | 0 | |
| 110 | _ _ | 0 | |
| 111 | _ _ | 0 | |

8 Blocks

00
01
10
11

Some
Data
Here

00000

00011

01000

01011

10000

10011

11000

11011

0101100
0101101
0101110
0101111

Some
Data
Here

32 Blocks

4.



Request Byte
Address

1000001

**Cache Miss**

| | | |
|---|---|---|
| 000 | _ _ | 0 | |
| 001 | _ _ | 0 | |
| 010 | _ _ | 0 | |
| 011 | 0 1 | 1 | Mem(01011) |
| 100 | _ _ | 0 | |
| 101 | _ _ | 0 | |
| 110 | _ _ | 0 | |
| 111 | _ _ | 0 | |

8 Blocks

| |
|---|
| 00000 |
| |
| 00011 |
| |
| 01000 |
| |
| 01011 |
| |
| 10000 |
| |
| 10011 |
| |
| 11000 |
| |
| 11011 |
| |

32 Blocks

5.



Request Byte
Address

1000001

| 000 | 1 0 | 1 | Mem(10000) |
|-----|-----|---|------------|
| 001 | _ _ | 0 | |
| 010 | _ _ | 0 | |
| 011 | 0 1 | 1 | Mem(01011) |
| 100 | _ _ | 0 | |
| 101 | _ _ | 0 | |
| 110 | _ _ | 0 | |
| 111 | _ _ | 0 | |

8 Blocks

| 00000 | |
|-------|---|
| 00011 | |
| 01000 | |
| 01011 | |
| 10000 | |
| 10011 | |
| 11000 | |
| 11011 | |

0101100
0101101
0101110
0101111

Some
Data
Here

32 Blocks

6.



Request Byte
Address

0101100

Cache Hit 🙂

| | | | |
|---|---|---|---|
| 000 | 1 0 | 1 | Mem(10000) |
| 001 | _ _ | 0 | |
| 010 | _ _ | 0 | |
| 011 | 0 1 | 1 | Mem(01011) |
| 100 | _ _ | 0 | |
| 101 | _ _ | 0 | |
| 110 | _ _ | 0 | |
| 111 | _ _ | 0 | |

8 Blocks

00000

00011

01000

01011

10000

10011

11000

11011

32 Blocks

7.



Request Byte
Address

0100011

Cache Miss

| | | | |
|---|---|---|---|
| 000 | 1 0 | 1 | Mem(10000) |
| 001 | _ _ | 0 | |
| 010 | _ _ | 0 | |
| 011 | 0 1 | 1 | Mem(01011) |
| 100 | _ _ | 0 | |
| 101 | _ _ | 0 | |
| 110 | _ _ | 0 | |
| 111 | _ _ | 0 | |

8 Blocks

00000

00011

01000

01011

10000

10011

11000

11011

0100000
0100001
0100010
0100011

Some
Data
Here

32 Blocks

7.



Request Byte
Address

0100011

replace

| | | | |
|---|---|---|---|
| 000 | 0 1 | 1 | Mem(01000) |
| 001 | _ _ | 0 | |
| 010 | _ _ | 0 | |
| 011 | 0 1 | 1 | Mem(01011) |
| 100 | _ _ | 0 | |
| 101 | _ _ | 0 | |
| 110 | _ _ | 0 | |
| 111 | _ _ | 0 | |

8 Blocks

00000

00011

01000

01011

10000

10011

11000

11011

0100000
0100001
0100010
0100011

Some
Data
Here

32 Blocks

## Cache Performance

MissRate = #CacheMisses / #CacheAccesses

Average Memory
Access Time

$$\text{AMAT} = \text{HitTime} + \text{MissRate} \times \text{MissPenalty}$$

Average Memory Access Time
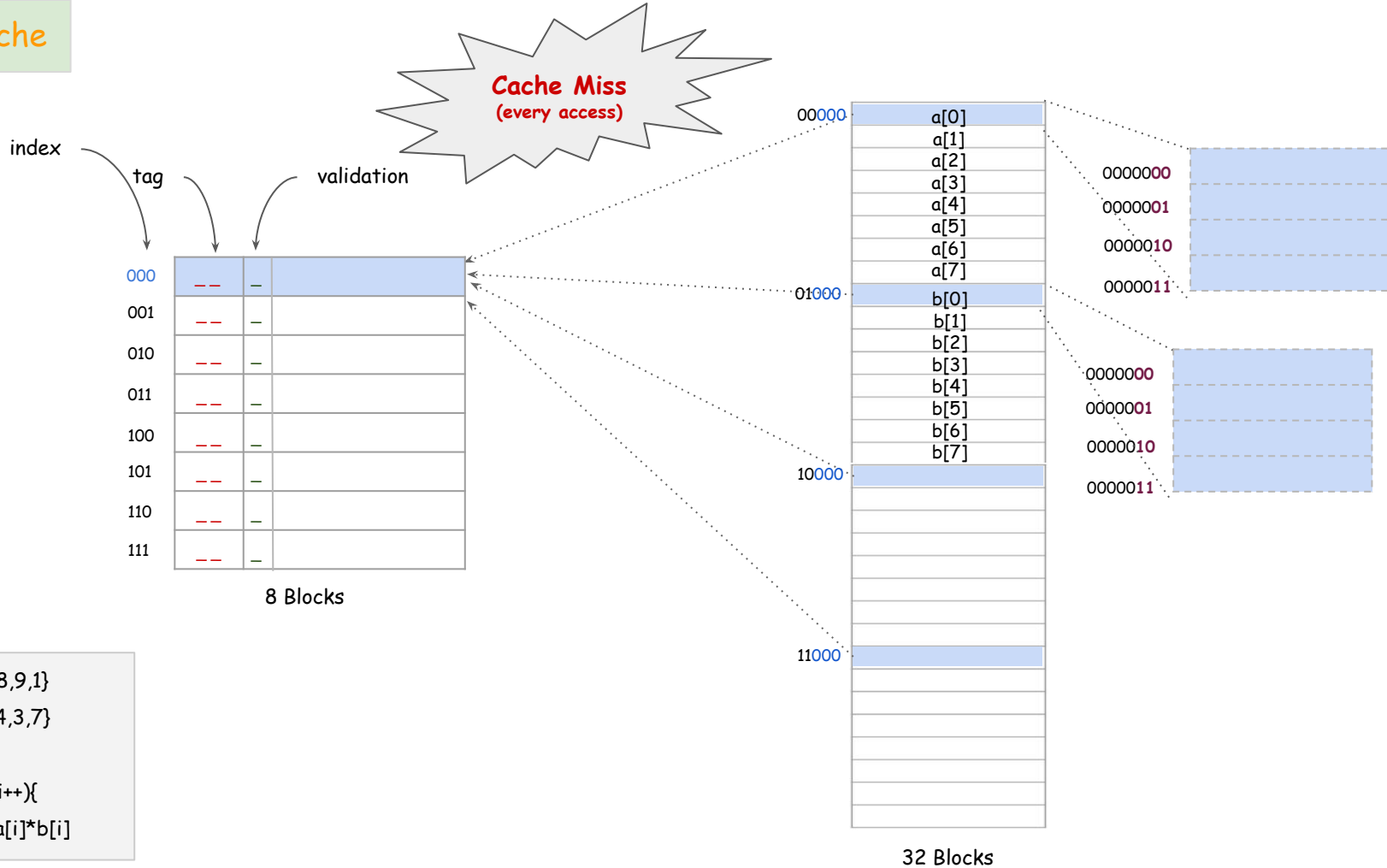
AMAT = HitTime + MissRate x MissPenalty
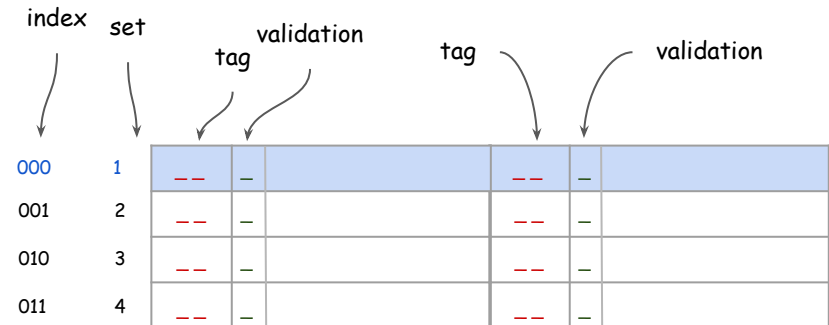
MissRate = #CacheMisses / #CacheAccesses

## Cache Miss Categories (3Cs)

- Compulsory (First Access is always a miss)

- Capacity (program working set is larger than cache capacity)

- Conflict (several blocks are mapped to same block frame)

# Associative Cache

**Cache Miss**
**(every access)**

index

tag

validation

| | | | |
|---|---|---|---|
| 000 | _ _ | _ | |
| 001 | _ _ | _ | |
| 010 | _ _ | _ | |
| 011 | _ _ | _ | |
| 100 | _ _ | _ | |
| 101 | _ _ | _ | |
| 110 | _ _ | _ | |
| 111 | _ _ | _ | |

8 Blocks

int a={3,4,5,6,2,8,9,1}
int b={5,7,1,8,2,4,3,7}
product=0
for(int i=0; i<10;i++){
        product+=a[i]*b[i]

| | |
|---|---|
| 00000 | a[0] |
| | a[1] |
| | a[2] |
| | a[3] |
| | a[4] |
| | a[5] |
| | a[6] |
| | a[7] |
| 01000 | b[0] |
| | b[1] |
| | b[2] |
| | b[3] |
| | b[4] |
| | b[5] |
| | b[6] |
| | b[7] |
| 10000 | |
| | |
| | |
| | |
| | |
| 11000 | |

32 Blocks

0000000
0000001
0000010
0000011

0000000
0000001
0000010
0000011

# Associative Cache

index
set
tag
validation
tag
validation

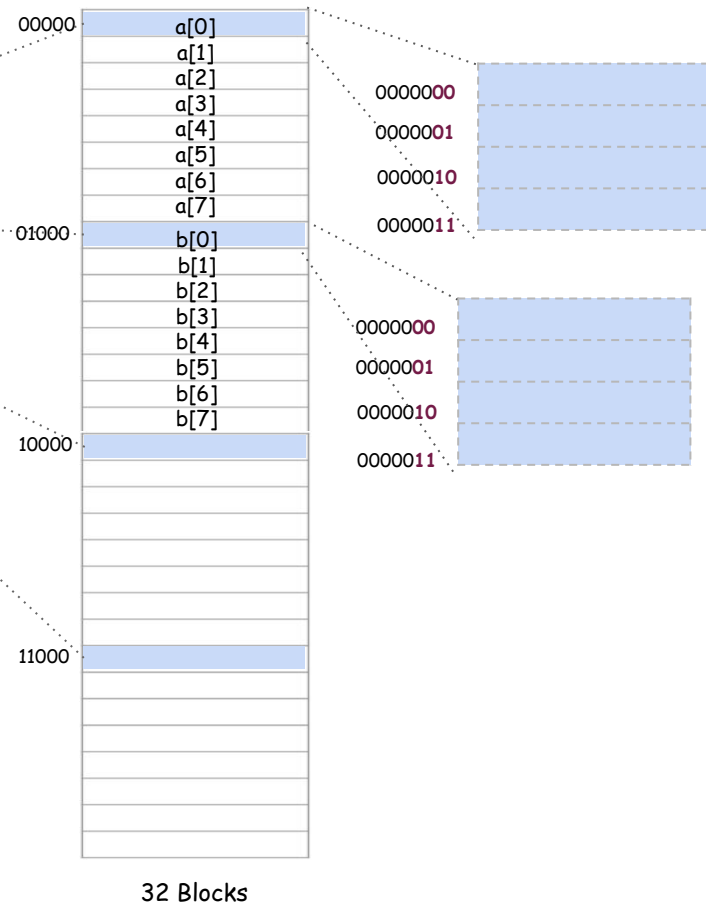| | | | | | | |
|---|---|---|---|---|---|---|
| 000 | 1 | __ | _ | | __ | _ | |
| 001 | 2 | __ | _ | | __ | _ | |
| 010 | 3 | __ | _ | | __ | _ | |
| 011 | 4 | __ | _ | | __ | _ | |

8 Blocks

2-way set associative

```
int a={3,4,5,6,2,8,9,1}
int b={5,7,1,8,2,4,3,7}
product=0
for(int i=0; i<10;i++){
        product+=a[i]*b[i]
```

00000    a[0]
         a[1]
         a[2]
         a[3]
         a[4]
         a[5]
         a[6]
         a[7]
01000    b[0]
         b[1]
         b[2]
         b[3]
         b[4]
         b[5]
         b[6]
         b[7]
10000


11000

32 Blocks

0000000
0000001
0000010
0000011

0000000
0000001
0000010
0000011

# Associative Cache

Request Address

index    validation    validation
     tag    tag

| 00 | | | | | |
| 01 | | | | | |
| 10 | | | | | |
| 11 | | | | | |

8 Blocks

## 2-way set associative

| Block Address | | Offset |
|---|---|---|
| Tag | index | |

Given ByteAddress (0x59), find its location in cache (which block address)?
(Each block has 4 bytes, and the cache is 2-way set associative and has 4 sets in total)

00000
01000
10000
11000

32 Blocks

0000000
0000001
0000010
0000011

0000000
0000001
0000010
0000011

**Exercise:**

Consider a 64 Blocks 2-way associative cache and a block size of 16 bytes. To what Block number does byte address $(0x4B0)_{16}=(0000010010110000)_2$ map to (assuming 16 bits address)?

**Exercise Solution:**

Consider a 64 Blocks 2-way associative cache and a block size of 16 bytes. To what Block number does byte address $(0x4B0)_{16} = (0000010010110000)_2$ map (assuming 16 bits address)?
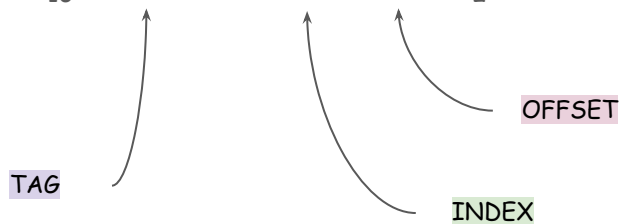
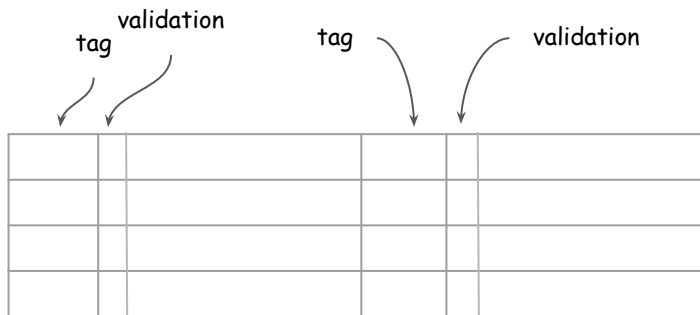CacheSize is 64 blocks (2-ways) $\Rightarrow$     so the is INDEX 5 bits which comes from 64/2 = 32 = $2^5$

BlockSize is 16 bytes = $2^4$         $\Rightarrow$     so the OFFSET is 4 bits

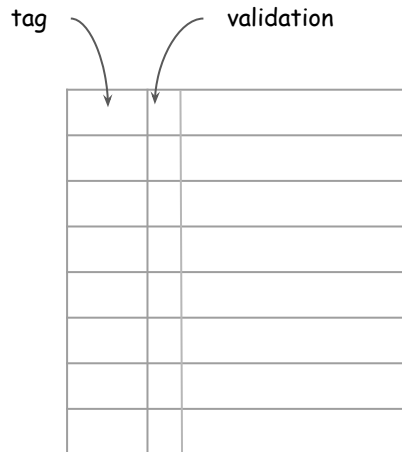Address is $(0x4B0)_{16} = (0000010010110000)_2$

OFFSET

TAG

INDEX

# Associative Cache

## 2-ways set associative

tag  validation  tag  validation

#Sets ≠ #Blocks

## 1-way set associative (direct mapping)

tag  validation

#Sets = #Blocks

## 8-ways set associative (fully associative)

tag  validation  tag  validation  tag  validation  tag  validation  tag  validation  tag  validation  tag  validation  tag  validation

#Sets = 1
(cacheIndex=0)

CacheSize = #Sets x Ways x #BlockSize

# Associative Cache

Request Byte Address

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

| index | tag | v | data | | tag | v | data | | tag | v | data | | tag | v | data |
|-------|-----|---|------|---|-----|---|------|---|-----|---|------|---|-----|---|------|
| 000 | | | | | | | | | | | | | | | |
| 001 | | | | | | | | | | | | | | | |
| 010 | | | | | | | | | | | | | | | |
| 011 | | | | | | | | | | | | | | | |
| 100 | | | | | | | | | | | | | | | |
| 101 | | | | | | | | | | | | | | | |
| 110 | | | | | | | | | | | | | | | |
| 111 | | | | | | | | | | | | | | | |

=   =   =   =

4-1 multiplexer

Hit

Data

**What is the total size of this cache ?**
**(offset=4)**

## Cache Performance

MissRate = #CacheMisses / #CacheAccesses

Average Memory
Access Time

$$\text{AMAT} = \text{HitTime} + \text{MissRate} \times \text{MissPenalty}$$

Which cache organization reduces misses?

What type of misses can this organization reduce (compulsory, capacity, conflict) ?

tag    validation    validation    validation    validation

tag    tag    tag

4-ways set associative
(fully associative)

tag    validation    tag    validation

2-ways set associative

index    tag    validation

1   2   3   4

1-way set associative
(direct mapping)