# Lecture 5:
# Memory Hierarchy
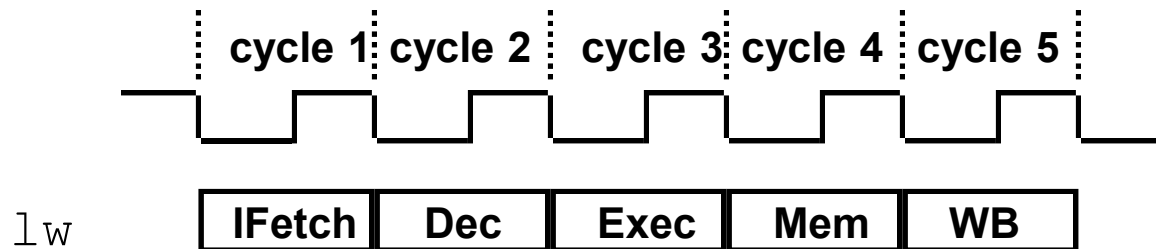
**ECE 411 COMPUTER ORGANIZATION AND DESIGN**

# Watch this

Click the chip

# Review: Instruction Type vs # of Required Cycles

| cycle 1 | cycle 2 | cycle 3 | cycle 4 | cycle 5 |
|---------|---------|---------|---------|---------|

`lw`

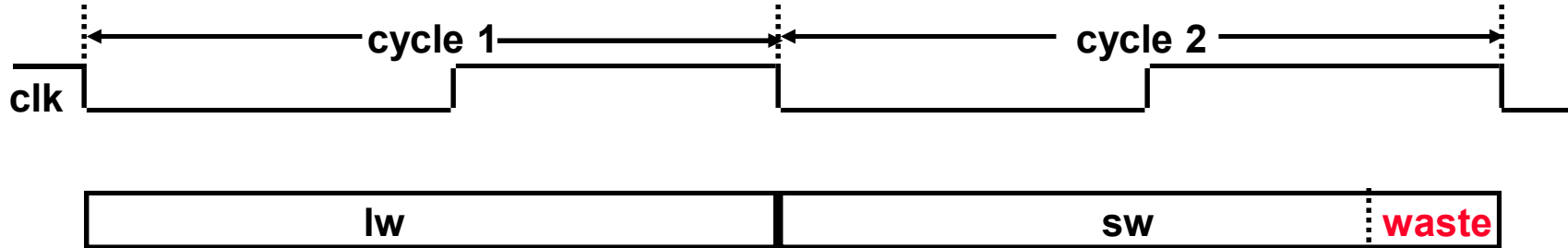| IFetch | Dec | Exec | Mem | WB |
|--------|-----|------|-----|-----|

❑ IFetch: Instruction Fetch and Update PC

❑ Dec: Instruction Decode, Register Read, Sign Extend Offset

❑ Exec: Execute R-type; Calculate Memory Address; Branch Comparison; Branch and Jump Completion

❑ Mem: Memory Read; Memory Write Completion; R-type Completion (RegFile write)

❑ WB:  Memory Read Completion (RegFile write)

*INSTRUCTIONS TAKE FROM 3 - 5 CYCLES!*

# Review: Single Cycle Pros & Cons

❑ uses the clock cycle inefficiently – the clock cycle must be timed to accommodate the <span style="color:red">slowest</span> instruction

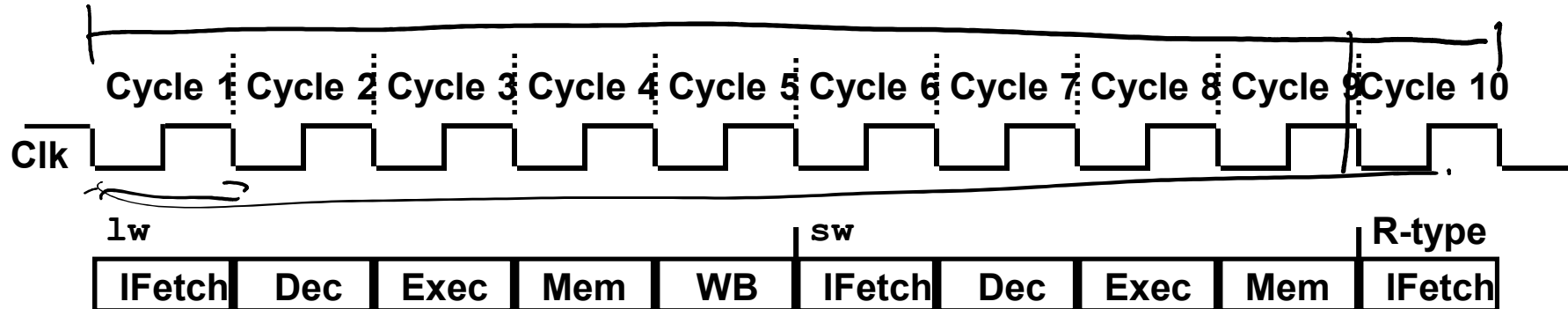  ● especially problematic for more complex instructions like floating point multiply



❑ may be wasteful of area since some functional units (e.g., adders) must be duplicated since they can not be shared during a clock cycle

but

❑ is simple and easy to understand

# Review: Multicycle Pros & Cons

❑ uses the clock cycle efficiently – the clock cycle is timed to accommodate the slowest instruction step

Cycle 1  Cycle 2  Cycle 3  Cycle 4  Cycle 5  Cycle 6  Cycle 7  Cycle 8  Cycle 9  Cycle 10

Clk

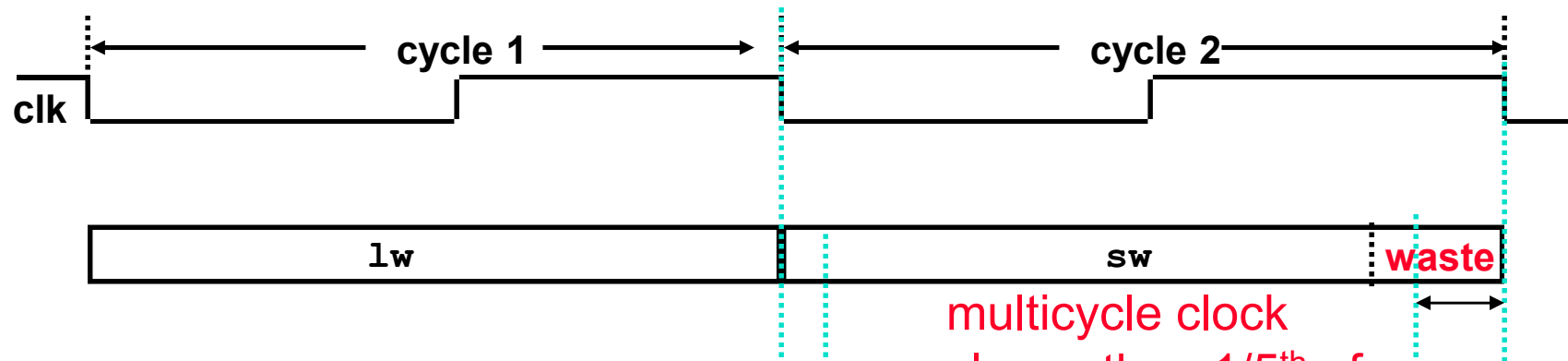| lw | | | | | sw | | | | R-type |
|---|---|---|---|---|---|---|---|---|---|
| IFetch | Dec | Exec | Mem | WB | IFetch | Dec | Exec | Mem | IFetch |

❑ multicycle implementations allow functional units to be used more than once per instruction as long as they are used on different clock cycles
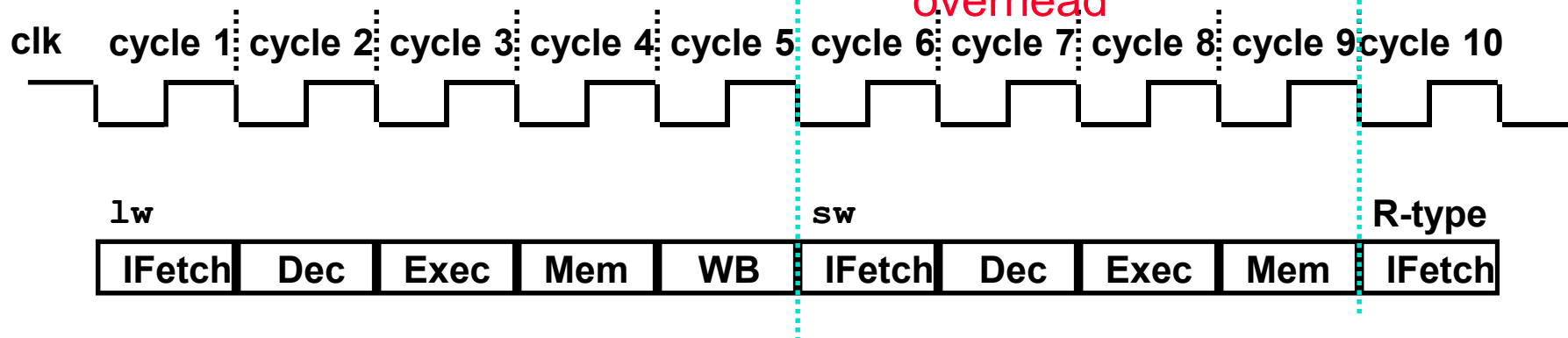
but

❑ requires additional internal state registers, more muxes, and more complicated (FSM) control

ECE 411 COMPUTER ORGANIZATION AND DESIGN

# Review: Single Cycle vs. Multiple Cycle Timing

**single cycle implementation:**

clk

cycle 1 — cycle 2

lw  sw  <span style="color:red">waste</span>

<span style="color:red">multicycle clock slower than 1/5th of single cycle clock due to state register overhead</span>

**multiple cycle implementation:**

clk  cycle 1  cycle 2  cycle 3  cycle 4  cycle 5  cycle 6  cycle 7  cycle 8  cycle 9  cycle 10

| lw | | | | | sw | | | | R-type |
|---|---|---|---|---|---|---|---|---|---|
| IFetch | Dec | Exec | Mem | WB | IFetch | Dec | Exec | Mem | IFetch |

# Review: Will multicycle design be faster?

❑ let's assume $t_{setup}$ + $t_{cq}$ time for registers = 0.1 ns

- single cycle design:
  - clock cycle time = 4.7 + 0.1 = 4.8 ns
  - time/inst = 1 cycle/inst × 4.8 ns/cycle = 4.8 ns/inst
- multicycle design:
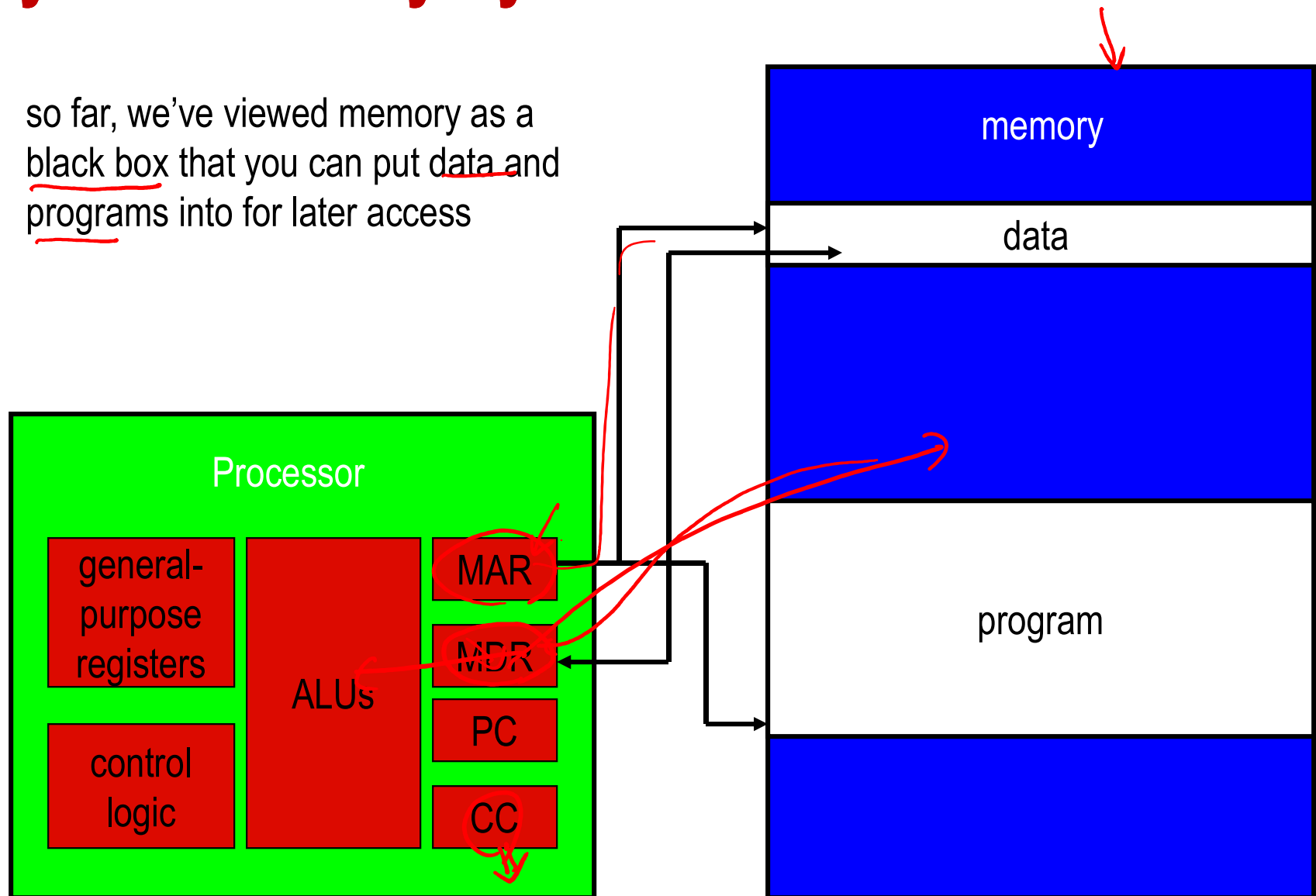  - clock cycle time = 1.0 + 0.1 = 1.1
  - time/inst = CPI × 1.1 ns/cycle (depends on the types or mixture of instructions!)

| | I Fetch | Decode, R-Read | ALU | PC update | D Memory | R-Write | Total (ns) |
|---|---|---|---|---|---|---|---|
| Add | 1 | 1 | .9 | - | - | .8 | 3.7 |
| Load | 1 | 1 | .9 | - | 1 | .8 | 4.7 |
| Store | 1 | 1 | .9 | - | 1 | - | 3.9 |
| beq | 1 | 1 | .9 | .1 | - | - | 3.0 |

# Physical Memory Systems

so far, we've viewed memory as a black box that you can put data and programs into for later access

**memory**

data

**Processor**

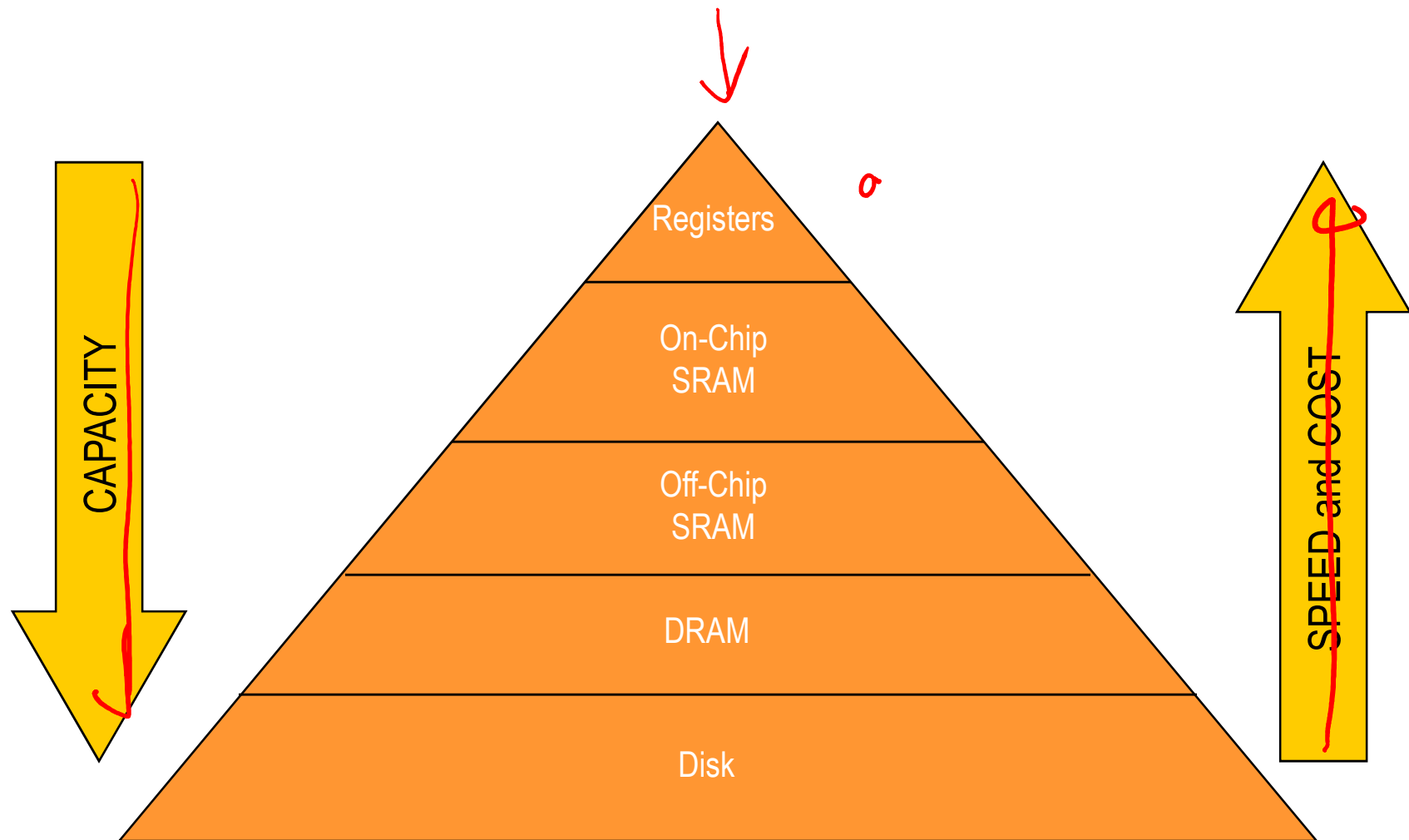| general-purpose registers | ALUs | MAR |
| control logic | | MDR |
| | | PC |
| | | CC |

program

# Types of Memories

- SSD: $0.75 /GB, HD: $0.1-0.2/GB

- DRAM: $20-25/GB

- more on bandwidth later

*Need understand the trade-off among size, speed, cos[t]*

| Type | Size | Latency | Cost/bit |
|------|------|---------|----------|
| Register | < 1KB | < 1ns _(1ns)_ | $$$$ |
| On-chip SRAM | 8KB-6MB _(2ns)_ | < 2ns | $$$ |
| Off-chip SRAM | 1Mb – 16Mb _(long)_ | < 10ns | $$ |
| DRAM | 64MB – 1TB _(100ns)_ | < 100ns | $ |
| Disk (SSD, HD) | 40GB – 1PB _(20m)_ | < 20ms | < $1/GB |

# Memory Hierarchy

CAPACITY

SPEED and COST

Registers

On-Chip
SRAM

Off-Chip
SRAM

DRAM

Disk

# Why Does a Hierarchy Work?

● locality of reference
  ✓ temporal locality
    ○ reference same memory location many times (close together, in time)
  ✓ spatial locality
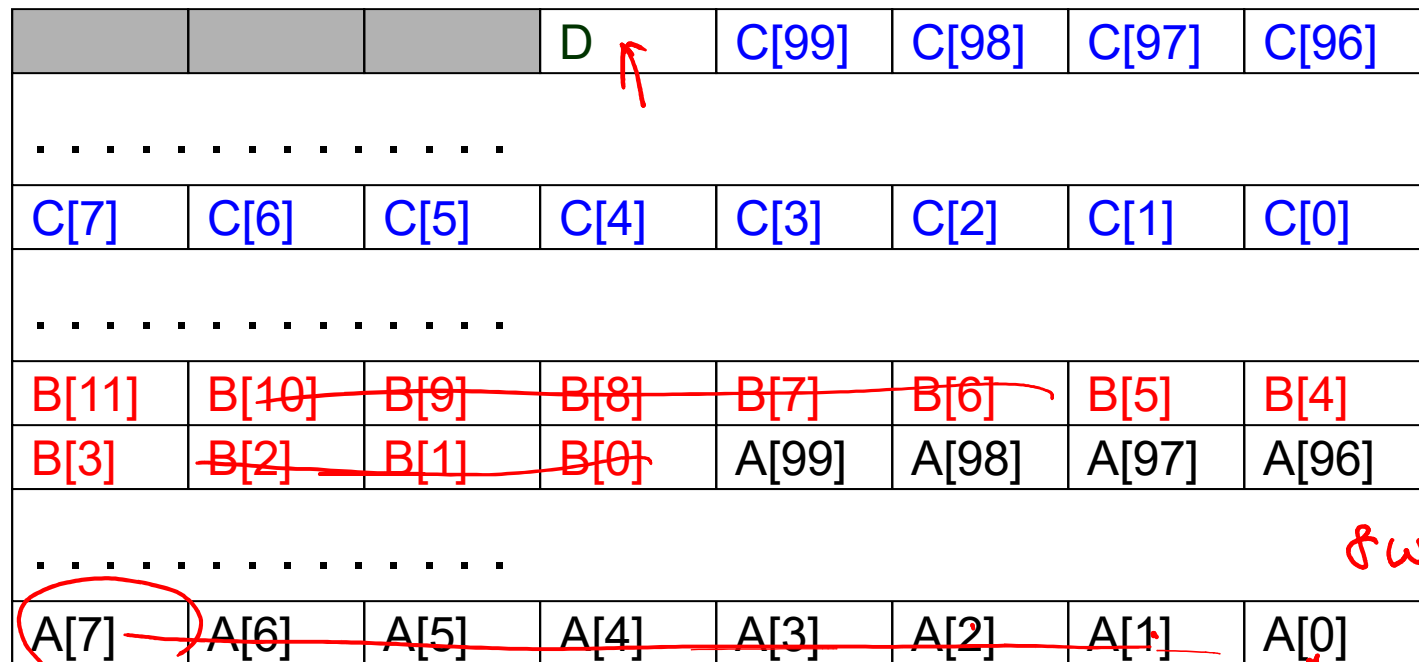    ○ reference near neighbors around the same time

# Example of Locality

```
int A[100], B[100], C[100], D;
for (i=0; i<100; i++) {
        C[i] = A[i] * B[i] + D;
}
```
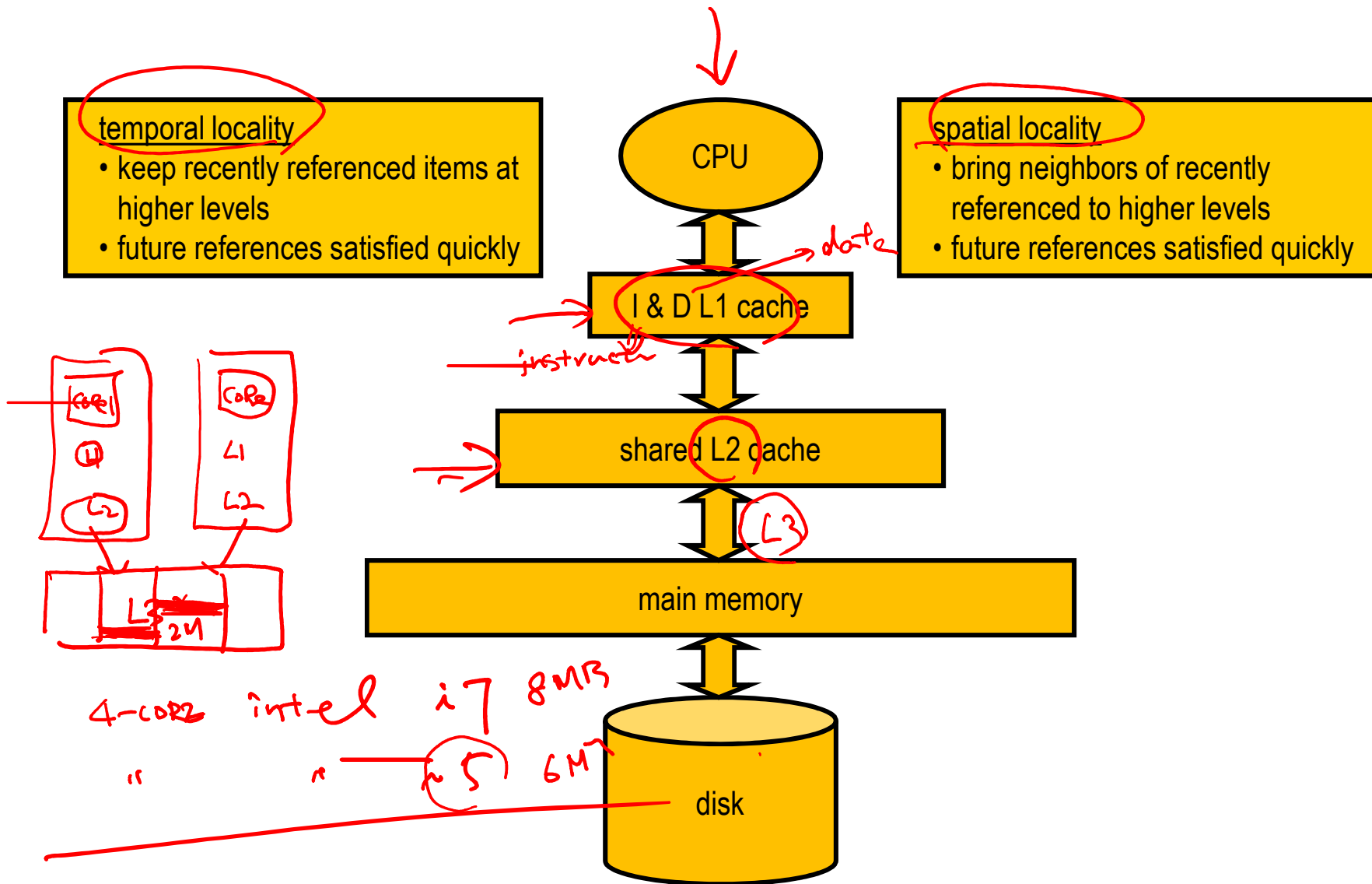
| | | | D | C[99] | C[98] | C[97] | C[96] |
|---|---|---|---|---|---|---|---|
| . . . . . . . . . . . . . . . . | | | | | | | |
| C[7] | C[6] | C[5] | C[4] | C[3] | C[2] | C[1] | C[0] |
| . . . . . . . . . . . . . . . . | | | | | | | |
| B[11] | B[10] | B[9] | B[8] | B[7] | B[6] | B[5] | B[4] |
| B[3] | B[2] | B[1] | B[0] | A[99] | A[98] | A[97] | A[96] |
| . . . . . . . . . . . . . . . . | | | | | | | |
| A[7] | A[6] | A[5] | A[4] | A[3] | A[2] | A[1] | A[0] |

a cache line (one fetch)

8 words

ECE 411 COMPUTER ORGANIZATION AND DESIGN

# Memory Hierarchy

**temporal locality**
- keep recently referenced items at higher levels
- future references satisfied quickly

CPU

data

I & D L1 cache

instruct

**spatial locality**
- bring neighbors of recently referenced to higher levels
- future references satisfied quickly

shared L2 cache

L3

main memory

disk

CORE1    CORE2
L1       L1
L2       L2
L3  24

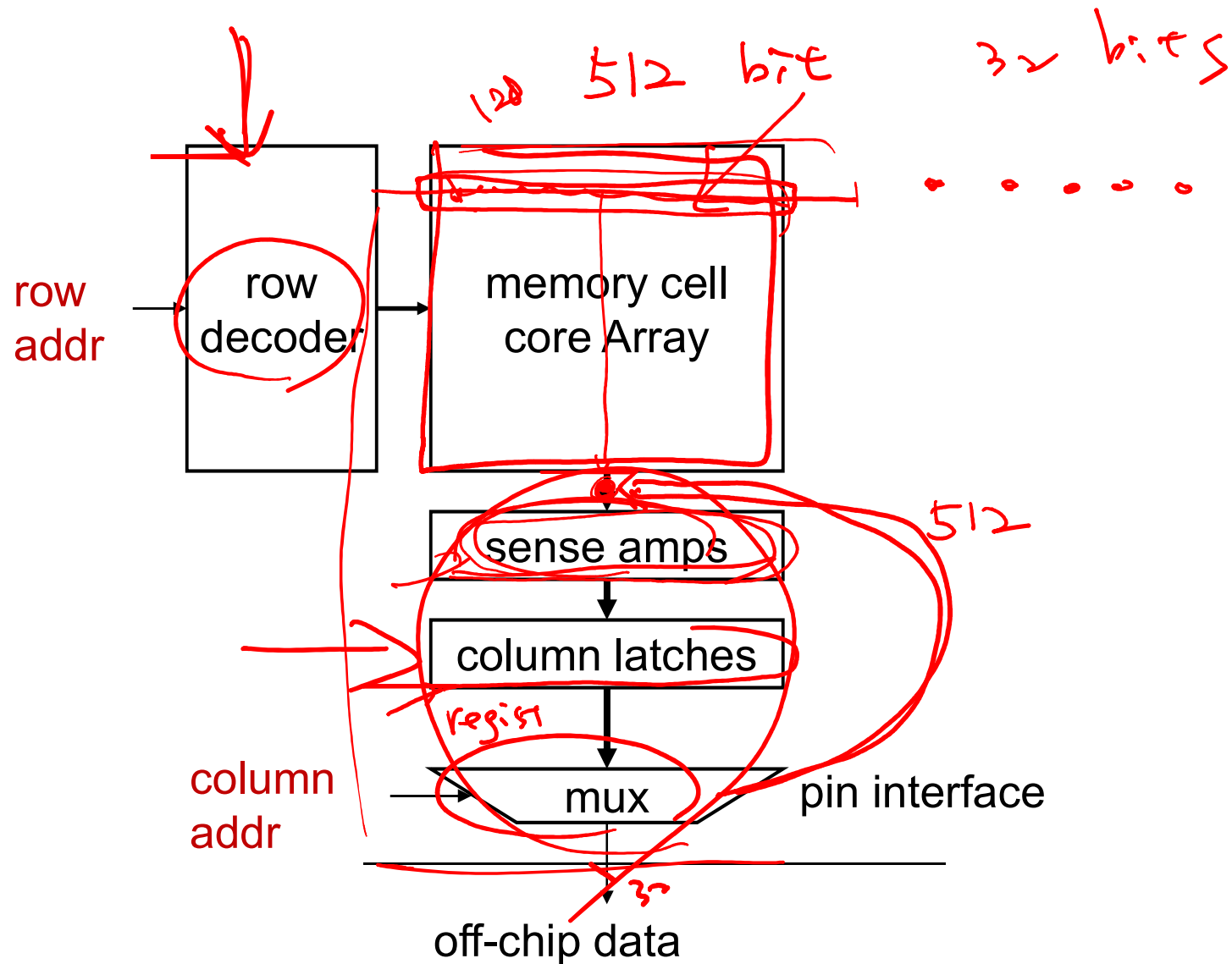4-core intel i7 8MB
"        "   ~5 6M7
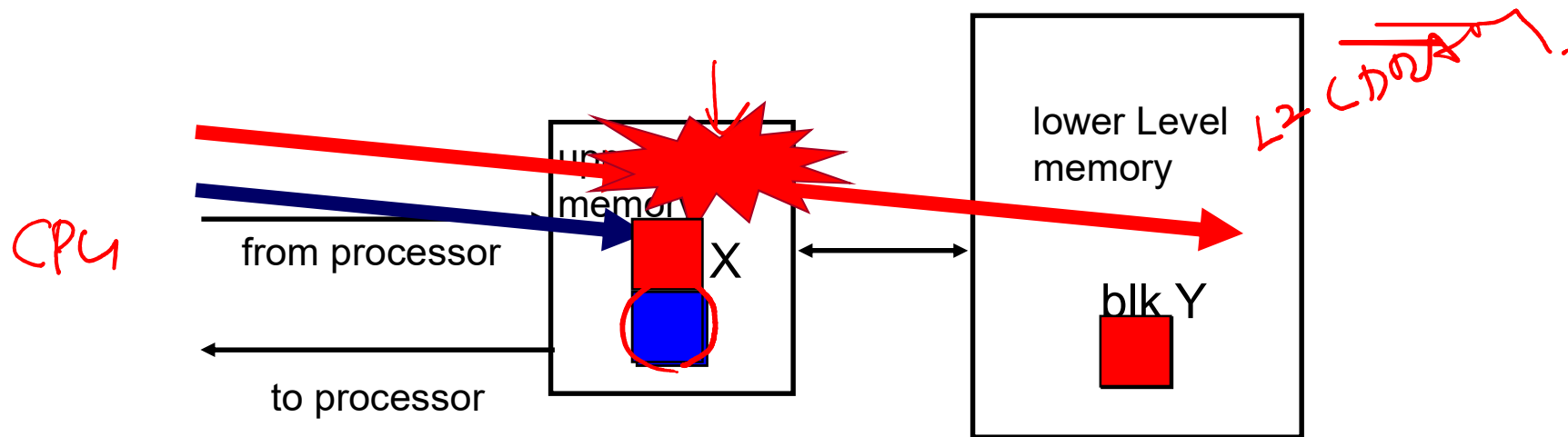
# Example: Intel Nehalem Memory Hierarchy



1. 4-way set associative instruction cache

2. 8-way set associative L1 data cache (32 KB)

3. 8-way set associative L2 data cache (256 KB)

4. 16-way shared L3 cache (8 MB)

5. 3 DDR3 memory connections

# Typical Memory Organization

**row addr** → **row decoder** → **memory cell core Array**

*handwritten:* 120  512 bit   32 bits

**sense amps** (512)

**column latches** (regist)

**column addr** → **mux** — pin interface

**off-chip data**

# Basic Cache Operation

upper memory

from processor

X

to processor

lower Level memory

blk Y

CPU

L2 CDDR

ECE 411 COMPUTER ORGANIZATION AND DESIGN

# Cache Terminology

- **hit**: data appears in some block
  - ✓ **hit rate**: the fraction of accesses found in the level   90 %
  - ✓ **hit time**: time to access the level (consists of RAM access time + time to determine hit)

- **miss**: data needs to be retrieved from a block in the lower level (e.g., block Y)
  - ✓ **miss rate**  = 1 - (hit rate)
  - ✓ **miss penalty**: time to replace a block in the upper level  + time to deliver the block to the processor

- **hit time << miss penalty**

  "1"              10 — 20              any hundred cycles.

# Average Memory Access Time

- average memory-access time
    = hit time + miss rate x miss penalty

*(handwritten: hit time(1 − miss rate) + miss rate)*
*miss rate x miss penalty*

- miss penalty: time to fetch a block from lower memory level
    - ✓ access time:  function of latency
    - ✓ transfer time: function of bandwidth b/w levels
        - transfer one "cache line/block" at a time
        - transfer at the size of the memory-bus width

# Memory Hierarchy Performance

- Average Memory Access Time (AMAT)
  - ✓ = hit time + miss rate × miss penalty
  - ✓ = $T_{hit}$(L1) + miss%(L1) × T(memory)

- example:
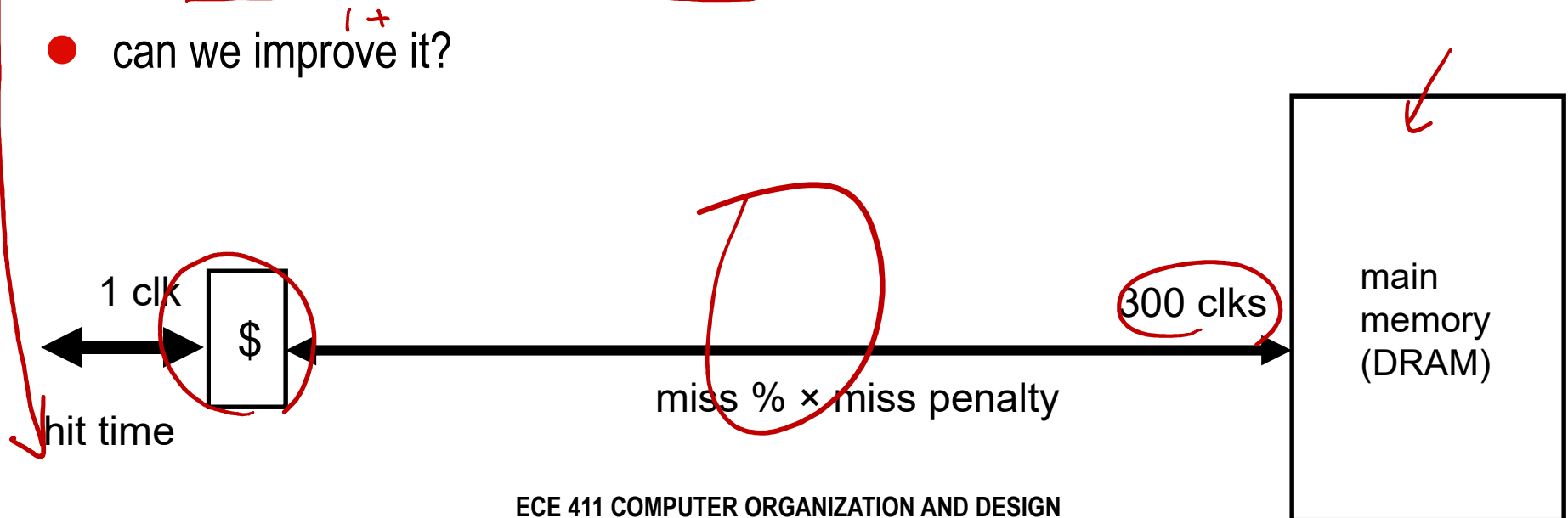  - ✓ cache hit = 1 cycle
  - ✓ miss rate = 10% = 0.1
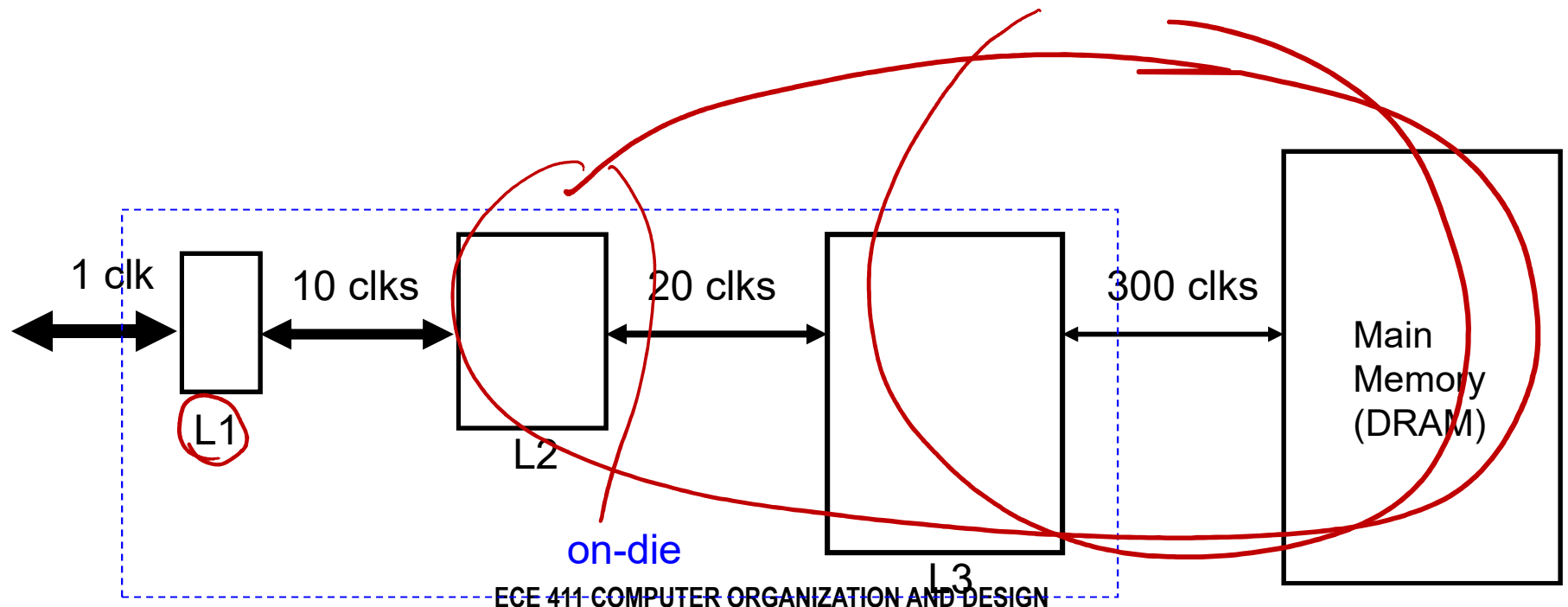  - ✓ miss penalty = 300 cycles
  - ✓ AMAT = 1 + 0.1 × 300 = 31 cycles

- can we improve it?

1 clk

$

hit time

miss % × miss penalty

300 clks

main memory (DRAM)

**ECE 411 COMPUTER ORGANIZATION AND DESIGN**

# Reducing Penalty: Multi-Level Cache

● Average Memory Access Time (AMAT)

$= T_{hit}(L1) + miss\%(L1)\times(T_{hit}(L2) + miss\%(L2)\times(T_{hit}(L3) + miss\%(L3)\times T(memory)\ )\ )$

$= T_{hit}(L1) + miss\%(L1)\times T_{miss}(L1)$

$= T_{hit}(L1) + miss\%(L1)\times\{T_{hit}(L2) + miss\%(L2)\times(T_{miss}(L2)\ \}$

$= T_{hit}(L1) + miss\%(L1)\times\{T_{hit}(L2) + miss\%(L2)\times(T_{miss}(L2)\ \}$

$= T_{hit}(L1) + miss\%(L1)\times\{T_{hit}(L2) + miss\%(L2)\times[\ T_{hit}(L3) + miss\%(L3)\times T(memory)\ ]\ \}$

1 clk    10 clks    20 clks    300 clks

L1    L2    L3    Main Memory (DRAM)

on-die

ECE 411 COMPUTER ORGANIZATION AND DESIGN

# AMAT Example

$$= T_{hit}(L1) + \text{miss\%}(L1) \times (T_{hit}(L2) + \text{miss\%}(L2) \times (T_{hit}(L3) + \text{miss\%}(L3) \times T(\text{memory})))$$

- Example:
  - ✓ miss rate L1=10%, $T_{hit}(L1)$ = 1 cycle
  - ✓ miss rate L2=5%, $T_{hit}(L2)$ = 10 cycles
  - ✓ miss rate L3=1%, $T_{hit}(L3)$ = 20 cycles
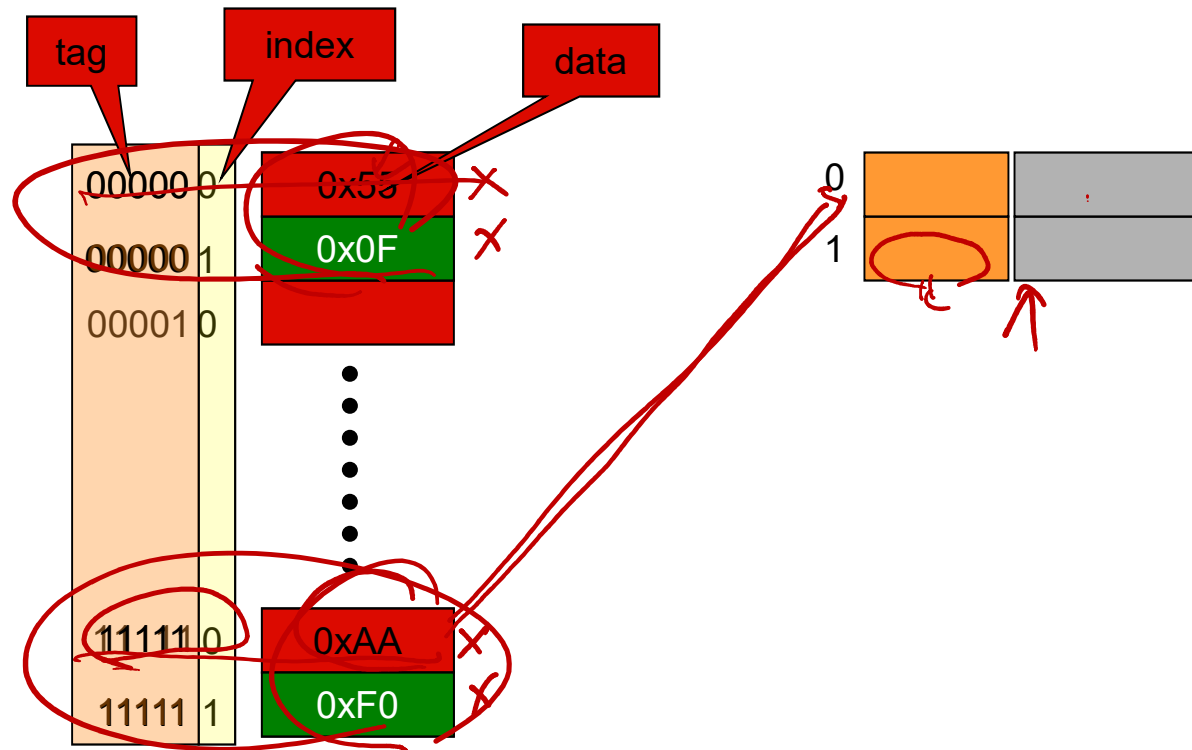  - ✓ T(memory) = 300 cycles

- AMAT = ?
  - ✓ 2.115 (compare to 31 with no multi-levels)
  - ✓ 14.7× speed-up!

# Types of Caches

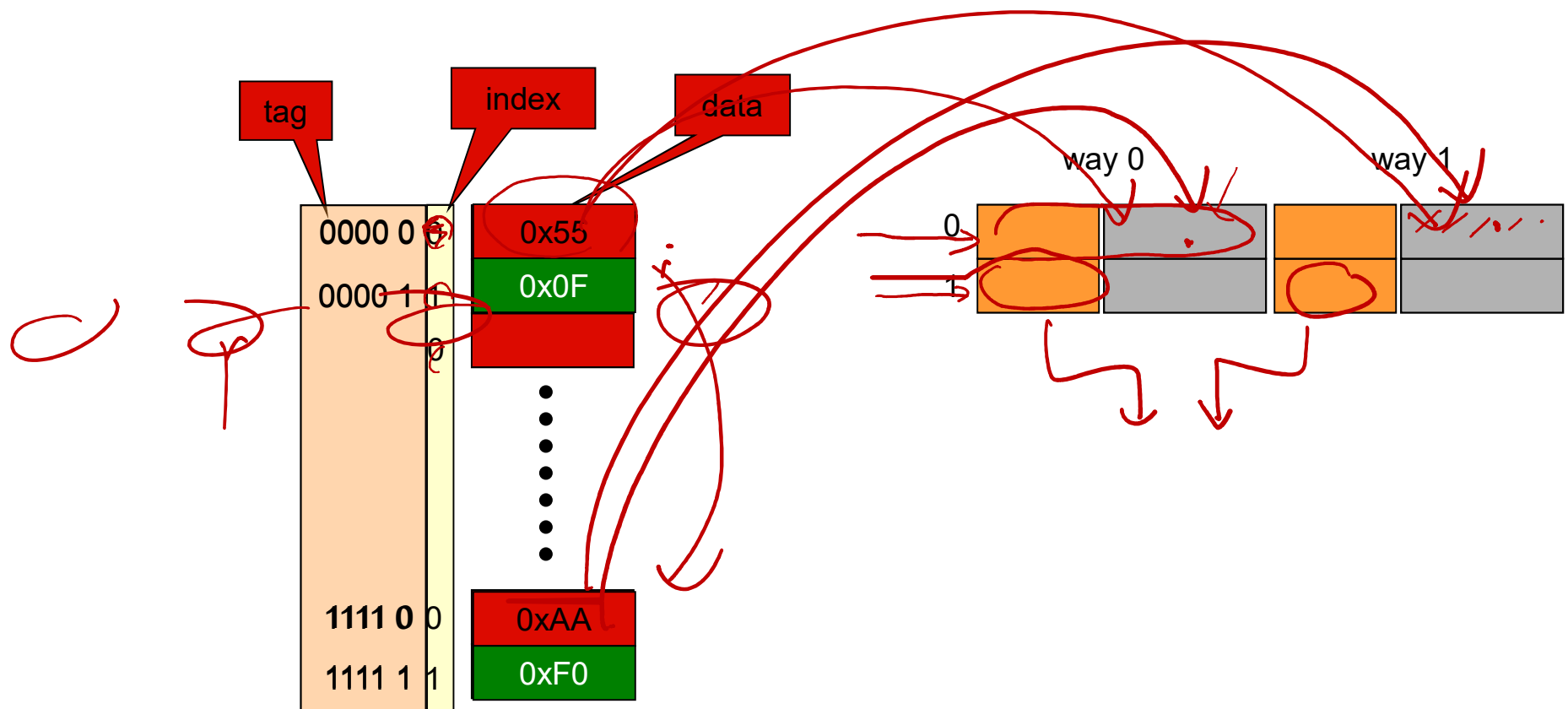| type of cache | mapping of data from memory to cache | complexity of searching the cache |
|---|---|---|
| direct mapped (DM) | a memory value can be placed at **a single corresponding location** in the cache | fast indexing mechanism |
| set-associative (SA) | a memory value can be placed in **any of a set of locations** in the cache | slightly more involved search mechanism |
| fully-associative (FA) | a memory value can be placed in **any location** in the cache | extensive hardware resources required to search (CAM) |

# Direct Mapping

- direct mapping:
  - ✓ a memory value can only be placed at a single corresponding location in the cache
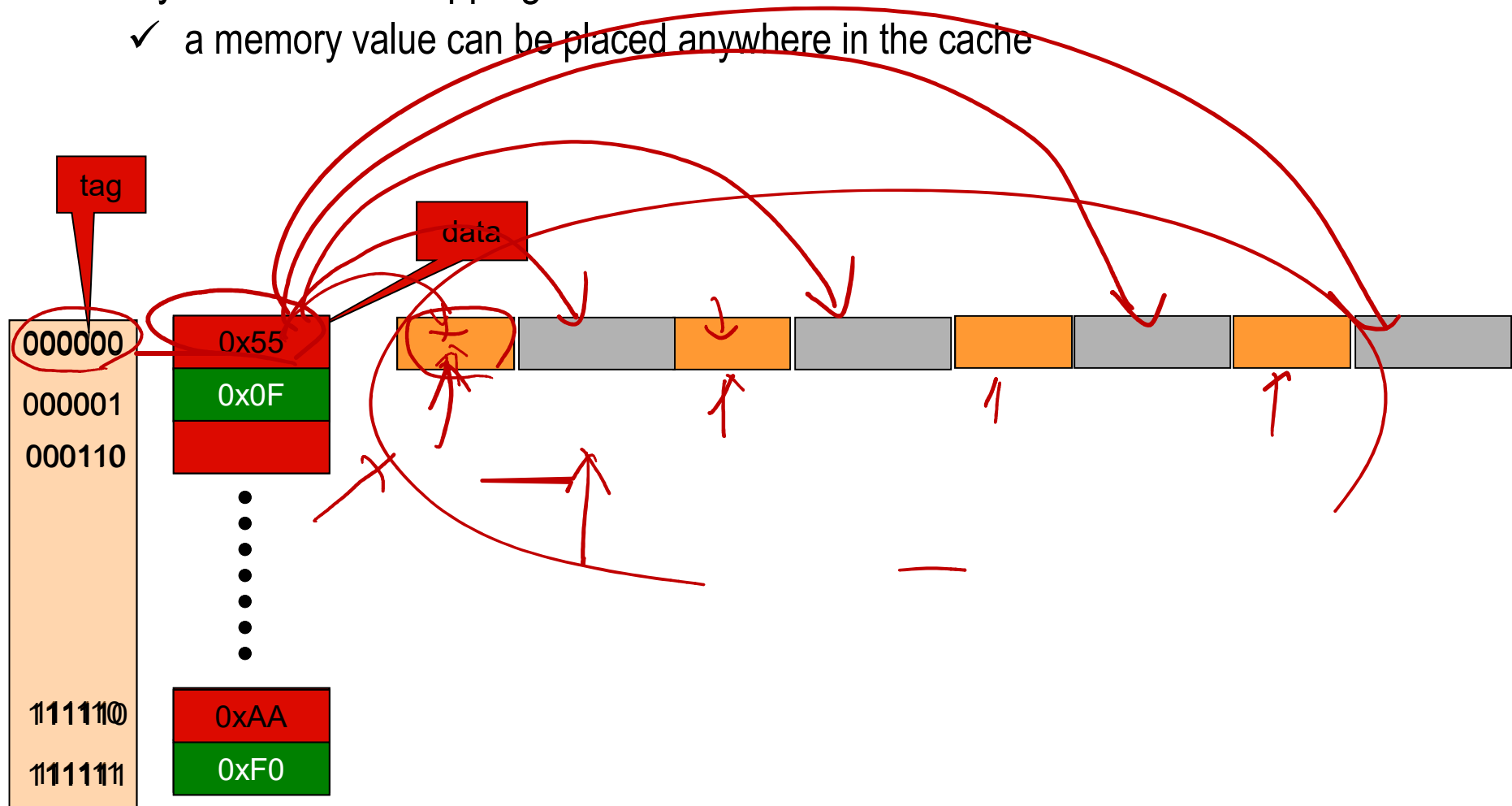
# Set Associative Mapping (2-Way)

- set-associative mapping:
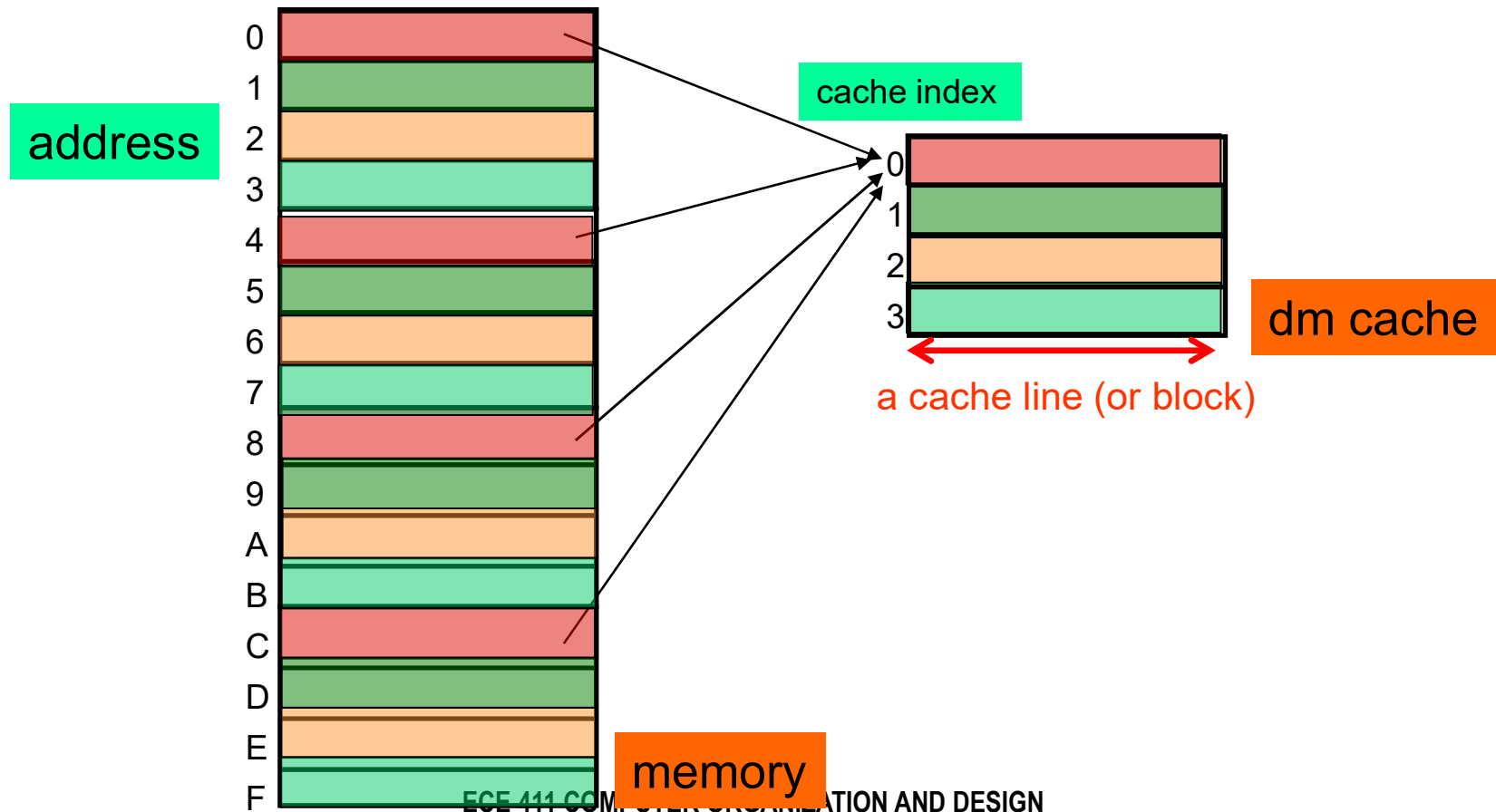  - ✓ a memory value can be placed in any location of a set in the cache

tag   index   data

way 0   way 1

| | 0x55 |
|---|---|
| | 0x0F |

0x55
0x0F

0000 0
0000 1

1111 0 0   0xAA
1111 1 1   0xF0

0xAA
0xF0

0
1

# Fully Associative Mapping

● fully-associative mapping:

   ✓ a memory value can be placed anywhere in the cache

tag

data

| 000000 | 0x55 |
| 000001 | 0x0F |
| 000110 | |

⋮
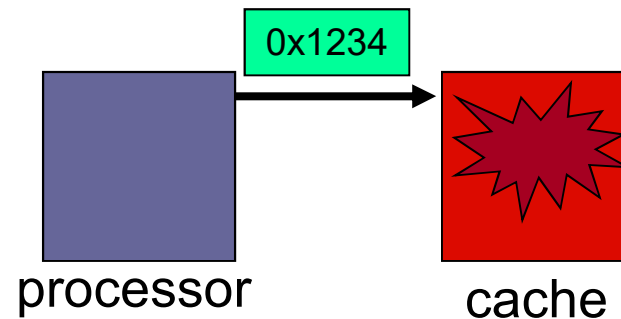
| 111110 | 0xAA |
| 111111 | 0xF0 |

# Direct Mapped Cache

- location 0 is occupied by data from (0, 4, 8, and C)
  - ✓ which one should we place in the cache?
  - ✓ how can we tell which one is in the cache?



address

cache index

dm cache
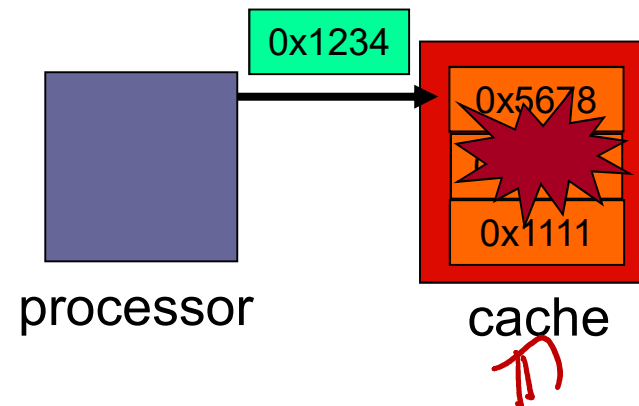
a cache line (or block)

memory

# Three Cs (Cache Miss Terms)

- compulsory misses:     *conflic misses*     *capacity miss*
  - ✓ cold start misses (caches do not have valid data at the start of the program)

0x1234

processor          cache

# Three Cs (Cache Miss Terms)

- capacity misses:
  - ✓ increase cache size

0x1234

0x5678

0x1111

processor

cache

# Three Cs (Cache Miss Terms)

- conflict misses:
  - ✓ increase cache size and/or associativity.
  - ✓ associative caches reduce conflict misses



processor          cache

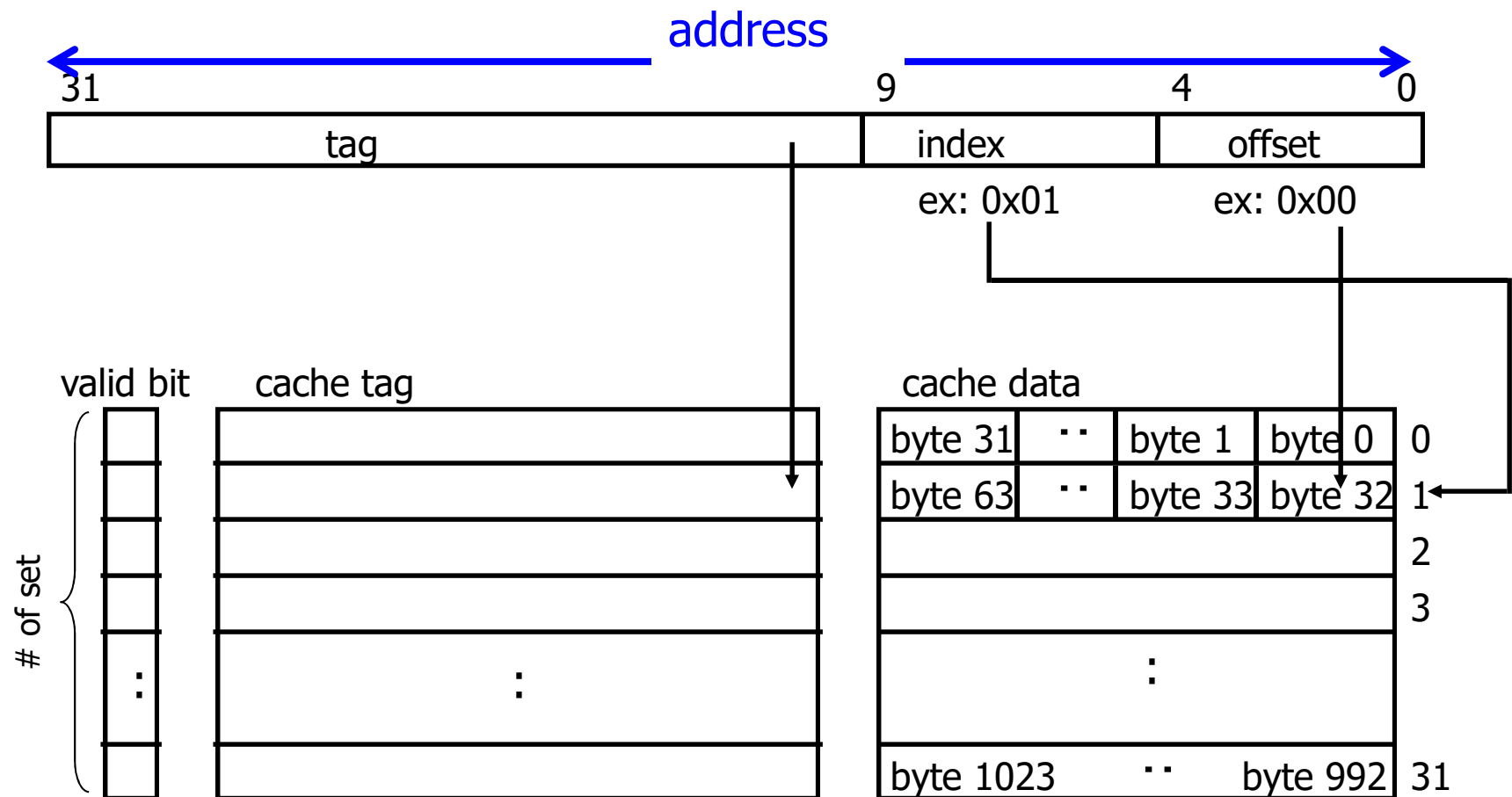# Four Central Questions in Designing a Cache

- P-I-R-W:
  - ✓ placement: where can a block of memory go?
  - ✓ identification: how do i find a block of memory?
  - ✓ replacement: how do i make space for new blocks?
  - ✓ write policy: how do i propagate changes?

- need to consider these for all levels of the memory hierarchy
  - ✓ L1/L2/L3 caches now

- main memory, disks have similar issues, addressed later

# Describing Caches: 7 Parameters

- access time: $T_{hit}$

- capacity
  - ✓ total amount of data the cache can hold
    - o # of blocks × block size

- block (line) size
  - ✓ the amount of data that gets moved into or out of the cache as a chunk
    - o analogous to page size in virtual memory

- replacement policy
  - ✓ what data is replaced on a miss?

- associativity
  - ✓ how many locations in the cache is a given address eligible to be placed in?

- unified, instruction, data
  - ✓ what type of data is kept in the cache? We'll cover this in more detail later

# Example: 1KB DM Cache, 32-byte Lines

- lowest M bits are offset (Line Size = 2M)

- index = log2 (# of sets)



address

31                                                    9          4          0

| tag | index | offset |

ex: 0x01        ex: 0x00

valid bit    cache tag                    cache data

| byte 31 | ·· | byte 1 | byte 0 | 0 |
| byte 63 | ·· | byte 33 | byte 32 | 1 |
|  |  |  |  | 2 |
|  |  |  |  | 3 |

# of set

:                          :                          :

| byte 1023 | ·· | byte 992 | 31 |

**ECE 411 COMPUTER ORGANIZATION AND DESIGN**

# Example of Caches

- given a 2MB, direct-mapped physical caches, line size=64bytes, and 52-bit physical address
  - ✓ tag size?

  - ✓ now change it to 16-way, tag size?

  - ✓ how about if it's fully associative, tag size?

**ECE 411 COMPUTER ORGANIZATION AND DESIGN**

# Announcement

- today's lecture: cache basics
  - ✓ Ch. 5.1 – 5.4 (HP1)

- next lecture: other cache topics (e.g., replacement policy)
  - ✓ Ch. 5.4 – 5.8 (HP1)

- MP assignment
  - ✓ MP1 due on 2/4 5pm