# Lecture 23: Cache, Memory

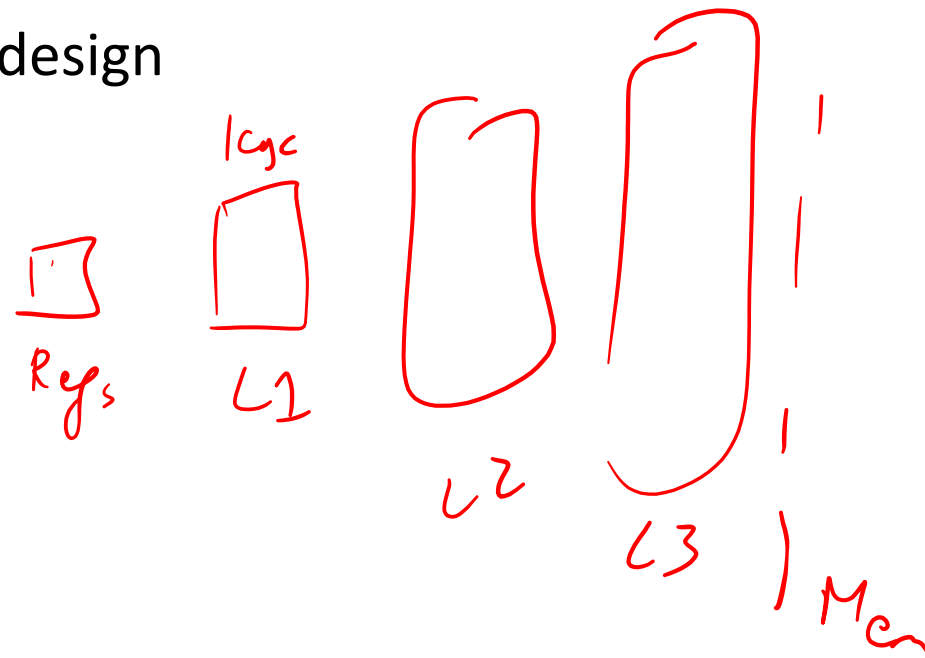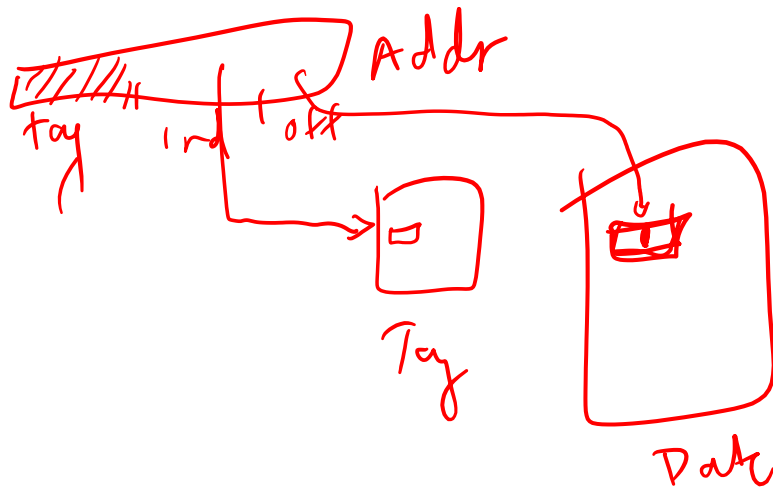*(handwritten, top right)* 40% of instr end up in L2
20% of in end up in L3

- Today's topics:

  - Example problems in cache design
  - Caching policies
  - Main memory system

*(handwritten notes)*

HW 9 due Tue/Wed

logic

Regs   L1   L2   L3 ) Mem

Addr
tag  ind  off

Tag

Data

# Example 2

Offset = address % blksize
Index = (address / blksize) % sets
Tag = address / (blksize × sets)

Show how the following addresses map to the cache and yield hits or misses.
The cache is direct-mapped, has 16 sets, and a 64-byte block size.
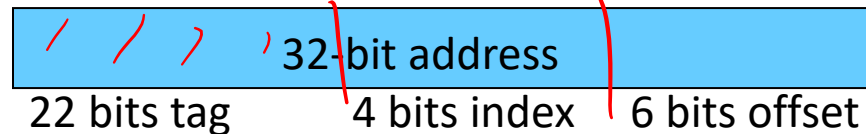Addresses:  8, 96, 32, 480, 976, 1040, 1096

Offset = address % 64  (address modulo 64, extract last 6)
Index = address/64 % 16    (shift right by 6, extract last 4)
Tag = address/1024        (shift address right by 10)

32-bit address

| | 22 bits tag | 4 bits index | 6 bits offset | |
|---|---|---|---|---|
| 8: | 0 | 0 | 8 | M |
| 96: | 0 | 1 | 32 | M |
| 32: | 0 | 0 | 32 | H |
| 480: | 0 | 7 | 32 | M |
| 976: | 0 | 15 | 16 | M |
| 1040: | 1 | 0 | 16 | M |
| 1096: | 1 | 1 | 8 | M |

# Example 3

- A pipeline has CPI 1 if all loads/stores are L1 cache hits
→ 40% of all instructions are loads/stores
  85% of all loads/stores hit in 1-cycle L1
  50% of all (10-cycle) L2 accesses are misses
  Memory access takes 100 cycles
  What is the CPI?

1000 instrs

→ 1000 cycles
(ignore pipeline warm-up)

400 are ld/st

85% of 400 = 340 are L1 hits

15% (60 instrs) are L1 misses
not stalls

60 × 10 cyc = 600 cyc accessing L2.  ←

→30 are L2 misses
30 × 100 cyc = 3000 cyc access memory. ←

Exec time
= 1000 + 600
    + 3000
= 4600 cyc

CPI = $\frac{4600}{1000}$ = 4.6

3

# Example 3

- A pipeline has CPI 1 if all loads/stores are L1 cache hits
  40% of all instructions are loads/stores
  85% of all loads/stores hit in 1-cycle L1
  50% of all (10-cycle) L2 accesses are misses
  Memory access takes 100 cycles
  What is the CPI?

Start with 1000 instructions
1000 cycles         (includes all 400 L1 accesses)
+ 400 (ld/st) x 15% x 10 cycles  (the L2 accesses)
+ 400 x 15% x 50% x 100 cycles  (the mem accesses)
=  4,600 cycles
CPI = 4.6

# Example 4

8B blocks

Assume that addresses are 8 bits long
How many of the following address requests
are hits/misses?
4, 7, 10, 13, 16, 24, 36, 4, 48, 64, 4, 36, 64, 4

LRU

Byte address

00010000

Tag

Tag array

Way-1     Way-2

| 0/4/7  64-71 | 32 - 39 |
| 8 - 15 | 40 - 47 |
| 16 - 23 | 48 - 55 |
| 24 - 31 | 56 - 63 / 8-byte blocks |

4 sets

2 ways

Data array

1,000
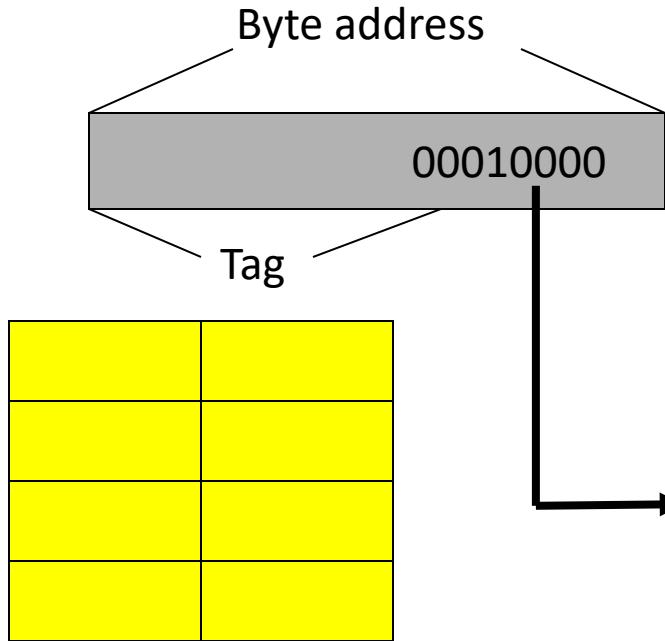index  3 offset bits
0 1 0  0 0 0

5

# Example 4

Byte address

00010000

Tag

Tag array

Assume that addresses are 8 bits long
How many of the following address requests are hits/misses?
4, 7, 10, 13, 16, 24, 36, 4, 48, 64, 4, 36, 64, 4
M H M  H  M  M  M H M  M  H  M  M  M

Way-1          Way-2

8-byte blocks

Data array

6

# Cache Misses

16-way 32MB L3

24-way 24MB L3

(11) way 22MB L3

st

- On a write miss, you may either choose to bring the block into the cache (write-allocate) or not (write-no-allocate)

- On a read miss, you always bring the block in (spatial and temporal locality) – but which block do you replace?
  - ➤ no choice for a direct-mapped cache
  - ➤ randomly pick one of the ways to replace
  - ➤ replace the way that was least-recently used (LRU)
  - ➤ FIFO replacement (round-robin)

Pseudo-LRU

8-way cache
recency list for
each set

MR: 3, 5, 1, 2, 6, 7 ···· LR

MR: 1, 3, 5, 2, 6, 7 ····

# Writes

- When you write into a block, do you also update the copy in L2?
  - ➤ write-through: every write to L1 → write to L2
  - ➤ write-back: mark the block as dirty, when the block gets replaced from L1, write it to L2

- Writeback coalesces multiple writes to an L1 block into one L2 write

- Writethrough simplifies coherency protocols in a multiprocessor system as the L2 always has a current copy of data

# Types of Cache Misses

3 C's

prefetching → compiler

- Compulsory misses: happens the first time a memory word is accessed – the misses for an infinite cache

  M  M        M  M
  A  A        A  A

  tiny → app developer

  Comp

- Capacity misses: happens because the program touched many other words before re-touching the same word – the misses for a fully-associative cache

  Capacity miss

  confl miss
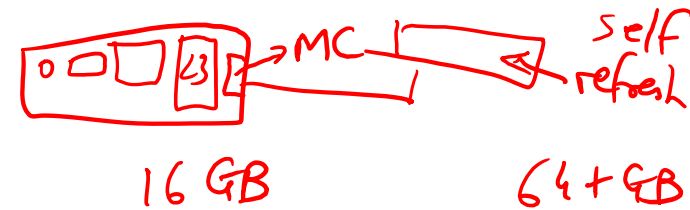
  hsv ← desig a cache with many ways

- Conflict misses: happens because two words map to the same location in the cache – the misses generated while moving from a fully-associative to a direct-mapped cache

# Off-Chip DRAM Main Memory

*(handwritten annotations: MC, self refresh, 16 GB, 64+GB, 1 TB)*

- Main memory is stored in <u>DRAM cells</u> that have much higher storage density

*(handwritten: Dynamic → many GBs of capacity, low cost, high density; Periodically refreshed)*

- DRAM cells lose their state over time – must be refreshed periodically, hence the name *Dynamic*

*(handwritten: DRAM cell = Capacitor; 8x higher density)*

- A number of DRAM chips are aggregated on a DIMM to provide high capacity – a DIMM is a module that plugs into a bus on the motherboard

*(handwritten: SRAM cell)*

*(handwritten: A – every ms; B – a few times every sec; every 64 ms)*

- DRAM access suffers from long access time and high energy overhead

*(handwritten: C – every min; D – even hour; E – Every day; Caches – SRAM cells → Caches, static, speed, less focus on density)*

10

# Memory Architecture

Processor

Memory Controller

Bank

Row Buffer

Address/Cmd

Data

DIMM

- DIMM: a PCB with DRAM chips on the back and front
- The memory system is itself organized into ranks and banks; each bank can process a transaction in parallel
- Each bank has a row buffer that retains the last row touched in a bank (it's like a cache in the memory system that exploits spatial locality) (row buffer hits have a lower latency than a row buffer miss)