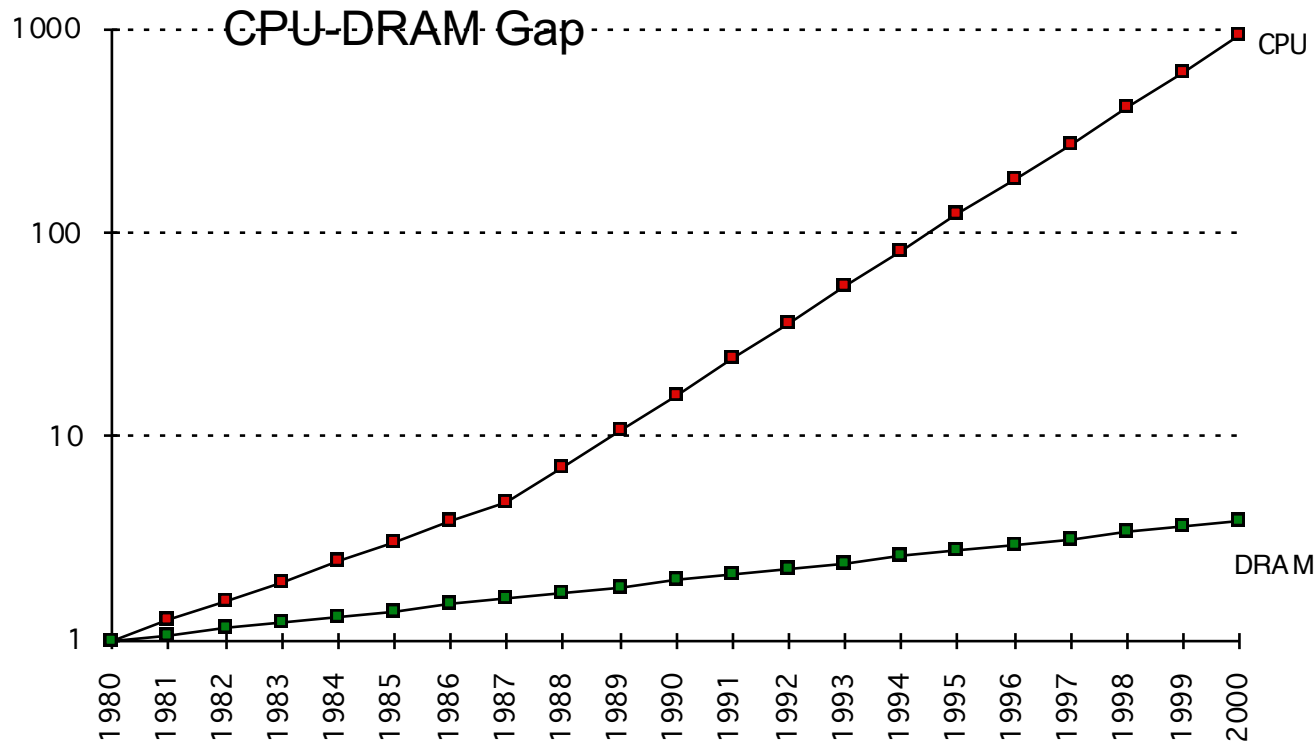# Lecture 8

## Memory Hierarchy

# Memory Hierarchy

- Keeping useful books close to you
- Keeping useful data close to the CPU
  - Memory Hierarchy

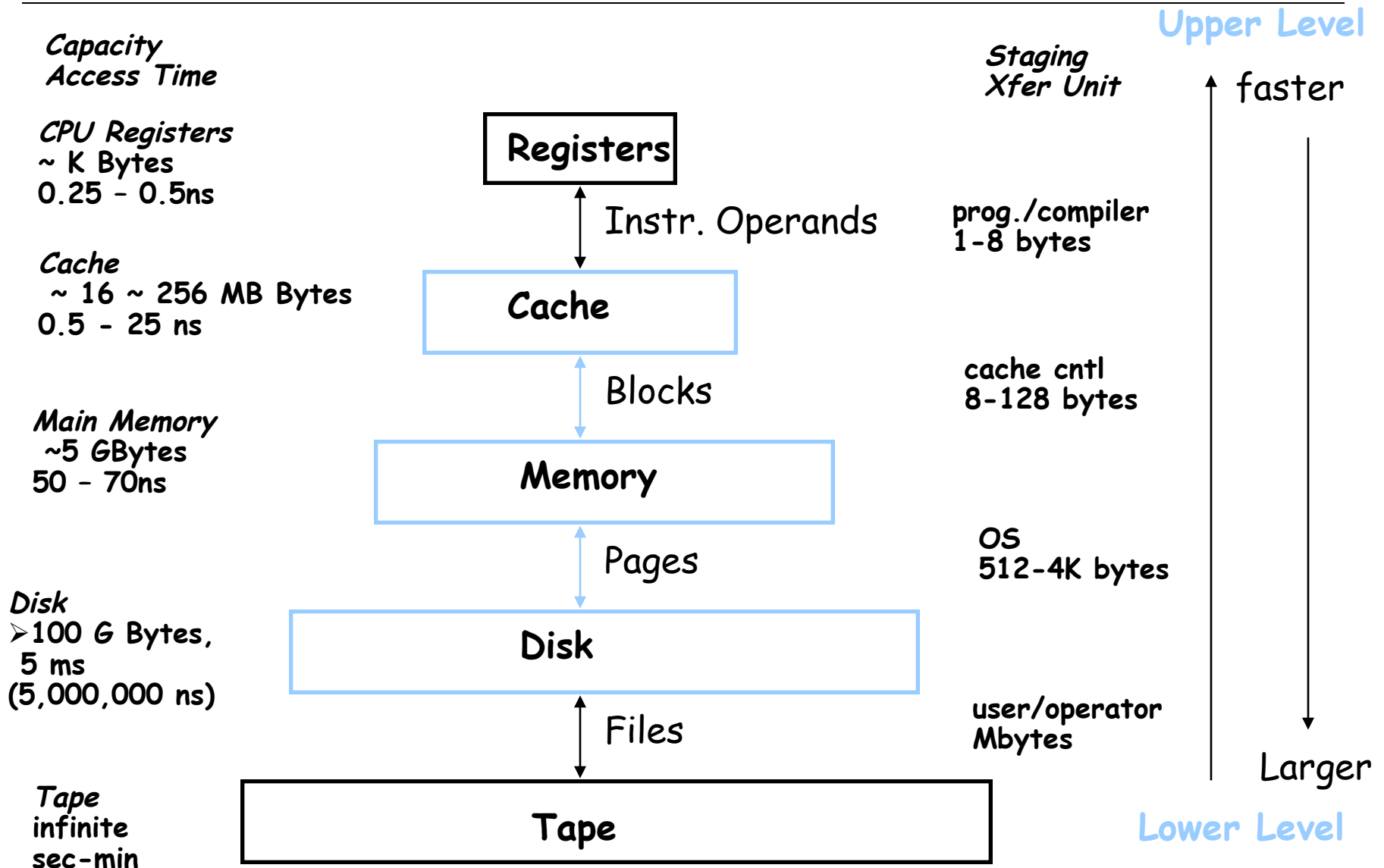# Why Cares about Memory Hierarchy?

- **Processor Only Thus Far in Course**
  - CPU cost/performance, ISA, Pipelined Execution



CPU-DRAM Gap

# Levels of the Memory Hierarchy

*Capacity*
*Access Time*

*Staging*
*Xfer Unit*

faster

*CPU Registers*
*~ K Bytes*
*0.25 – 0.5ns*

**Registers**

Instr. Operands

prog./compiler
1-8 bytes

*Cache*
*~ 16 ~ 256 MB Bytes*
*0.5 - 25 ns*

**Cache**

Blocks

cache cntl
8-128 bytes

*Main Memory*
*~5 GBytes*
*50 – 70ns*

**Memory**

Pages

OS
512-4K bytes

*Disk*
➢*100 G Bytes,*
*5 ms*
*(5,000,000 ns)*

**Disk**

Files

user/operator
Mbytes

*Tape*
**infinite**
**sec-min**

**Tape**

Larger

# General Principle

- **The Principle of Locality:**
  - Programs access a relatively small portion of the address space at any instant of time.

- **Two Different Types of Locality:**
  - **Temporal Locality** (Locality in Time):
    - If an item is referenced, it will tend to be referenced again soon (e.g., loops, reuse)

  - **Spatial Locality** (Locality in Space):
    - If an item is referenced, items whose addresses are close by tend to be referenced soon  (e.g., straightline code, array access)
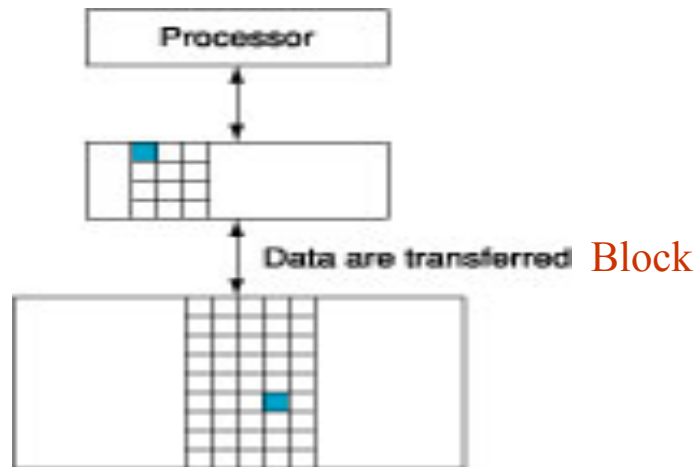
    ```
    For (i= 0; i++;i:<100)
       {
         A[i]=A[i] + x;
         x++;
       }
    ```

- **Locality + smaller HW is faster = memory hierarchy**
  - *Levels*: each smaller, faster, more expensive/byte than level below
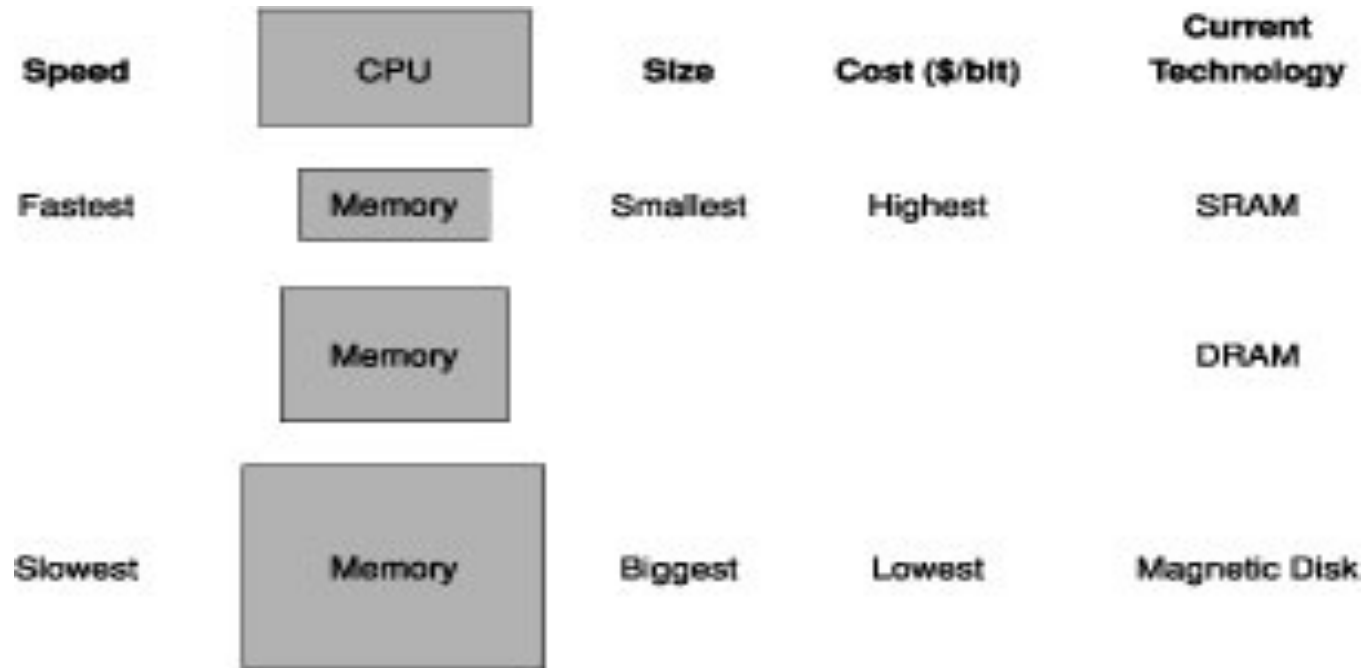  - *Inclusive*: data found in top also found in the bottom

5

# Memory Hierarchy: Terminology

- **Definitions**
    - *Upper* is closer to processor
    - *Block*: minimum unit that present or not in upper level
- **Hit**: data appears in some block in the upper level (example: Block X)
    - Hit Rate: the fraction of memory access found in the upper level
    - Hit Time: Time to access the upper level which consists of
        RAM access time + Time to determine hit/miss
- **Miss**: data needs to be retrieved from a block in the lower level (Block Y)
    - Miss Rate = 1 - (Hit Rate)
    - Miss Penalty: Time to replace a block in the upper level +
        Time to deliver the block the processor
- **Hit Time << Miss Penalty (500 instructions on Alpha 21264!)**



Data are transferred  Block
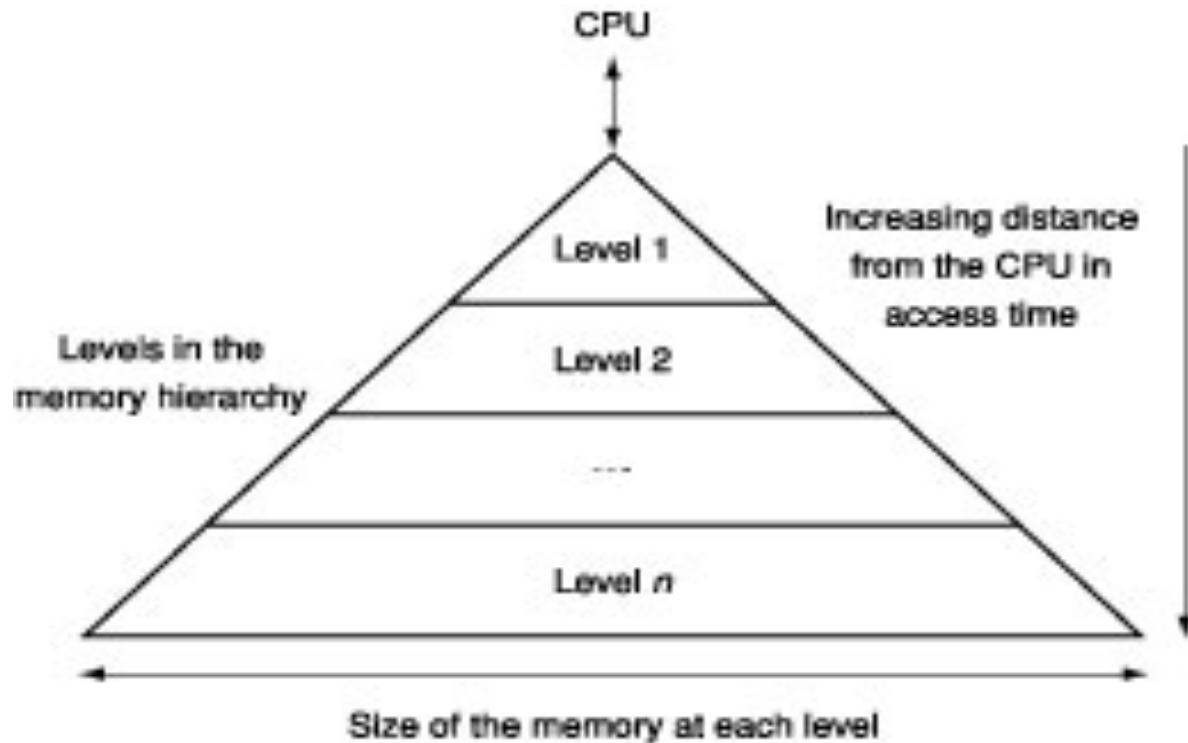
# Basic Structure of a Memory Hierarchy

| | | | | Current |
|---|---|---|---|---|
| Speed | CPU | Size | Cost ($/bit) | Technology |
| Fastest | Memory | Smallest | Highest | SRAM |
| | Memory | | | DRAM |
| Slowest | Memory | Biggest | Lowest | Magnetic Disk |

| Memory Technology | Typical access time | $ per GB in 2012 |
|---|---|---|
| SRAM | 0.5-5 ns | $500-$1000 |
| DRAM | 50 – 70 ns | $10-$20 |
| Flash Memory | 5,000 ~ 50,000ns | $0.75 -$1 |
| Magnetic disk | 5,000,000-20,000,000 ns | $0.50-$2 |

# Memory Hierarchy
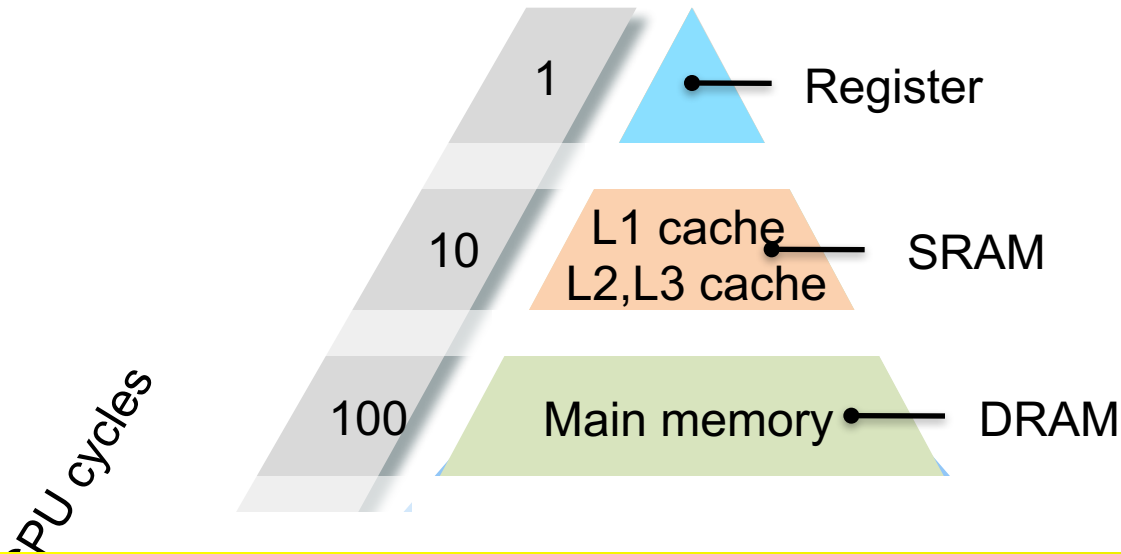


1. Keeping more recently accessed data items closer to the processor -> temporal locality
2. Moving blocks consisting of multiple contiguous words in memory to upper levels of the hierarchy -> spatial locality
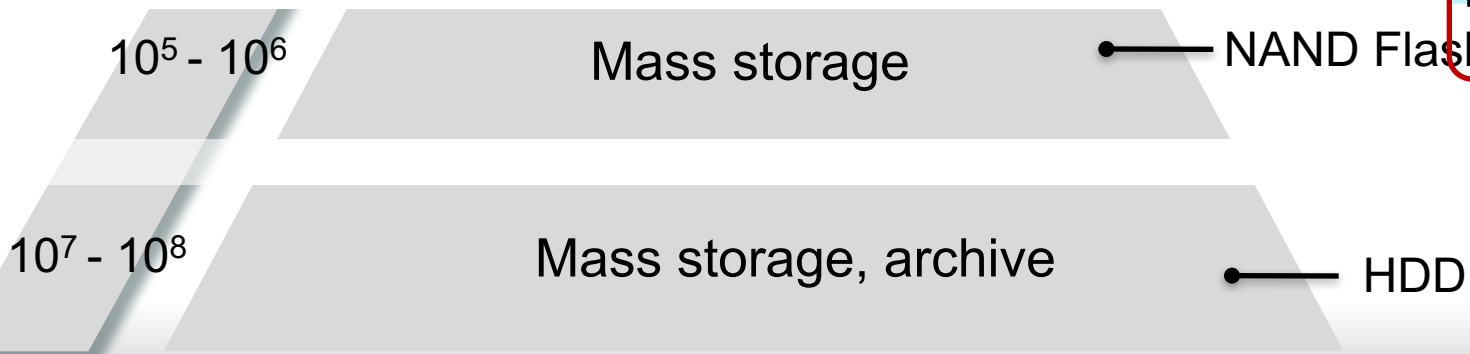3. Data cannot be present in level i unless it is also present in level i+1

# Storage Class Memory



CPU cycles

| | | |
|---|---|---|
| 1 | | Register |
| 10 | L1 cache / L2,L3 cache | SRAM |
| 100 | Main memory | DRAM |
| $10^5$ - $10^6$ | Mass storage | NAND Flash |
| $10^7$ - $10^8$ | Mass storage, archive | HDD |

Fast, byte-addressable

Large, low-cost, non-volatile

**SCM blurs the line between Memory and Storage**
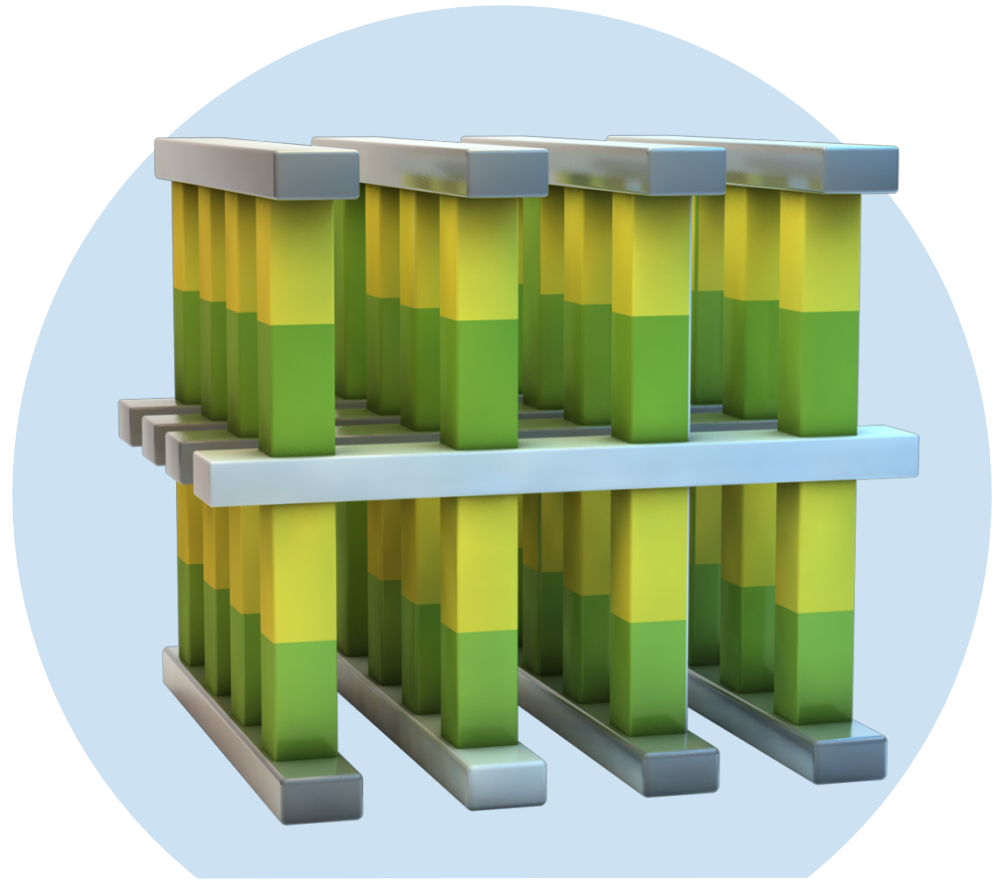
# Intel/Micron: Introducing 3D XPoint™
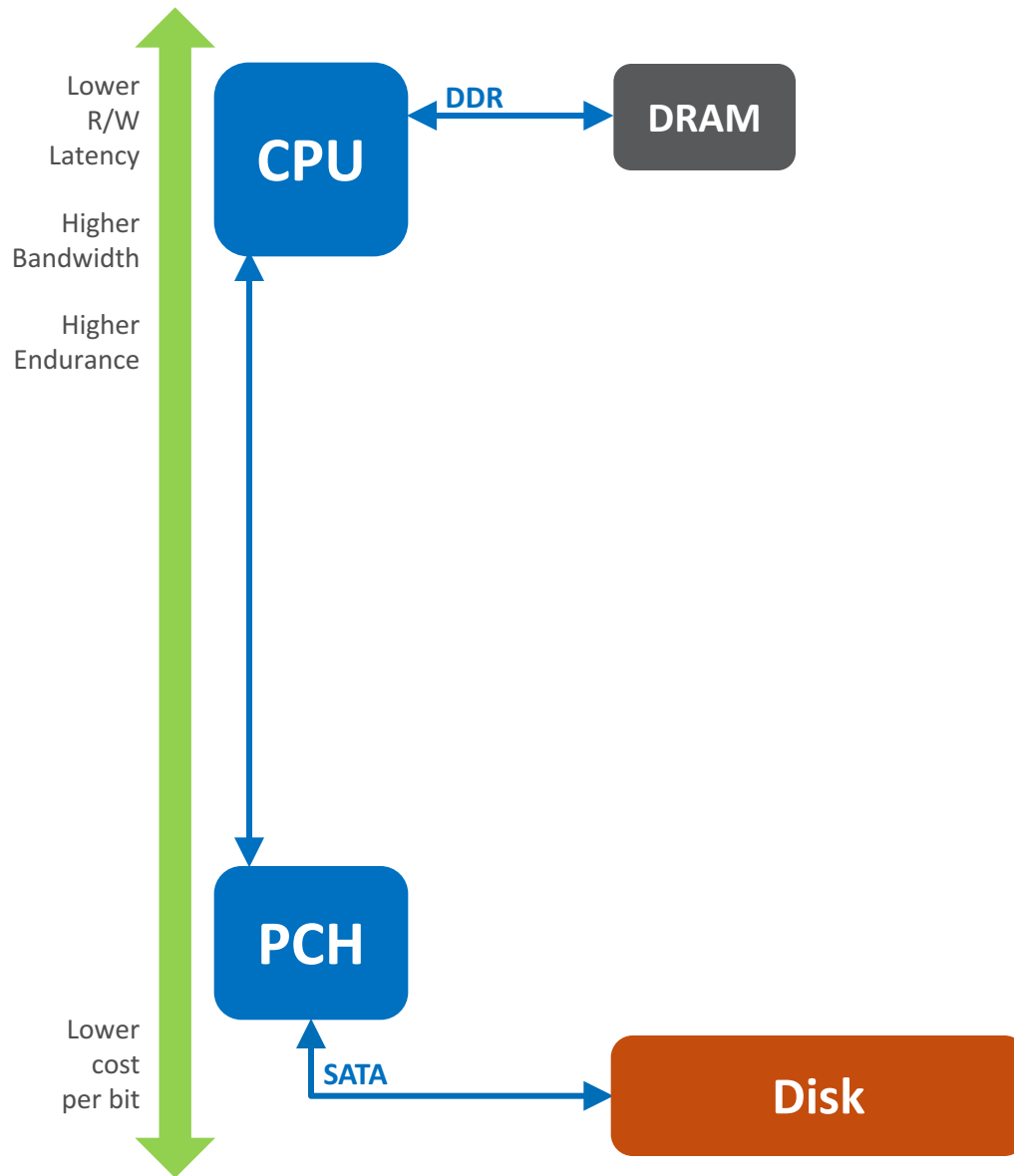
**1000X**
FASTER
THAN NAND

**1000X**
ENDURANCE
OF NAND

**10X**
DENSER
THAN CONVENTIONAL MEMORY


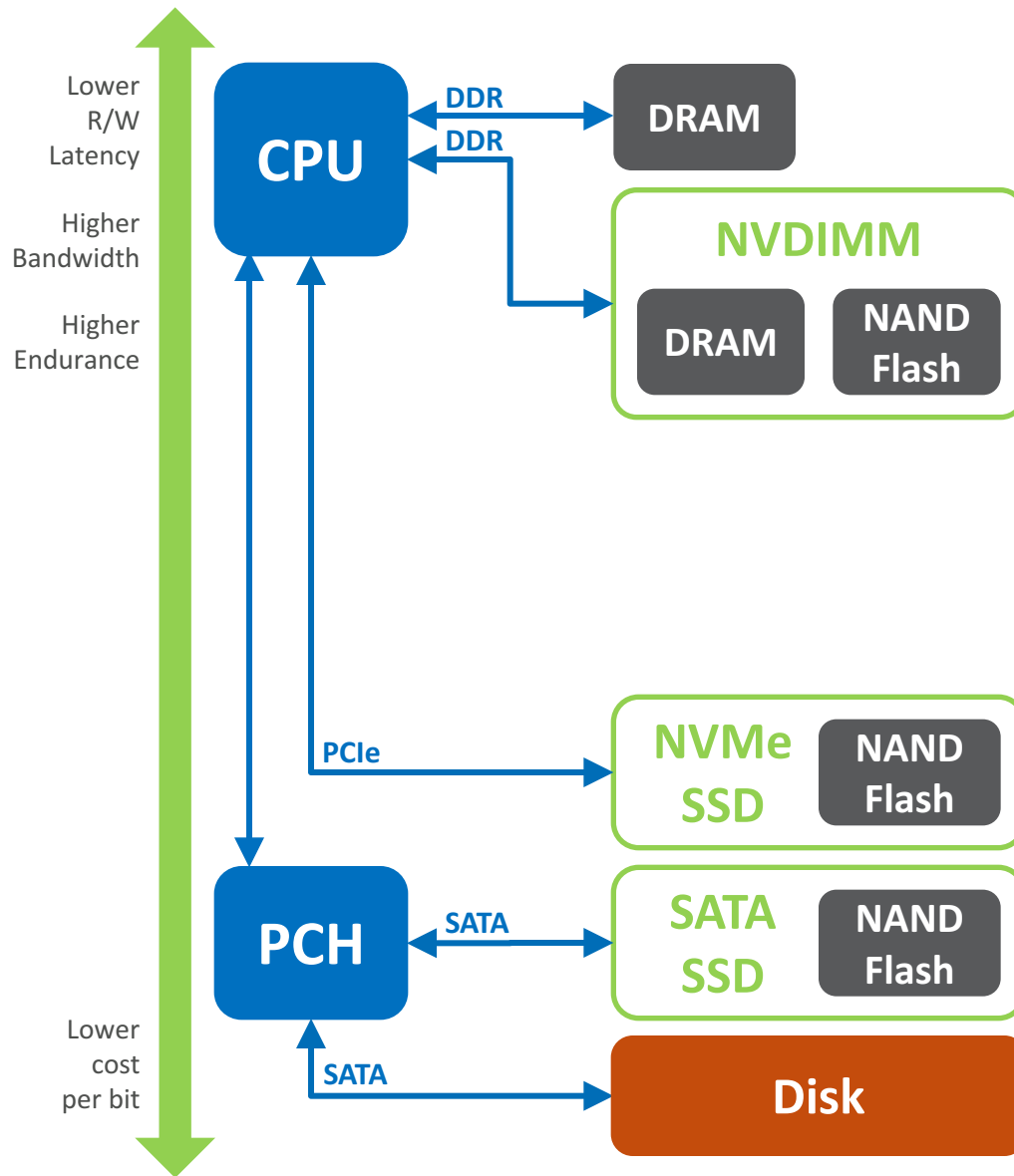
**3D XPoint**

# The Past: Nonvolatile Memories in Server Architectures



Lower R/W Latency

Higher Bandwidth

Higher Endurance

Lower cost per bit

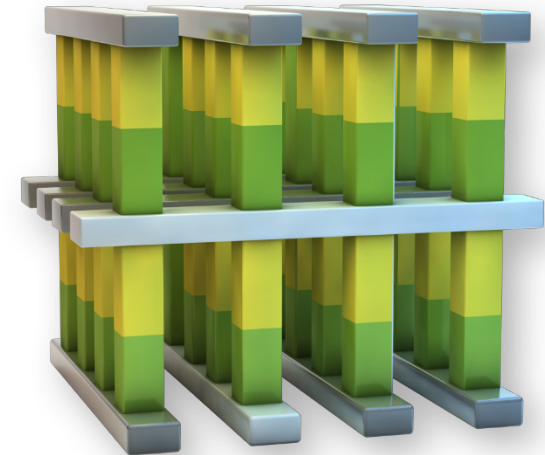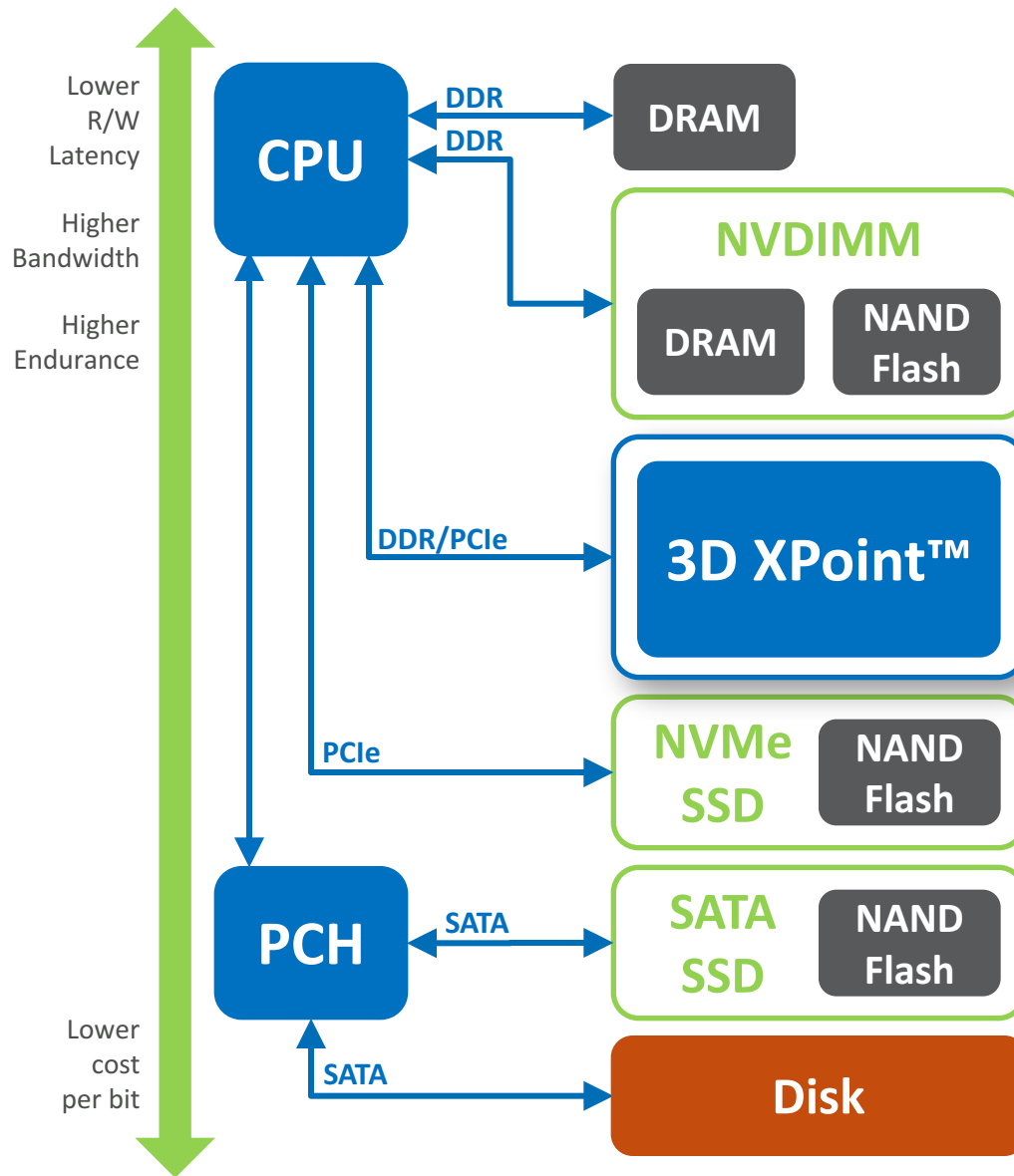**CPU** ↔ **DDR** ↔ **DRAM**

**PCH** — **SATA** → **Disk**

- In the old days of not long ago we had two primary types of memories in computers: DRAM and Hard Disk Drive (HDD)

- DRAM was fast and volatile and HDDs were slower, but nonvolatile

- Data moves from the HDD to DRAM where it is the fed to the processor

- The processor writes the result in DRAM and then it is stored back to disk to remain for future use

# The Present: Nonvolatile Memories in Server Architectures



- System performance increased as the speed of both the interface and the memory accesses improved

- NAND Flash considerably improved the nonvolatile response time

- SATA and PCIe made further optimization to the storage interface

- NVDIMM provides battery- or super capacitor-backed DRAM, operating at DRAM-like speeds and retains data when power is removed

# The Future: Nonvolatile Memories in Server Architectures



- Lower R/W Latency
- Higher Bandwidth
- Higher Endurance
- Lower cost per bit

CPU — DDR — DRAM

CPU — DDR — NVDIMM (DRAM, NAND Flash)

CPU — DDR/PCIe — 3D XPoint™

CPU — PCIe — NVMe SSD (NAND Flash)

PCH — SATA — SATA SSD (NAND Flash)

PCH — SATA — Disk

- 3D XPoint technology provides the benefit in the middle

- It is considerably faster than NAND Flash

- Performance can be realized on PCIe or DDR buses

- Lower cost per bit than DRAM while being considerably more dense

Micron®

# Basics of Caches

Before the reference to $X_n$

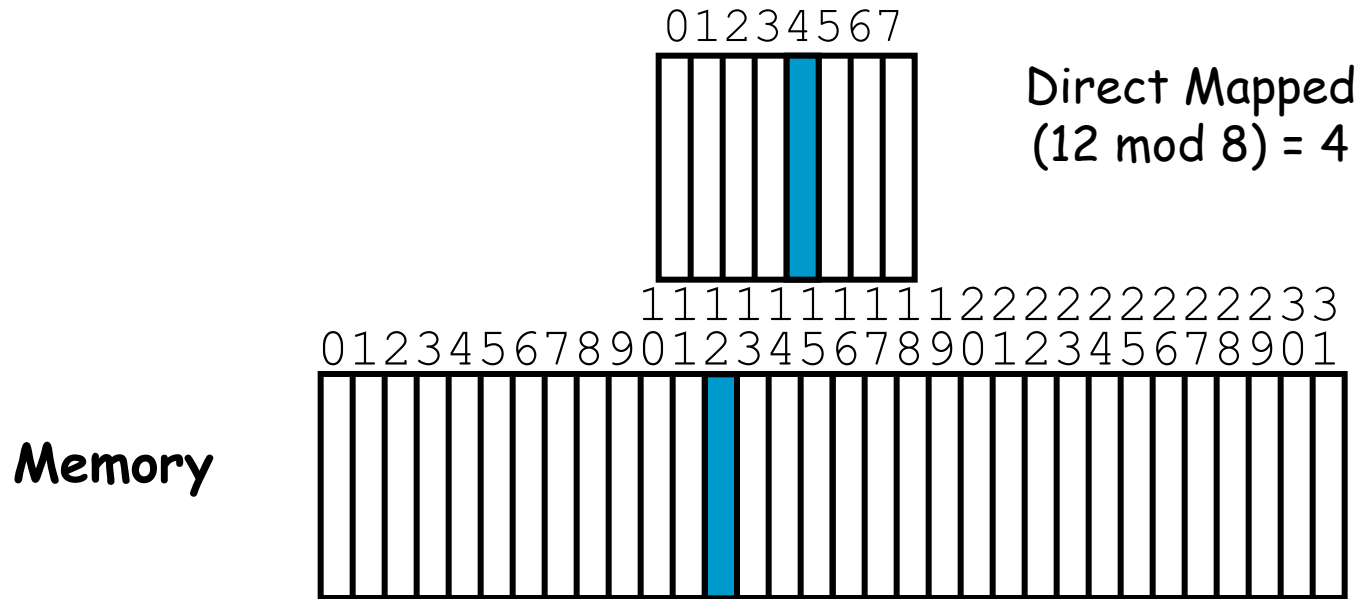| |
|---|
| $X_4$ |
| $X_1$ |
| $X_{n-2}$ |
| |
| $X_{n-1}$ |
| $X_2$ |
| |
| $X_3$ |

After the reference to $X_n$

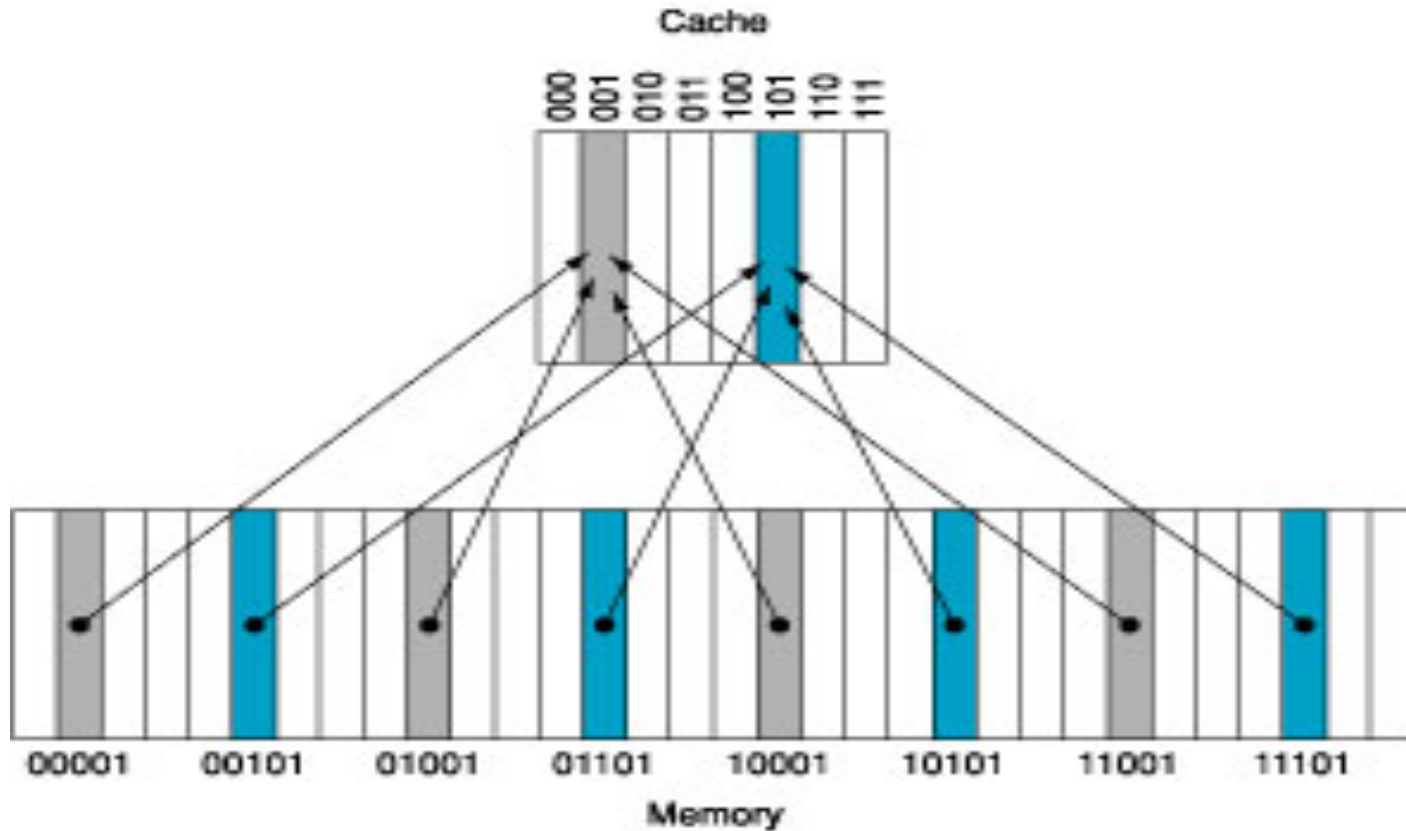| |
|---|
| $X_4$ |
| $X_1$ |
| $X_{n-2}$ |
| |
| $X_{n-1}$ |
| $X_2$ |
| $X_n$ |
| $X_3$ |

How do we find $Xn$ in cache?

# Direct-mapped cache

```
01234567
```

Direct Mapped
(12 mod 8) = 4

```
1111111111222222222233
0123456789012345678901234567890 1
```
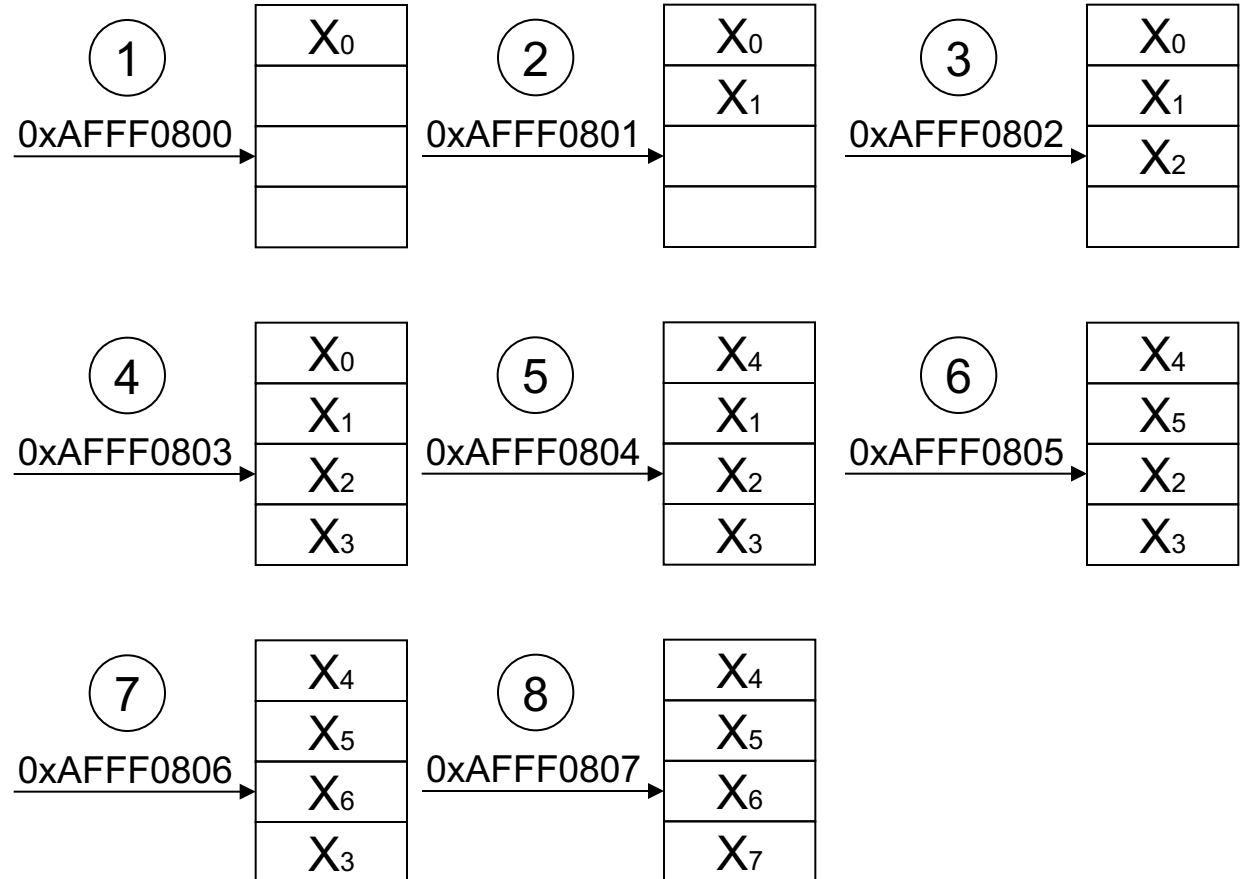
**Memory**

1. Multiple data items map to the same cache location
2. How do we know whether a requested word is in the cache or not?

16

# Example

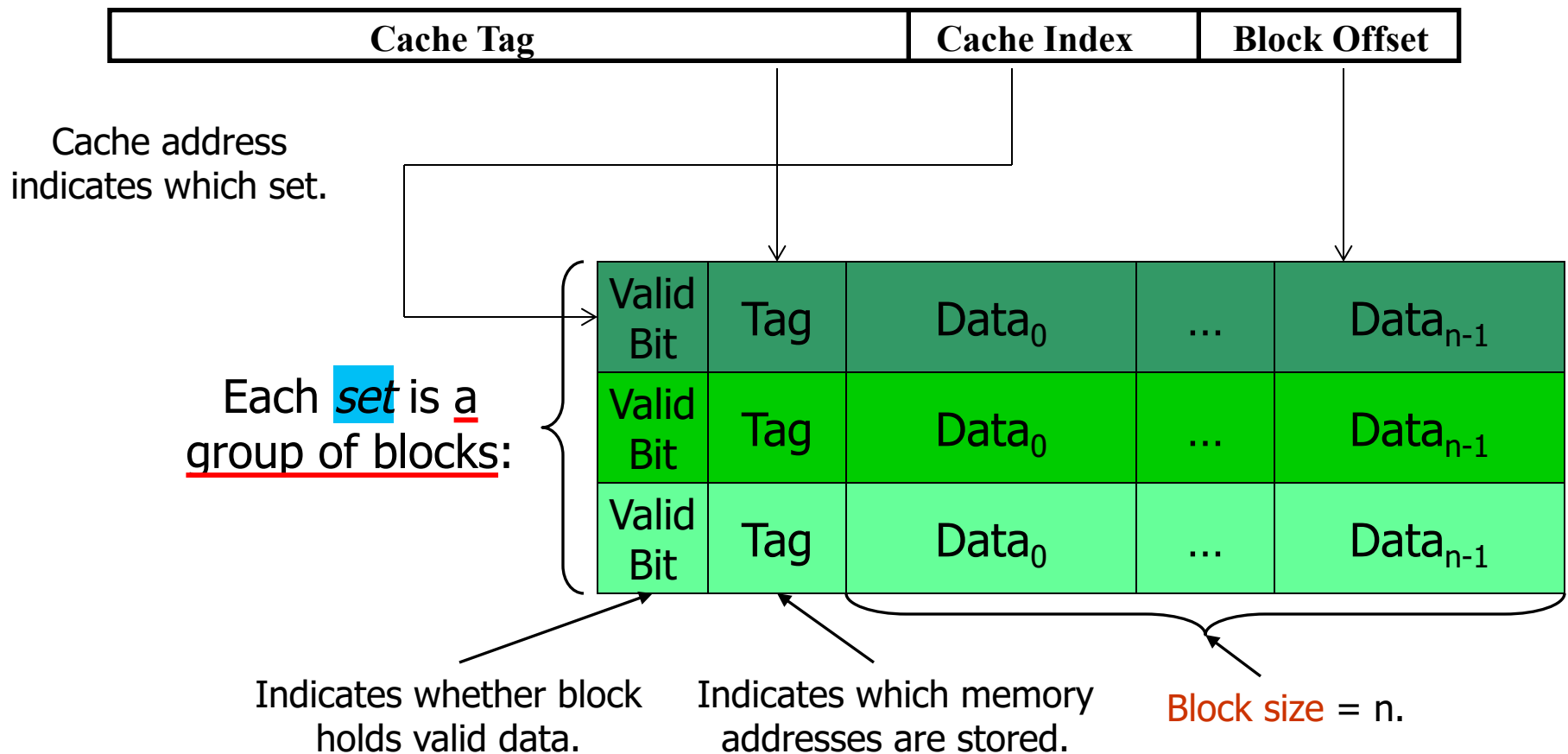| Address | Data |
|---|---|
| 0xAFFF0800 | $X_0$ |
| 0xAFFF0801 | $X_1$ |
| 0xAFFF0802 | $X_2$ |
| 0xAFFF0803 | $X_3$ |
| 0xAFFF0804 | $X_4$ |
| 0xAFFF0805 | $X_5$ |
| 0xAFFF0806 | $X_6$ |
| 0xAFFF0807 | $X_7$ |

1   0xAFFF0800 → 
| $X_0$ |
|---|
| |
| |
| |

2   0xAFFF0801 →
| $X_0$ |
|---|
| $X_1$ |
| |
| |

3   0xAFFF0802 →
| $X_0$ |
|---|
| $X_1$ |
| $X_2$ |
| |

4   0xAFFF0803 →
| $X_0$ |
|---|
| $X_1$ |
| $X_2$ |
| $X_3$ |

5   0xAFFF0804 →
| $X_4$ |
|---|
| $X_1$ |
| $X_2$ |
| $X_3$ |

6   0xAFFF0805 →
| $X_4$ |
|---|
| $X_5$ |
| $X_2$ |
| $X_3$ |

7   0xAFFF0806 →
| $X_4$ |
|---|
| $X_5$ |
| $X_6$ |
| $X_3$ |

8   0xAFFF0807 →
| $X_4$ |
|---|
| $X_5$ |
| $X_6$ |
| $X_7$ |

# How is a block found if it is in the upper level?

## Memory Address

| Cache Tag | Cache Index | Block Offset |
|---|---|---|

Cache address indicates which set.

Each *set* is a group of blocks:

| Valid Bit | Tag | $Data_0$ | ... | $Data_{n-1}$ |
|---|---|---|---|---|
| Valid Bit | Tag | $Data_0$ | ... | $Data_{n-1}$ |
| Valid Bit | Tag | $Data_0$ | ... | $Data_{n-1}$ |

Indicates whether block holds valid data.

Indicates which memory addresses are stored.

Block size = n.

# 1 KB Direct Mapped Cache, 32B blocks

- **How many sets?**

$$\frac{2^{10}}{2^5} = 32$$

| Valid Bit | Cache Tag | | Cache Data | | | | |
|---|---|---|---|---|---|---|---|
| | | | Byte 31 | •• | Byte 1 | Byte 0 | 0 |
| | | | Byte 63 | •• | Byte 33 | Byte 32 | 1 |
| | | | | | | | 2 |
| | | | | | | | 3 |
| ⋮ | ⋮ | | | ⋮ | | | |
| | | | Byte 1023 | •• | | Byte 992 | 31 |

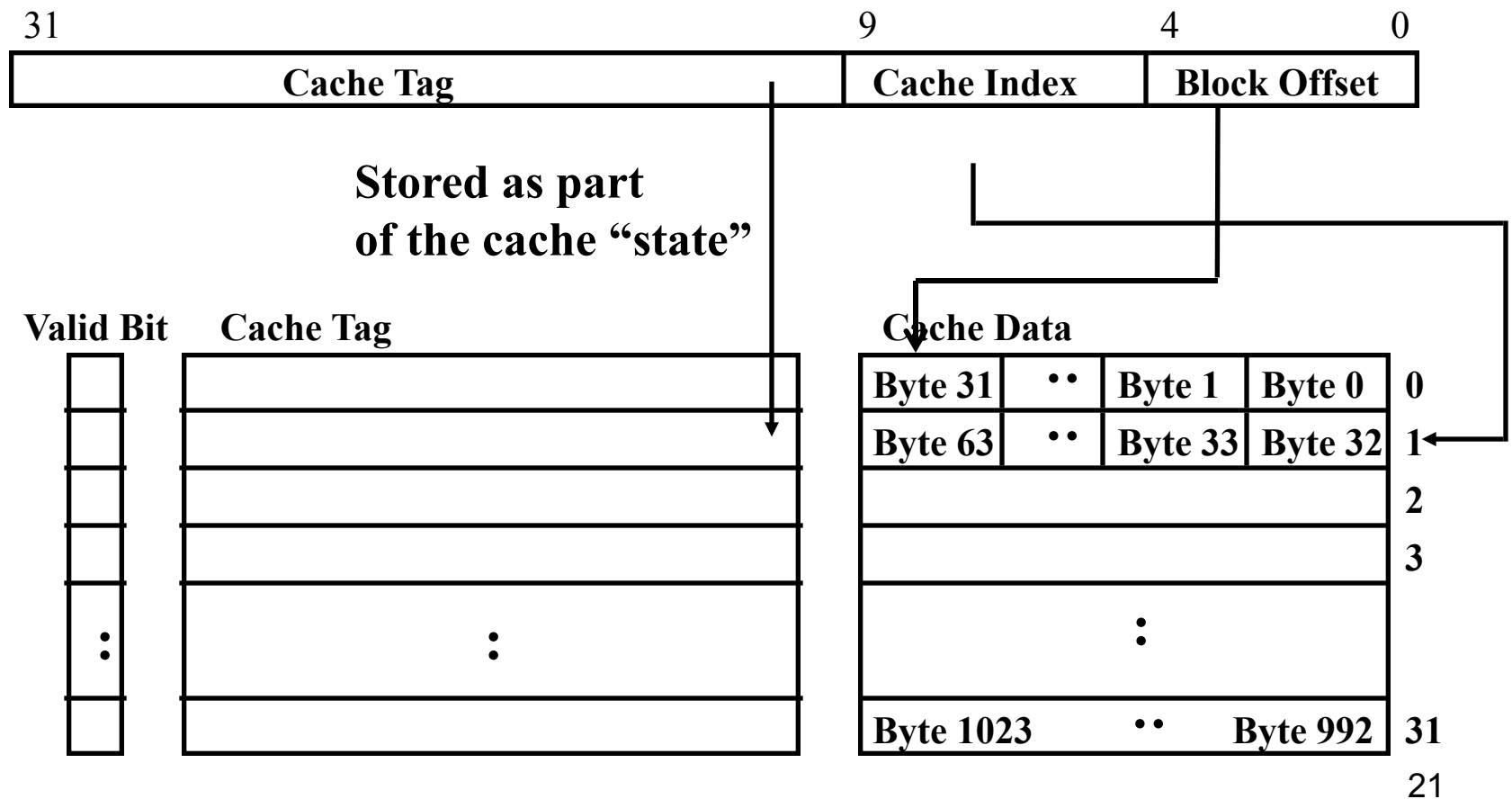# 1 KB Direct Mapped Cache, 32B blocks

■ What is the cache address and tag values for 0x00C0A01F?

```
0000,0000,1100,0000,1010,0000,0001,1111
```
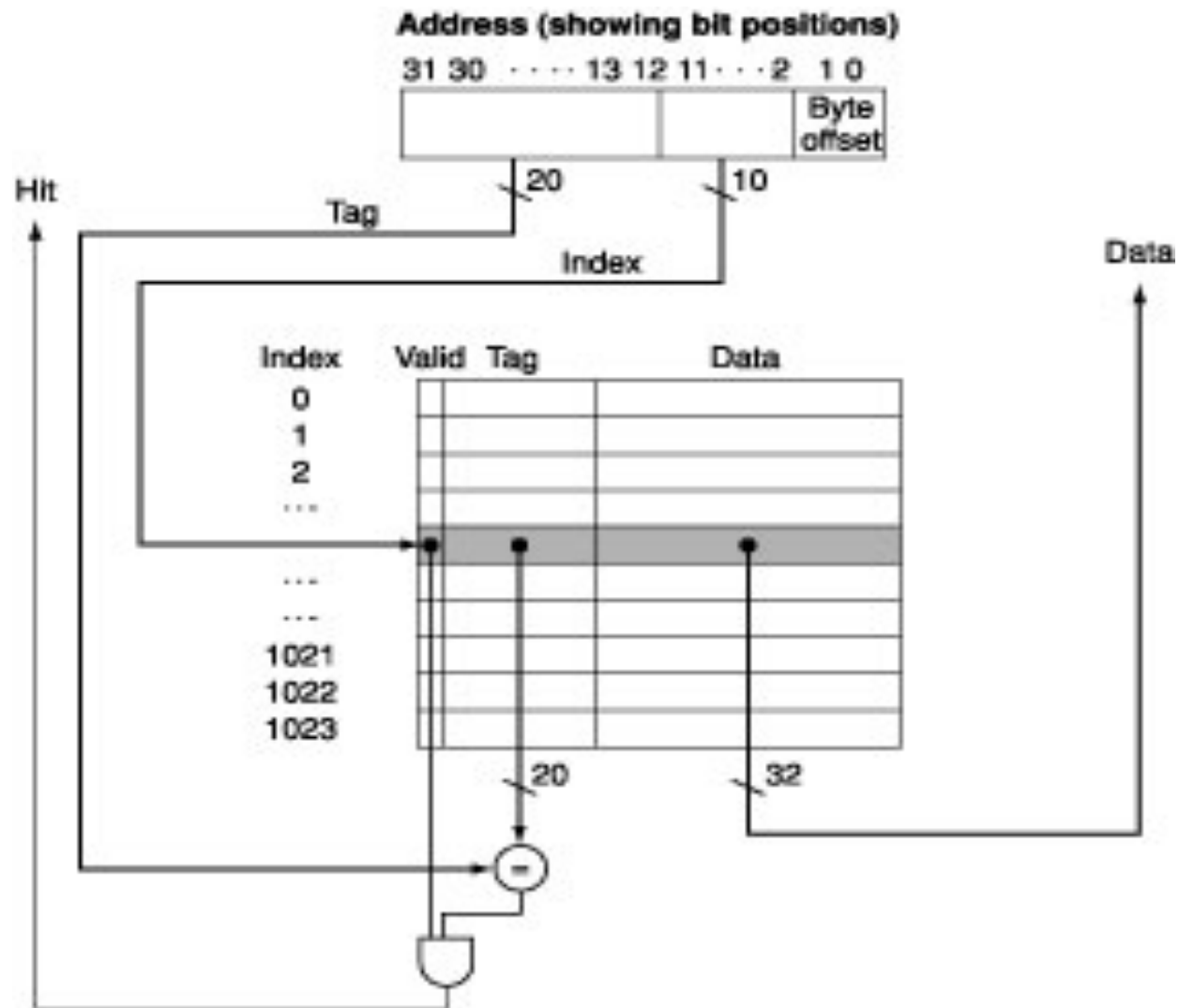
| 31 | | 9 | 4 | 0 |
|---|---|---|---|---|
| | **Cache Tag** | | **Cache Index** | **Block Offset** |

**Stored as part
of the cache "state"**

**Valid Bit**  **Cache Tag**  **Cache Data**

| Byte 31 | •• | Byte 1 | Byte 0 | 0 |
|---|---|---|---|---|
| Byte 63 | •• | Byte 33 | Byte 32 | 1 |
| | | | | 2 |
| | | | | 3 |
| | ⋮ | | | |
| Byte 1023 | •• | | Byte 992 | 31 |

# 1 KB Direct Mapped Cache, 32B blocks

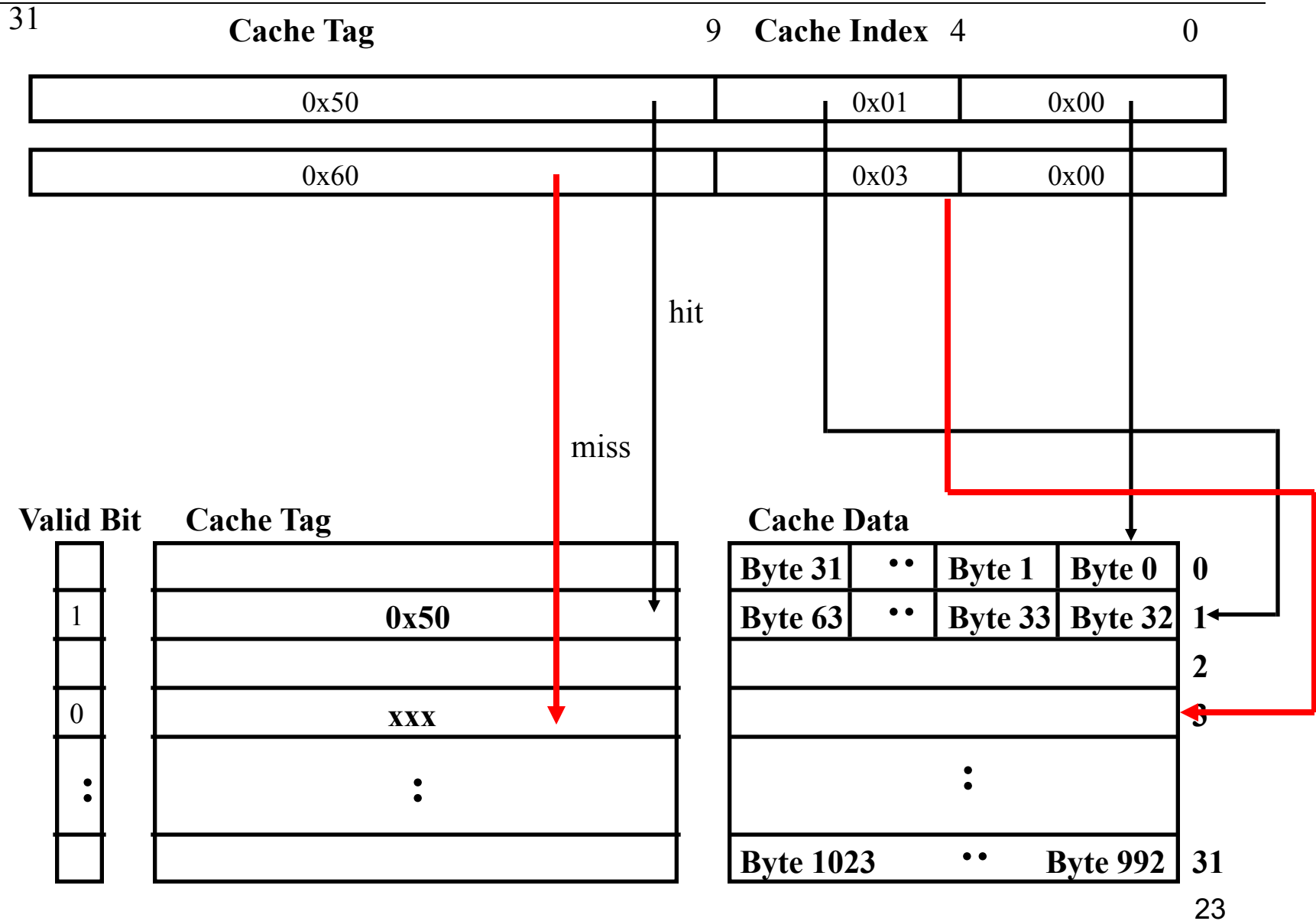■ For a 2 ** N byte direct-mapped cache:
  – The uppermost (32 - N) bits are always the Cache Tag
  – The lowest M bits are the Byte Select (Block Size = 2 ** M)

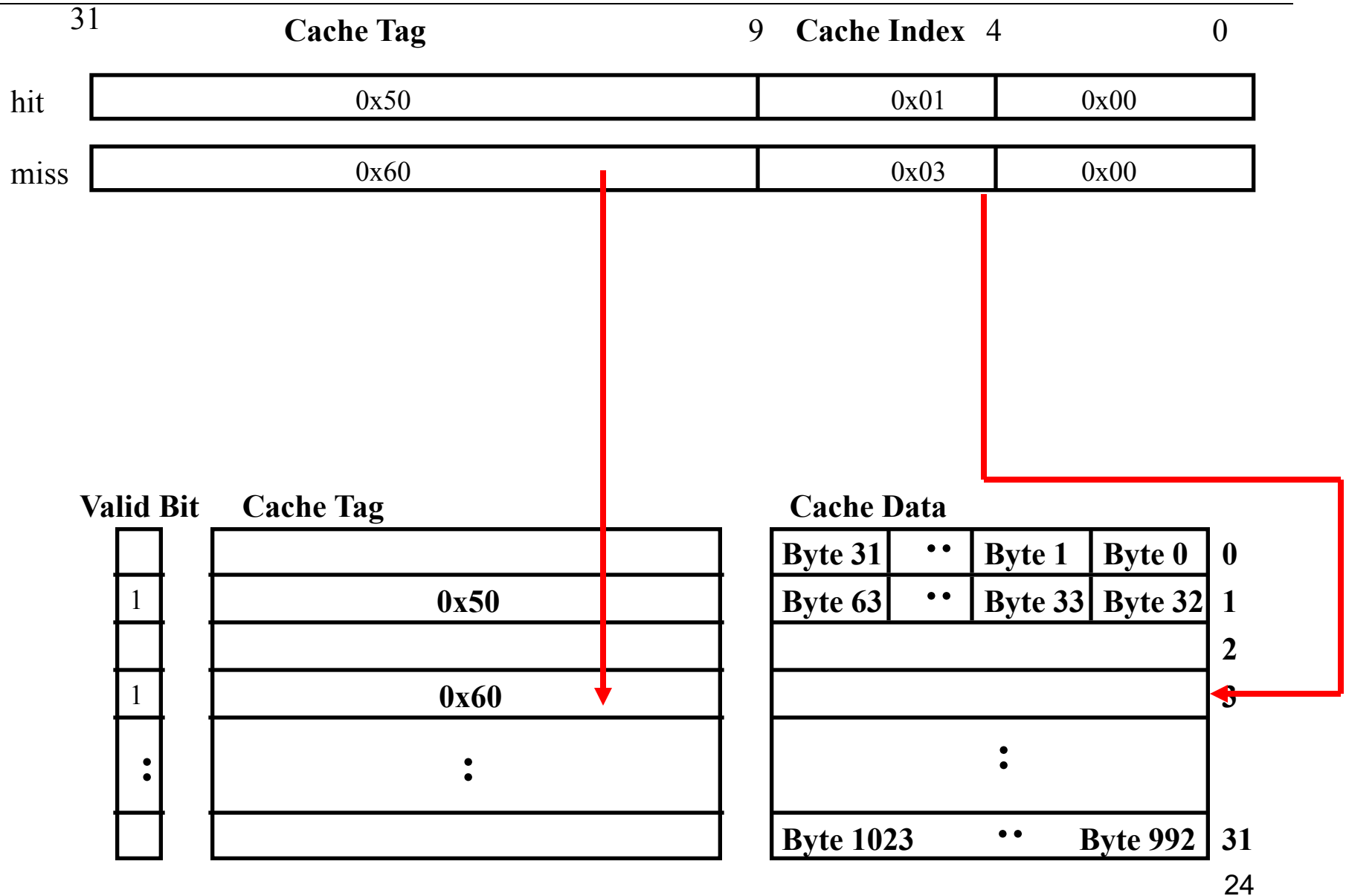| 31 | | 9 | 4 | 0 |
|---|---|---|---|---|
| **Cache Tag** | | **Cache Index** | **Block Offset** | |

**Stored as part of the cache "state"**

**Valid Bit**     **Cache Tag**                              **Cache Data**

| | | Byte 31 | •• | Byte 1 | Byte 0 | 0 |
|---|---|---|---|---|---|---|
| | | Byte 63 | •• | Byte 33 | Byte 32 | 1 |
| | | | | | | 2 |
| | | | | | | 3 |
| ⋮ | ⋮ | | ⋮ | | | |
| | | Byte 1023 | •• | | Byte 992 | 31 |

# Cache Access

# Example: 1 KB Direct Mapped Cache, 32B blocks

| 31 | Cache Tag | 9 | Cache Index | 4 | | 0 |
|---|---|---|---|---|---|---|
| | 0x50 | | 0x01 | | 0x00 | |
| | 0x60 | | 0x03 | | 0x00 | |

hit

miss

**Valid Bit**   **Cache Tag**

**Cache Data**

| Valid Bit | Cache Tag | | Cache Data | | | | |
|---|---|---|---|---|---|---|---|
| | | | Byte 31 | •• | Byte 1 | Byte 0 | 0 |
| 1 | **0x50** | | Byte 63 | •• | Byte 33 | Byte 32 | 1 |
| | | | | | | | 2 |
| 0 | **xxx** | | | | | | 3 |
| : | : | | | : | | | |
| | | | Byte 1023 | •• | | Byte 992 | 31 |

# Example: 1 KB Direct Mapped Cache, 32B blocks

| 31 | **Cache Tag** | 9 | **Cache Index** 4 | 0 |
|---|---|---|---|---|

| | | | | |
|---|---|---|---|---|
| hit | 0x50 | | 0x01 | 0x00 |
| miss | 0x60 | | 0x03 | 0x00 |

**Valid Bit**   **Cache Tag**

| Valid Bit | Cache Tag |
|---|---|
| | |
| 1 | **0x50** |
| | |
| 1 | **0x60** |
| ⋮ | ⋮ |
| | |

**Cache Data**

| Cache Data | | | | |
|---|---|---|---|---|
| **Byte 31** | •• | **Byte 1** | **Byte 0** | **0** |
| **Byte 63** | •• | **Byte 33** | **Byte 32** | **1** |
| | | | | **2** |
| | | | | **3** |
| | ⋮ | | | |
| **Byte 1023** | •• | | **Byte 992** | **31** |

24

# Example: 1 KB Direct Mapped Cache, 32B blocks

| 31 | Cache Tag | 9 | Cache Index 4 | 0 |
|---|---|---|---|---|

| | | | | |
|---|---|---|---|---|
| hit | 0x50 | | 0x01 | 0x00 |

| | | | | |
|---|---|---|---|---|
| miss | 0x60 | | 0x03 | 0x00 |

| | | | | |
|---|---|---|---|---|
| miss | 0x80 | | 0x01 | 0x00 |

miss

| Valid Bit | Cache Tag | | Cache Data | | | | |
|---|---|---|---|---|---|---|---|
| | | | Byte 31 | •• | Byte 1 | Byte 0 | 0 |
| 1 | **0x50** | | Byte 63 | •• | Byte 33 | Byte 32 | 1 |
| | | | | | | | 2 |
| 1 | **0x60** | | | | | | 3 |
| ⋮ | ⋮ | | | ⋮ | | | |
| | | | Byte 1023 | •• | | Byte 992 | 31 |

25

# Example: 1 KB Direct Mapped Cache, 32B blocks

| 31 | Cache Tag | 9 | Cache Index | 4 | 0 |
|---|---|---|---|---|---|

| | Cache Tag | Cache Index | |
|---|---|---|---|
| hit | 0x50 | 0x01 | 0x00 |
| miss | 0x60 | 0x03 | 0x00 |
| miss | 0x80 | 0x01 | 0x00 |

miss

| Valid Bit | Cache Tag | | Cache Data | | | | |
|---|---|---|---|---|---|---|---|
| | | | Byte 31 | •• | Byte 1 | Byte 0 | 0 |
| 1 | **0x80** | | Byte 63 | •• | Byte 33 | Byte 32 | 1 |
| | | | | | | | 2 |
| 1 | **0x60** | | | | | | 3 |
| ⋮ | ⋮ | | | ⋮ | | | |
| | | | Byte 1023 | •• | | Byte 992 | 31 |

26

# Exercise

■ Show the cache contents of an eight-word direct-mapped caches (1-word block size) after each reference for the following address trace (word addressing) :

$10110_{two}$, $11010_{two}$, $10110_{two}$, $110101_{two}$, $10000_{two}$, $00011_{two}$, $10000_{two}$, $100010_{two}$

# Exercise

- How many total bits are required for a direct-mapped cache with 16 KB of data and 4-word blocks, assuming a 32-bit address?

  - # of sets = ?
  - # of data bits for each set =?
  - # of tag bits for each set = ?
  - Valid bit for each set = 1
  - total cache bits = # of set x (valid bit (1-bit) + tag bits + data bits)

# Exercise

- Consider a cache with 64 blocks and a block size of 16 bytes. What block number does byte address 1200 map to?

- Block address = $\lfloor 1200/16 \rfloor$ = 75
- Block number = 75 modulo 64 = 11

| 31 | 10 | 9 | 4 | 3 | 0 |
|---|---|---|---|---|---|
| Tag | | Index | | Offset | |
| 22 bits | | 6 bits | | 4 bits | |

# Block Size

# Block Size (cont.)

- **Advantage of larger block size**
  - take advantage of spatial locality
- **Disadvantage**
  - Too few blocks in cache => high competition
  - Longer cache miss penalty
    - Early restart
      - Resume execution as soon as the requested word of the block is returned
    - Requested word first (critical word first)

requested word → | 1 | } block

# Handling Cache Misses

Cache miss =>

    Stall the entire pipeline & fetch the requested word

Steps to handle an instruction cache miss:

1. Send the original PC value (PC-4) to the memory.

2. Instruct main memory to perform a read and wait for the memory to complete its access.

3. Write the cache entry, putting the data from memory in the data portion of the entry, writing the upper bits of the address (from the ALU) into the tag field, and turning the valid bit on.

4. Restart the instruction execution at the first step, which will refresh the instruction, this time finding it in the cache.

Note that the control of the cache on data access is essentially identical as Instruction access shown above.

# Handling Writes

- *Write through*—The information is written to both the block in the cache and to the block in the lower-level memory.

- *Write back*—The information is written only to the block in the cache. The modified cache block is written to main memory only when it is replaced.
  - is block clean or dirty?

- Pros and Cons of each?
  - WT:
    - Good: read misses cannot result in writes & data coherency
    - Bad: write stall
  - WB:
    - no repeated writes to same location
    - Write new data to cache & write modified block to the lower level of memory hierarchy

# Write Buffer for Write Through



- A Write Buffer is needed between the Cache and Memory
  - Processor: writes data into the cache and the write buffer
  - Memory controller: write contents of the buffer to memory
- Write buffer is just a FIFO:
  - Typical number of entries: 4
  - Works fine if:  Store frequency (w.r.t. time) << 1 / DRAM write cycle
- Memory system designer's nightmare:
  - Store frequency (w.r.t. time)  >  1 / DRAM write cycle
  - Write buffer saturation
- Note: many write-back caches also include write buffers that are used to reduce the miss penalty

# Write Miss Policy

- Why there is a "write miss policy", but no "read miss policy"?

# Write Miss Policy

- ## Write allocate (fetch on write)
  - The block is loaded on a write miss

- ## No-write allocate (write-around)
  - The block is modified in the lower level and not loaded into the cache

|  | Write through | Write back |
|---|---|---|
| Write allocate | hit: write to cache/memory<br>miss: load block into cache; write to cache/memory | hit: write to cache, set dirty bit.<br>miss: load block into cache; write to cache;set dirty bit |
| Write around | hit: write to cache/memory<br>miss: write to memory | hit: write to cache, set dirty bit.<br>miss: write to memory |

# Example: Intrinsity FastMath Processor

- **Intrinsity FastMATH**
  - embedded microprocessor using the MIPS architecture
  - 12-stage pipeline
  - Separate instruction/data caches (split cache), 16 KB, 16-word blocks
  - Offer both write-through and write-back
  - One-entry write buffer.

- **Miss rate of Intrinsity FastMATH for SPEC2000:**
  - Instruction miss rate : 0.4%

  Q: Why instruction miss rate is lower than data miss rate?
  A: Because instruction is much more sequential than data and the spatial location is good.

  - Data miss rate: 11.4%

  - Effective combined miss rate : 3.2%

Q1: Why is data miss rate higher than instruction miss rate?

Q2: Why is combined miss rate is lower than data miss rate?

# Split cache vs. Combined cache

- Combined cache – higher cache hit rate & lower cache bandwidth

- Split cache – lower cache hit rate & higher cache bandwidth

# Example: Intrinsity FastMath Processor (cont.)

■ The 16 KB caches in the Intrinsity FastMATH each contain 256 blocks with 16 words per block.

# Memory Design to Support Cache

CPU

Cache

Bus

Memory

a. One-word-wide
 memory organization

CPU

Multiplexor

Cache

Bus

Memory

b. Wide memory organization

CPU

Cache

Bus

| Memory bank 0 | Memory bank 1 | Memory bank 2 | Memory bank 3 |

c. Interleaved memory organization

access

RAM

request

Cache/CPU

One-word bus

a. One-word-wide
memory organization

access

RAM

request

Cache/CPU

Two-word bus

b. Wide memory organization

access

Bank 0
Bank 1
Bank 2
Bank 3

request

Cache/CPU

One-word bus

c. Interleaved memory
organization

buffer

Assume

    1 memory bus cycle to send the address
    15 memory bus cycles for each DRAM access    15: assumption
    1 memory bus cycle to send a word of data
    <u>4-word block</u> & on-word-wide memory bank

What is the cache miss penalty?

a. One-word-wide memory organization



request

RAM     Cache/CPU

One-word bus

b. Wide memory organization

access

RAM     Cache/CPU

Two-word bus

c. Interleaved memory organization

Bank 0
Bank 1
Bank 2
Bank 3

request

Cache/CPU

One-word bus

Number of access DRAM

Send request     Number of send data back

$1 + 4 \times 15 + 4 \times 1 = 65$

讀4次word，每次花15個cycles

$1 + 2 \times 15 + 2 \times 1 = 33$

讀4words / 2words bus = 2次，
每次花15個cycles

$1 + 1 \times 15 + 4 \times 1 = 20$

# Cache Performance

- CPU time = (CPU execution clock cycles + Memory stall clock cycles) x clock cycle time

- Memory stall clock cycles = Read-stall cycles + write-stall cycles

- Read-stall cycles = # of Read  x   Read miss rate X Read miss penalty

- Write-stall cycles = ( # of Writes X Write miss rate X Write miss penalty) +

  <mark>Write buffer stalls.</mark>

- Memory-stall clock cycles =  # Memory accesses  X Miss rate  X Miss penalty

$$= \frac{\text{Instructions}}{\text{program}} \text{ x } \frac{\text{Misses}}{\text{Instruction}} \text{ x  Miss penalty}$$

- Average memory access time = Hit time + Miss rate x Miss penalty

# Example

- I-Cache miss rate = 2% & D-Cache miss rate = 4%
- Base CPI 2.0
- Miss penalty = 100 cycles
- Frequency of loads and stores is 36%.
- Compare the performance with a perfect cache

- I-cache stall cycles = I x 2% x 100 = 2.00 X I
- D-cache stall cycles = I X 36% X 4% X 100 = 1.44 X I
- CPU time with stalls = I X (2 + 1.44 + 2) X Clock-cycle

$$= I \ X \ 5.44 \ X \ Clock\text{-}cycle$$

- CPU time with perfect caches

$$= I \ X \ CPI_{perfect} \ X \ Clock\text{-}cycle$$
$$= I \ \ X \ 2 \ X \ Clock\text{-}cycle$$

單位(S)

CPU time = cycles x clock-cycle

= (CPU cycles + memory stall cycles) x clock-cycle

# Cache Performance with Increased Clock rate

- How much faster will the computer be with 2x clock rate, assuming the same miss rate as the previous example.

Stall time remains same!

- I-cache stall cycles = I x 2% x 200 = 4.00 X I

- D-cache stall cycles = I X 36% X 4% X 200 = 2.88 X I

- CPU time with stalls = I X (4 + 2.88 + 2) X Clock-cycle

    = I X 8.88 X Clock-cycle

$$\frac{Performance\ with\ fast\ clock}{Performance\ with\ slow\ clock} = \frac{Execution\ time\ with\ slow\ clock}{Execution\ time\ with\ fast\ clock}$$

$$= \frac{IC \times CPI_{slow\ clock} \times Clock\ Cycle}{IC \times CPI_{fast\ clock} \times \dfrac{Clock\ Cycle}{2}}$$

$$= \frac{5.44}{8.88 \times \dfrac{1}{2}} = 1.23$$

44

# How to Improve Cache Performance?

1. Reduce miss rate  -> Increasing associativity
2. Reduce miss penalty -> multi-level cache
3. Reduce hit time -> small cache

**Average memory access time = hit time + miss-rate X miss-penalty**

# Reducing Cache Misses by More Flexible Placement of Blocks

if 4-Way, there are 2 sets and each with 4 blocks

Full Mapped

Direct Mapped
(12 mod 8) = 4

2-Way Assoc
(12 mod 4) = 0

0 1 2 3 4 5 6 7

0 1 2 3 4 5 6 7

0 1 2 3 4 5 6 7

**Cache**

set 0  set 1  set 2  set 3

Search

Search

Search

1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 3 3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1

**Memory**

# Possible Associativity Structures

(direct mapped)

| Block | Tag | Data |
|---|---|---|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |

Two-way set associative

| Set | Tag | Data | Tag | Data |
|---|---|---|---|---|
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |

Four-way set associative

| Set | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|---|---|---|---|---|---|---|---|---|
| 0 | | | | | | | | |
| 1 | | | | | | | | |

Eight-way set associative (fully associative)

| Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | | | |

## An 8-block cache

```
float dot_prod(float x[SIZE],
               float y[SIZE])
{
  float sum = 0.0;
  int   i;

  for (i = 0; i < SIZE; i++)
    sum += x[i]*y[i];

  return sum;
}
```

Assume:

Direct-mapped cache.

`x[i]` and `y[i]` map to same blocks.

? What is the hit rate? ?

Under these assumptions, every access is a cache miss.
Hit rate = 0%.

What can we do?  Increasing the associativity

# Exercise

- **Three small caches, each consisting of four one-word blocks**
  - Direct mapped cache
  - Two-way set associative cache
  - Fully associative cache

- **Find the number of misses for each cache for the following sequence**
  - 0,8,0,6,8

# Exercise (cont.)

- The direct mapped cache

| Block Address | Cache Set |
|---------------|-----------|
| 0 | (0 modulo 4) = 0 |
| 6 | (6 modulo 4) = 2 |
| 8 | (8 modulo 4) = 0 |

| Address of memory block accessed | Hit or miss | Contents of cache blocks after reference | | | |
|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 |
| 0 | Miss | Memory[0] | | | |
| 8 | Miss | Memory[8] | | | |
| 0 | Miss | Memory[0] | | | |
| 6 | Miss | Memory[0] | | Memory[6] | |
| 8 | miss | Memory[8] | | Memory[6] | |

# Exercise (cont.)

- The two-way set associative cache

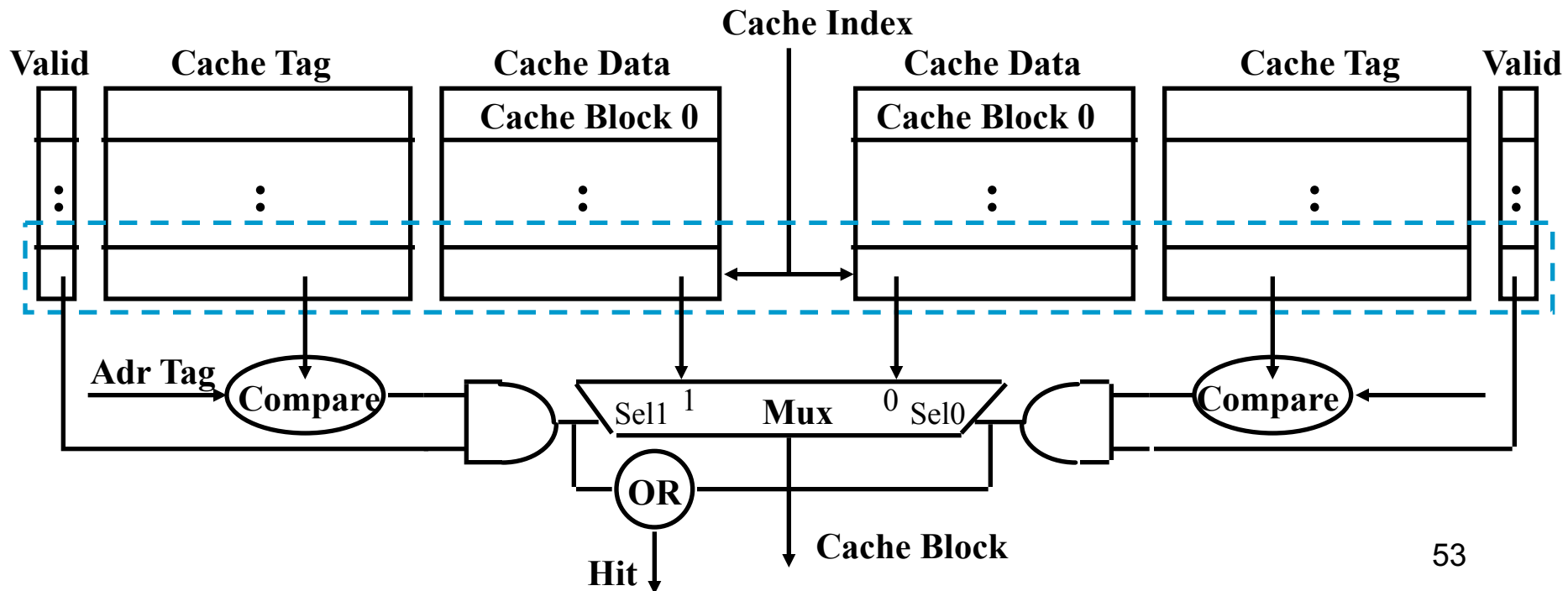| Block Address | Cache Set |
|---|---|
| 0 | (0 modulo 2) = 0 |
| 6 | (6 modulo 2) = 0 |
| 8 | (8 modulo 2) = 0 |

| Address of memory block accessed | Hit or miss | Contents of cache blocks after reference | | | |
|---|---|---|---|---|---|
| | | Set 0 | Set 0 | Set 1 | Set 1 |
| 0 | Miss | Memory[0] | | | |
| 8 | Miss | Memory[0] | Memory[8] | | |
| 0 | Hit | Memory[0] | Memory[8] | | |
| 6 | Miss | Memory[0] | Memory[6] | | |
| 8 | miss | Memory[8] | Memory[6] | | |

replacement policy: 6進來時，應該要踢掉誰？

# Exercise (cont.)

- The fully associative cache

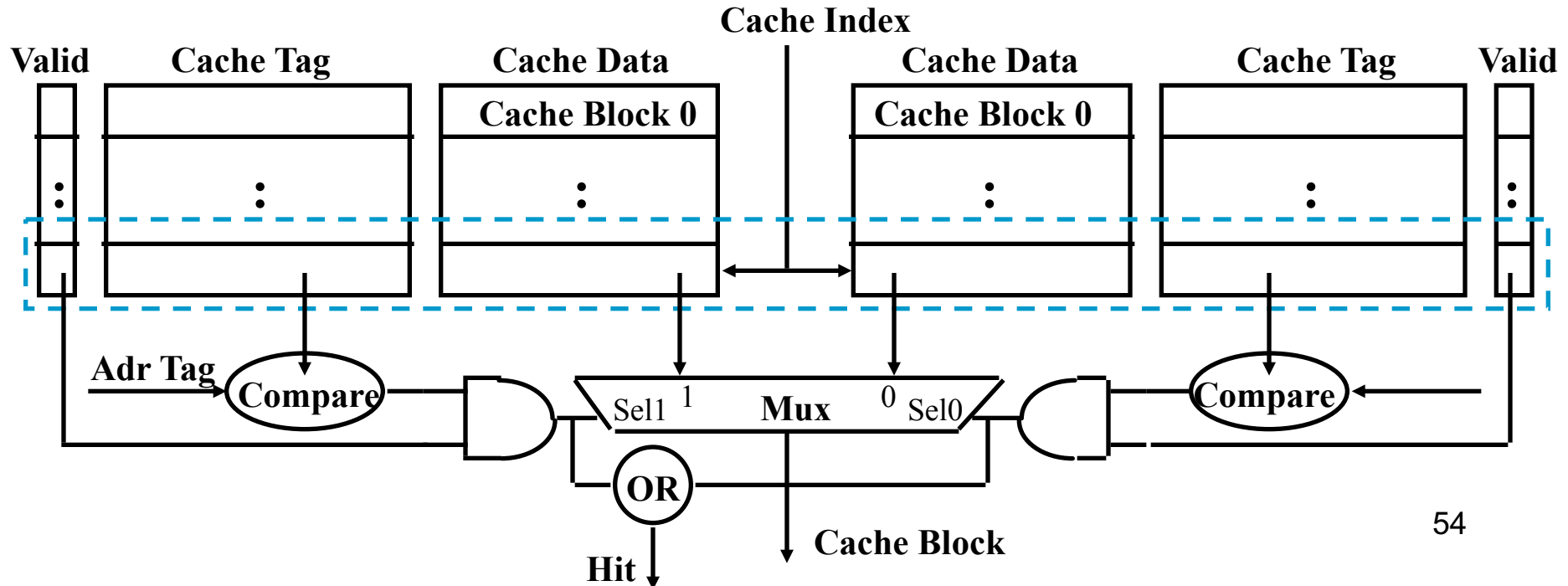| Address of memory block accessed | Hit or miss | Contents of cache blocks after reference | | | |
|---|---|---|---|---|---|
| | | Block 0 | Block 1 | Block 2 | Block 3 |
| 0 | Miss | Memory[0] | | | |
| 8 | Miss | Memory[0] | Memory[8] | | |
| 0 | Hit | Memory[0] | Memory[8] | | |
| 6 | Miss | Memory[0] | Memory[8] | Memory[6] | |
| 8 | Hit | Memory[0] | Memory[8] | Memory[6] | |

# 2-way Set Associative Cache

- 1K, 2-way,32B block
  - 16 sets, each set has 2 blocks
  - Index = 4 bits, tag = 23 bits
- Increasing associativity shrinks index, expands tag
- How to find a data in 2-way cache
  - Cache Index selects a "set" from the cache
  - The two tags in the set are compared in parallel
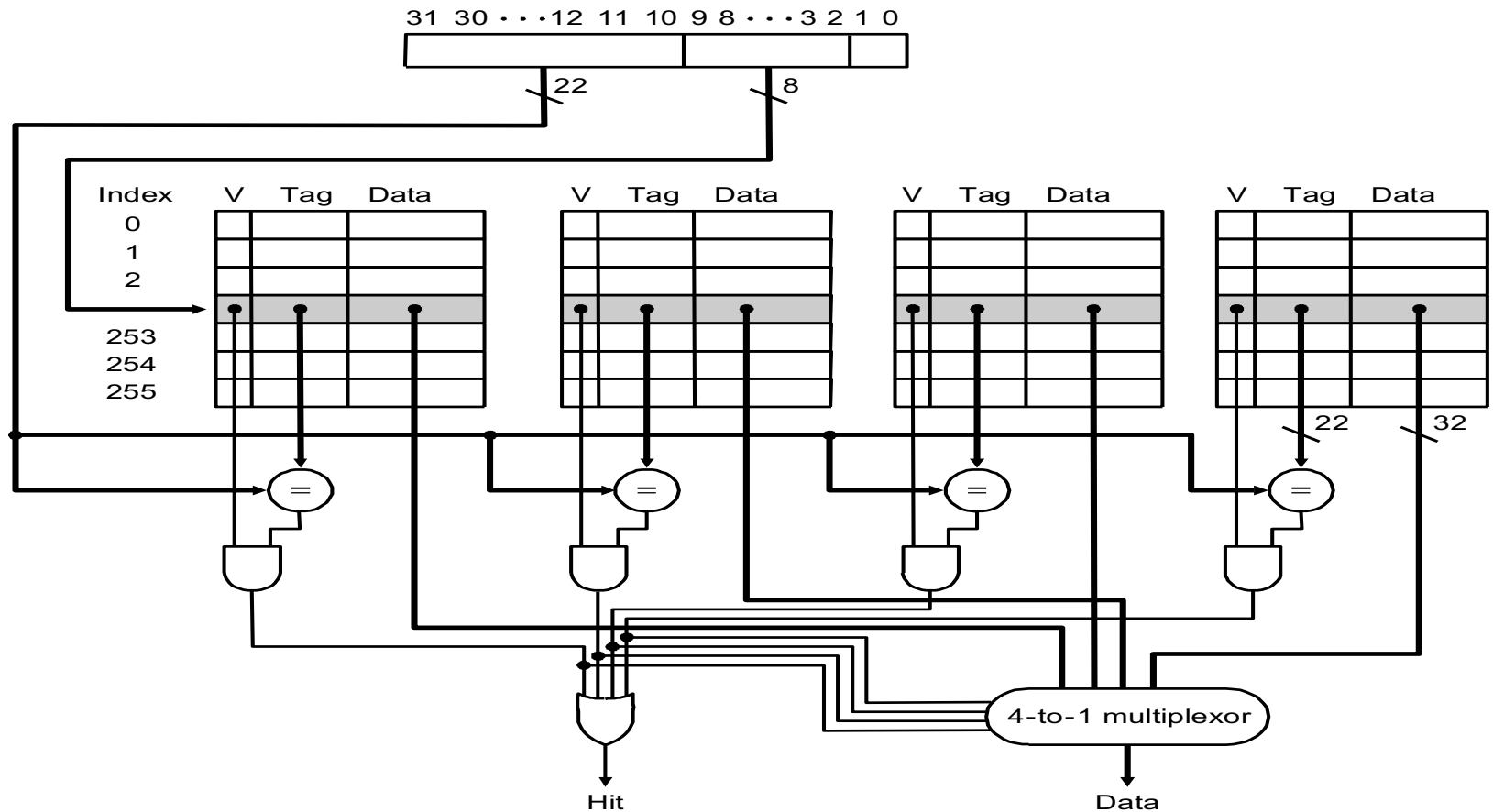  - Data is selected based on the tag result

# Disadvantage of Set Associative Cache

- **N-way Set Associative Cache vs. Direct Mapped Cache:**
  - N comparators vs. 1
  - Extra MUX delay for the data
  - Data comes AFTER Hit/Miss
- **In a direct mapped cache, cache block is available BEFORE Hit/Miss:**
  - Possible to assume a hit and continue.  Recover later if miss.

**Cache Index**

| Valid | Cache Tag | Cache Data | | Cache Data | Cache Tag | Valid |
|---|---|---|---|---|---|---|

Cache Block 0

Cache Block 0

**Adr Tag**

**Compare**

Sel1  1  **Mux**  0  Sel0

**Compare**

**OR**

**Hit**

**Cache Block**

# A 4-Way Set-Associative Cache

# Replacement Policy: Choosing Which Block to Replace

- **Easy for direct mapped**

- **Set associative or fully associative:**
    - Random
    - FIFO
    - LRU (Least Recently Used):
        - Hardware keeps track of the access history and replace the block that has not been used for the longest time

# Effects of Associativity

| Associativity | Data miss rate |
|:---:|:---:|
| 1 direct-map | 10.3% |
| 2 way | 8.6% |
| 4 | 8.3% |
| 8 | 8.1% |

Data Source: Spec2000  benchmarks

# Tag Size vs. Associtivity

■ With the same cache capacity, increasing associtivity increases or decreasing tag bits?

(direct mapped)

| Block | Tag | Data |
|-------|-----|------|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |

Two-way set associative

| Set | Tag | Data | Tag | Data |
|-----|-----|------|-----|------|
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |

Four-way set associative

| Set | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|-----|-----|------|-----|------|-----|------|-----|------|
| 0 | | | | | | | | |
| 1 | | | | | | | | |

Eight-way set associative (fully associative)

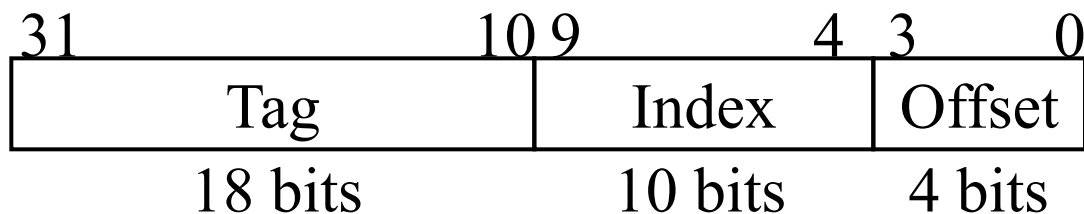| Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data | Tag | Data |
|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|-----|------|
| | | | | | | | | | | | | | | | |

# Exercise: Size of Tags vs. Set Associativity

- Find # of sets, total # of tag bits for direct-mapped, two-way and four way.  (Capacity- 4K blocks, four-word block, 32 bit address)
  - tag+index = 32- 4 = 28
  - The direct mapped cache
    - 4K sets
    - $\log_2$(4K) = 12 bits of index
    - Total (28-12) $\times$ 4K = 16 $\times$ 4K = 64 Kbits of tag
  - The two-way set-associative cache
    - 2K sets
    - Total (28-11) $\times$ 2 $\times$ 2K = 34 $\times$ 2K = 68 Kbits of tag
  - The four-way set-associative cache
    - 1K sets
    - Total (28-10) $\times$ 4 $\times$ 1K = 72 $\times$ 1K = 72 Kbits of tag
  - The fully associative cache
    - One set with 4K blocks
    - Tag is 28 bits
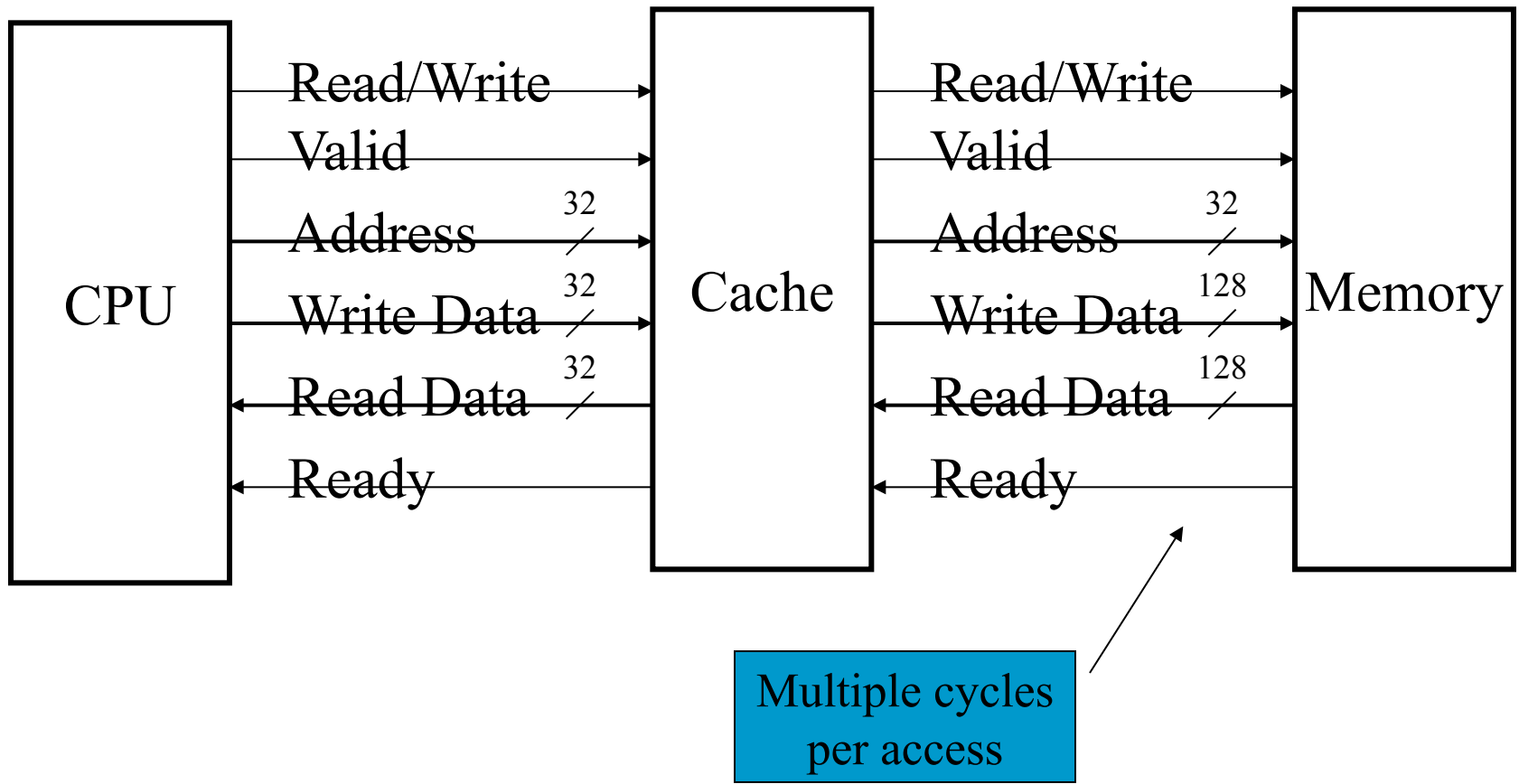    - Total 28 $\times$ 4K $\times$ 1 = 112K tag bits

# Cache Control

- **Example cache characteristics**
  - Direct-mapped, write-back, write allocate
  - Block size: 4 words (16 bytes)
  - Cache size: 16 KB (1024 blocks)
  - 32-bit byte addresses
  - Valid bit and dirty bit per block
  - Blocking cache 遇到cache miss就stall
    - CPU waits until access is complete
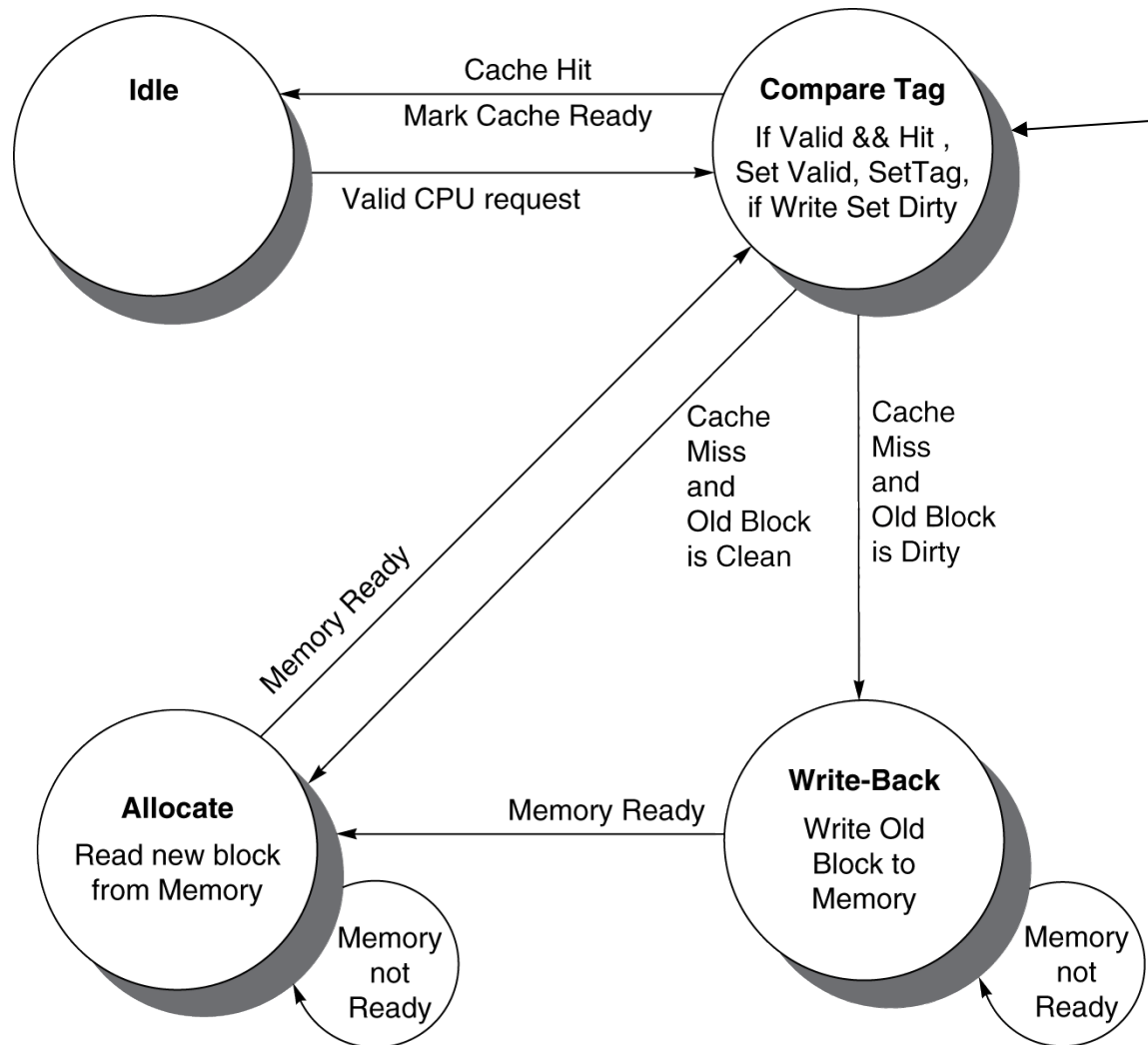
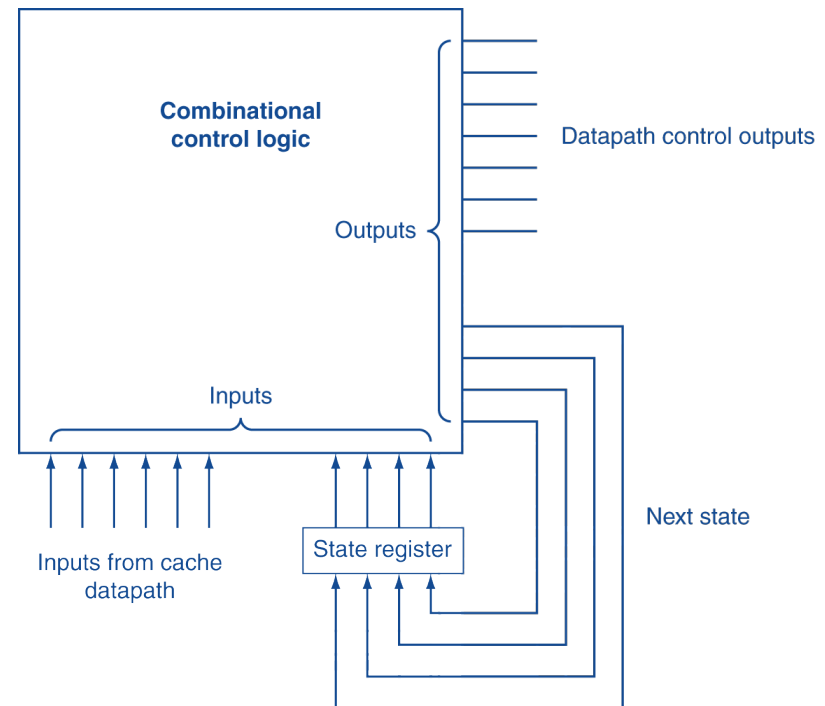| 31 | | 10 | 9 | | 4 | 3 | | 0 |
|---|---|---|---|---|---|---|---|---|
| | Tag | | | Index | | | Offset | |
| | 18 bits | | | 10 bits | | | 4 bits | |

# Interface Signals



CPU

Read/Write →

Valid →

Address $\overset{32}{/}$ →

Write Data $\overset{32}{/}$ →

Read Data $\overset{32}{/}$ ←

Ready ←

Cache

Read/Write →

Valid →

Address $\overset{32}{/}$ →

Write Data $\overset{128}{/}$ →

Read Data $\overset{128}{/}$ ←

Ready ←

Memory

Multiple cycles per access

# Cache Controller FSM

**Idle**

Cache Hit
Mark Cache Ready

Valid CPU request

**Compare Tag**
If Valid && Hit ,
Set Valid, SetTag,
if Write Set Dirty

Could partition
into separate
states to reduce
clock cycle
time

Cache
Miss
and
Old Block
is Clean

Cache
Miss
and
Old Block
is Dirty

Memory Ready

**Allocate**
Read new block
from Memory

Memory Ready

**Write-Back**
Write Old
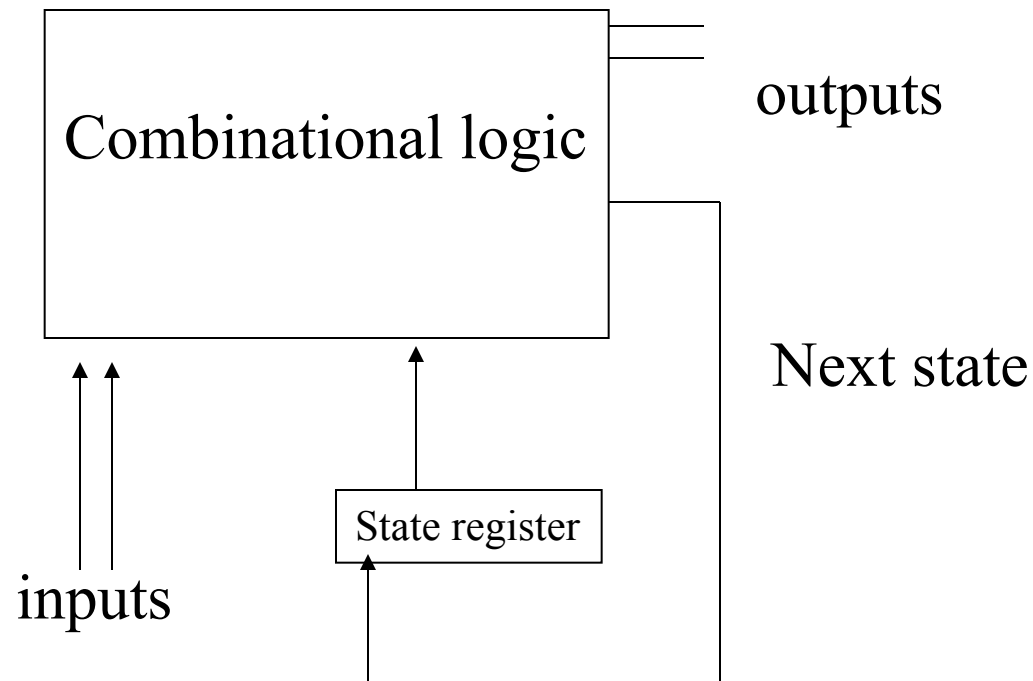Block to
Memory

Memory
not
Ready

Memory
not
Ready

# Finite State Machines

■ **Use a FSM to sequence control steps**

■ **Set of states, transition on each clock edge**
  – State values are binary encoded
  – Current state stored in a register
  – Next state $= f_n$ (current state, current inputs)

■ **Control output signals $= f_o$ (current state)**

# Finite State Implementation



How to implement the combinational logic?
- PLA (Programmable Logic Array)
- ROMs

# Miss Penalty Reduction: Multi- Level Cache

- **Larger cache vs. CPU time**
  - Add another level of cache
  - The L2 cache is much larger than L1

- **L2 Equations**

  $AMAT = \text{Hit Time}_{L1} + \text{Miss Rate}_{L1} \times \text{Miss Penalty}_{L1}$

  $\text{Miss Penalty}_{L1} = \text{Hit Time}_{L2} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2}$

  $AMAT = \text{Hit Time}_{L1} + \text{Miss Rate}_{L1} \times (\text{Hit Time}_{L2} + \text{Miss Rate}_{L2} \times \text{Miss Penalty}_{L2})$

- **Definitions:**
  - *Local miss rate*— misses in this cache divided by the total number of memory accesses *to this cache* (Miss rate$_{L2}$)
  - *Global miss rate*—misses in this cache divided by the total number of memory accesses *generated by the CPU* (Miss Rate$_{L1}$ x Miss Rate$_{L2}$)

L1

L2

Main Memory

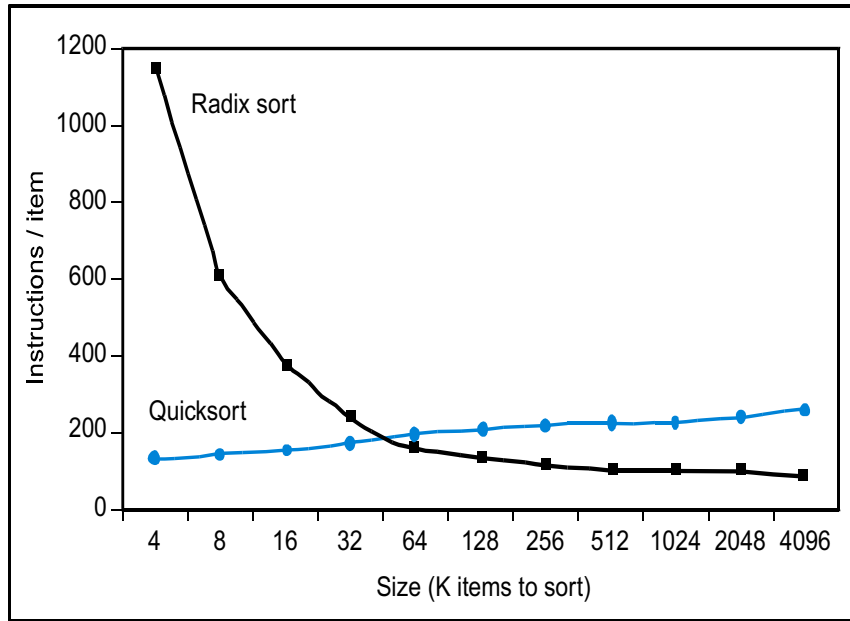# Performance of Multilevel Caches

- Suppose we have a processor with the following parameters:
  - Base CPI=1.0, if hit in the L1 cache. L1 cache miss rate is 2%. Clock rate is 5 GHz.  Memory access time is 100 ns, including all the miss handling.
  - Miss penalty to main memory is 100ns/0.2 = 500 clock cycle.
  - Total CPI = 1.0 + 2%×500 = 11.0

- If we add a L2 cache that has 5 ns access time. L2 global miss rate = 0.5%
  - Miss penalty to L2 is 5ns/0.2 = 25 clock cycle
  - Total CPI = 1.0+ 2%×25 + 0.5%×500 = 4.0

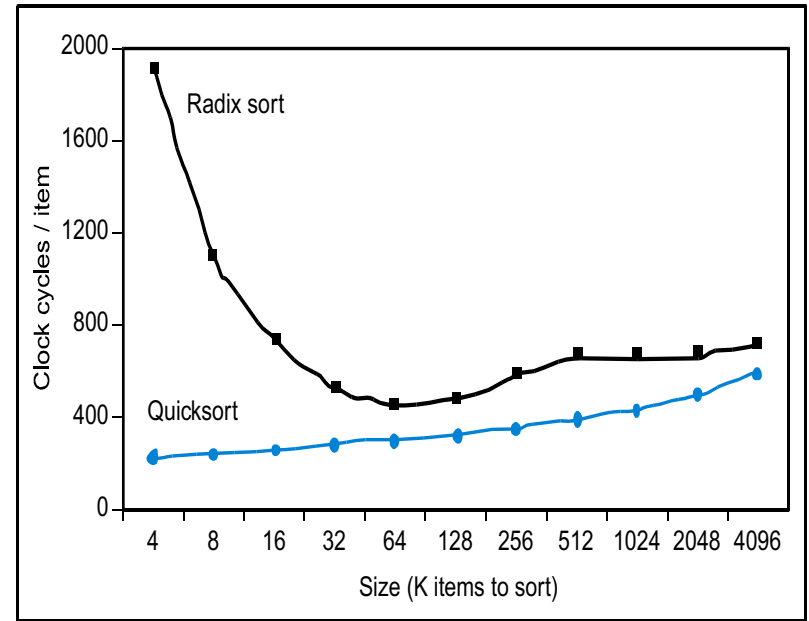# L2 Cache Design Principle

- **L2 not tied to CPU clock cycle**
  - Different design consideration from L1
  - Hits are less important than misses
  - Larger cache, higher associativity and larger blocks
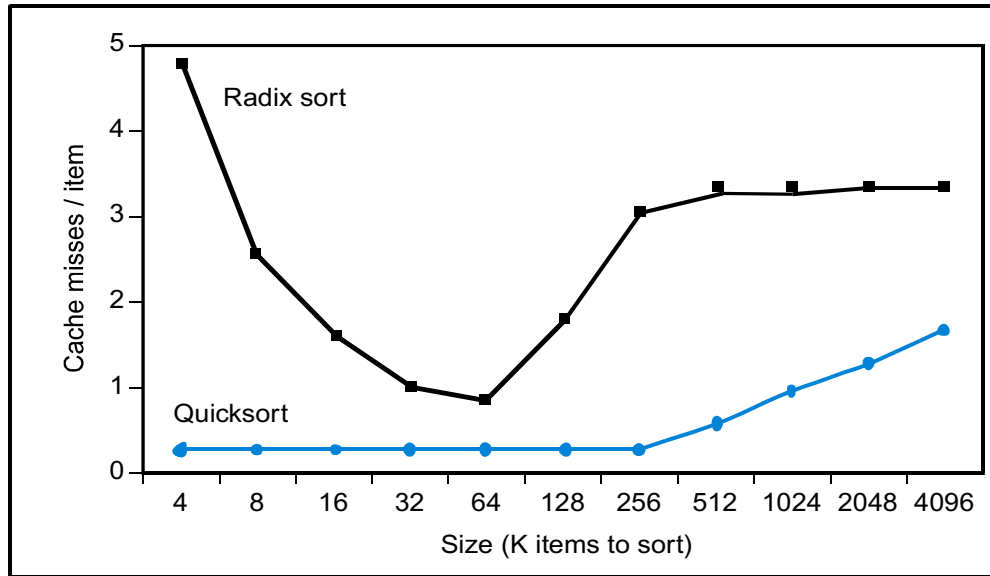
# Caches vs. Performance

included memory performance!



Theoretical behavior of
Radix sort vs. Quicksort
(instruction/item)

Observed behavior of
Radix sort vs. Quicksort
(clock cycles/item)

# Caches vs. Performance (cont.)



Cache behavior
Radix sort vs. Quicksort
(cache misses /item)

- **Memory system performance is often critical factor**
  - Multilevel caches, pipelined processors, make it harder to predict outcomes
  - Need experimental data