# Gerarchia di Memoria

# Review:  Major Components of a Computer

Processor

Control

Datapath

Memory

Devices

Input

Output

Cache

Main Memory

Secondary Memory (Disk)

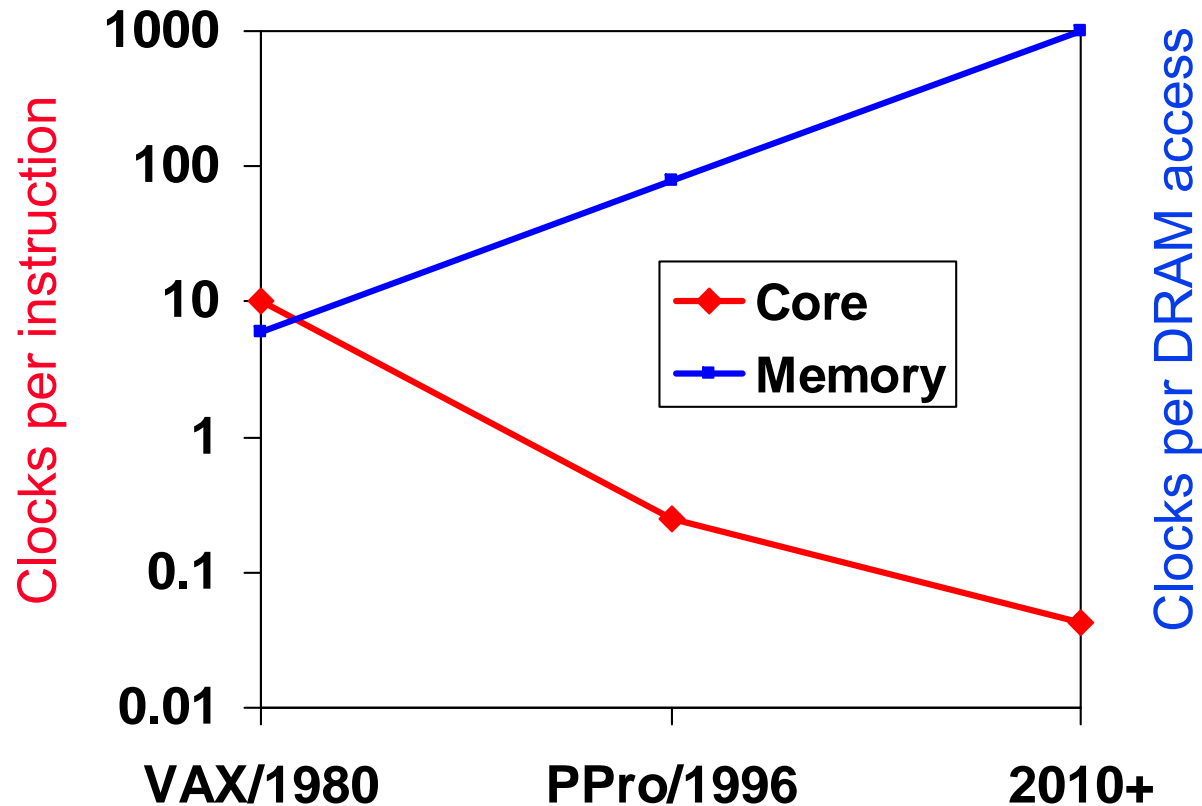# Processor-Memory Performance Gap

# The "Memory Wall"
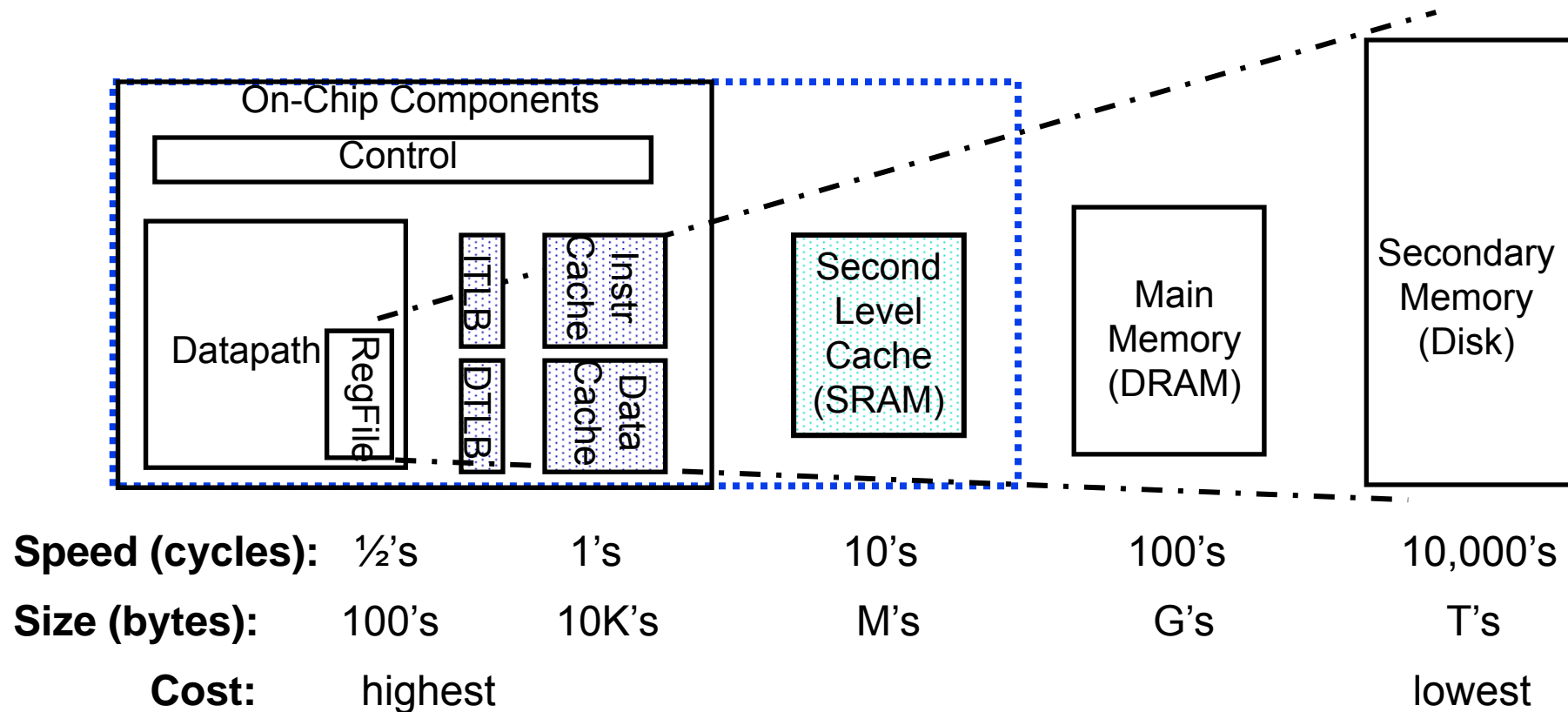
❑ Processor vs DRAM speed disparity continues to grow



❑ Good memory hierarchy (cache) design is increasingly important to overall performance

# The Memory Hierarchy Goal

❑ Fact:  Large memories are slow and fast memories are small

❑ How do we create a memory that gives the illusion of being large, cheap and fast (most of the time)?

  ● With hierarchy
  ● With parallelism

# A Typical Memory Hierarchy

❑ Take advantage of the principle of locality to present the user with as much memory as is available in the *cheapest* technology at the speed offered by the *fastest* technology



| | On-Chip Components | | | Second Level Cache (SRAM) | Main Memory (DRAM) | Secondary Memory (Disk) |
|---|---|---|---|---|---|---|

| **Speed (cycles):** | ½'s | 1's | 10's | 100's | 10,000's |
|---|---|---|---|---|---|
| **Size (bytes):** | 100's | 10K's | M's | G's | T's |
| **Cost:** | highest | | | | lowest |

# Memory Hierarchy Technologies

❑ Caches use *SRAM* for speed and technology compatibility

- Fast (typical access times of 0.5 to 2.5 nsec)
- Low density (6 transistor cells), higher power, expensive
- Static: content will last "forever" (as long as power is left on)

❑ Main memory uses *DRAM* for size (density)

- Slower (typical access times of 50 to 70 nsec)
- High density (1 transistor cells), lower power, cheaper
- Dynamic: needs to be "refreshed" regularly (e.g. ~ every 8 ms)
  - consumes1% to 2% of the active cycles of the DRAM
- Addresses divided into 2 halves (row and column)
  - *RAS* or *Row Access Strobe* triggering the row decoder
  - *CAS* or *Column Access Strobe* triggering the column selector

# The Memory Hierarchy:  Why Does it Work?

❑ Temporal Locality (locality in time)

- If a memory location is referenced then it will tend to be referenced again soon

⇒ Keep most recently accessed data items closer to the processor

❑ Spatial Locality (locality in space)

- If a memory location is referenced, the locations with nearby addresses will tend to be referenced soon

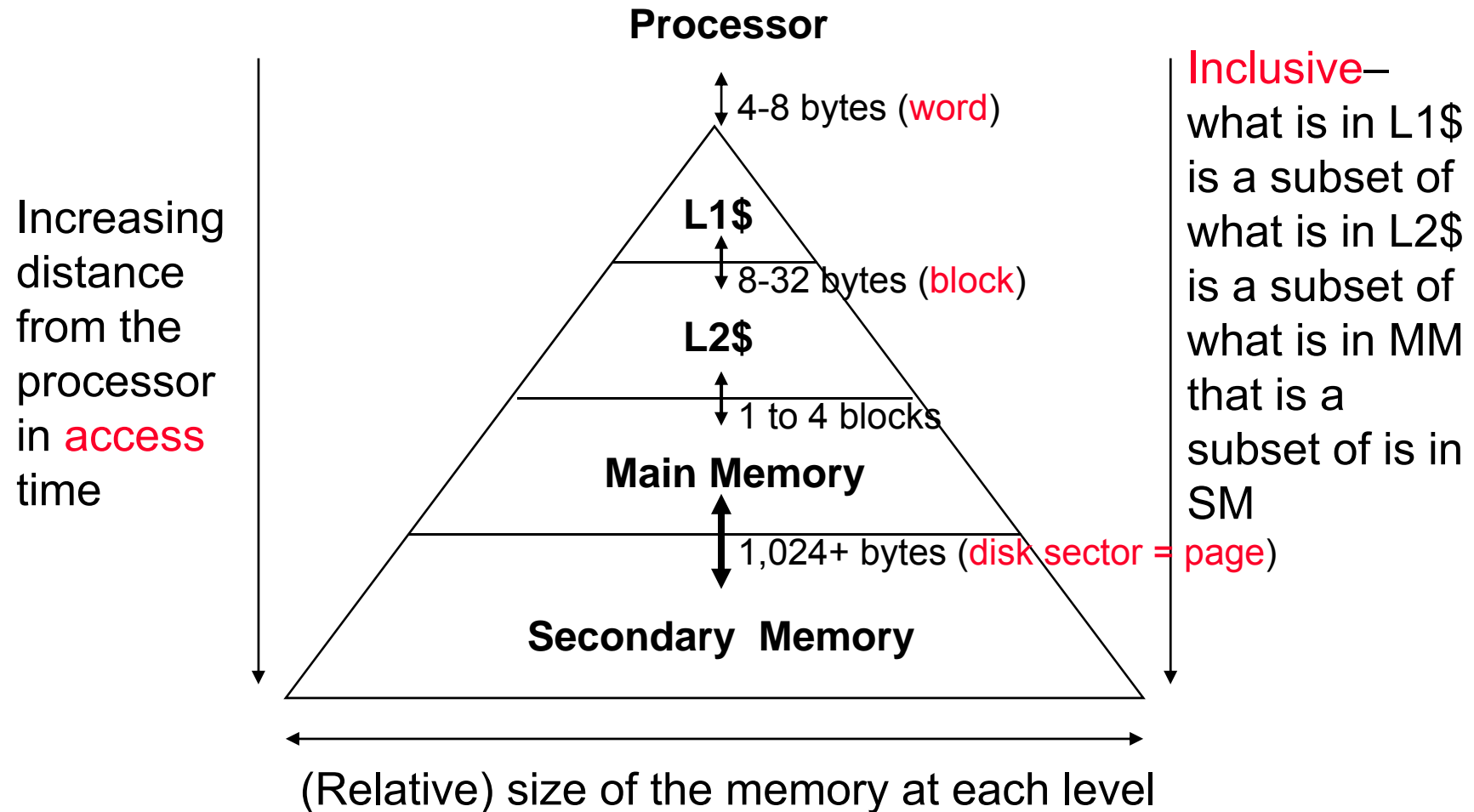⇒ Move blocks consisting of contiguous words closer to the processor

# The Memory Hierarchy:  Terminology

❑ Block (or line): the minimum unit of information that is present (or not) in a cache

❑ Hit Rate: the fraction of memory accesses found in a level of the memory hierarchy

● Hit Time: Time to access that level which consists of

Time to access the block + Time to determine hit/miss

❑ Miss Rate: the fraction of memory accesses *not* found in a level of the memory hierarchy    $\Rightarrow$   1 - (Hit Rate)

● Miss Penalty: Time to replace a block in that level with the corresponding block from a lower level which consists of

Time to access the block in the lower level + Time to transmit that block to the level that experienced the miss + Time to insert the block in that level + Time to pass the block to the requestor

## Hit Time << Miss Penalty

# Characteristics of the Memory Hierarchy

**Processor**

Increasing distance from the processor in access time

↕ 4-8 bytes (word)

**L1$**

↕ 8-32 bytes (block)

**L2$**

↕ 1 to 4 blocks

**Main Memory**

↕ 1,024+ bytes (disk sector = page)

**Secondary Memory**

Inclusive– what is in L1$ is a subset of what is in L2$ is a subset of what is in MM that is a subset of is in SM

(Relative) size of the memory at each level

# How is the Hierarchy Managed?

- ❑ registers ↔ memory
  - by compiler (programmer?)

- ❑ cache ↔ main memory
  - by the cache controller hardware

- ❑ main memory ↔ disks
  - by the operating system (virtual memory)
  - virtual to physical address mapping assisted by the hardware (TLB)
  - by the programmer (files)

# Cache Basics

❑ Two questions to answer (in hardware):

- Q1: How do we know if a data item is in the cache?
- Q2: If it is, how do we find it?

❑ Direct mapped

- Each memory block is mapped to exactly one block in the cache
  - lots of lower level blocks must share blocks in the cache

- Address mapping (to answer Q2):

  (block address) modulo (# of blocks in the cache)

- Have a tag associated with each cache block that contains the address information (the upper portion of the address) required to identify the block (to answer Q1)

# Caching:  A Simple First Example

**Main Memory**

**Cache**

Index Valid  Tag     Data

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

0000xx
0001xx
0010xx
0011xx
0100xx
0101xx
0110xx
0111xx
1000xx
1001xx
1010xx
1011xx
1100xx
1101xx
1110xx
1111xx

One word blocks
Two low order bits
define the byte in the
word (32b words)

Q2: How do we find it?

Use next 2 low order
memory address bits
– the index – to
determine which
cache block (i.e.,
modulo the number of
blocks in the cache)

Q1: Is it there?

Compare the cache
tag to the high order 2
memory address bits to
tell if the memory block
is in the cache

(block address) modulo (# of blocks in the cache)

# Caching: A Simple First Example

**Main Memory**

**Cache**

Index Valid  Tag     Data

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

0000xx
0001xx
0010xx
0011xx
0100xx
0101xx
0110xx
0111xx
1000xx
1001xx
1010xx
1011xx
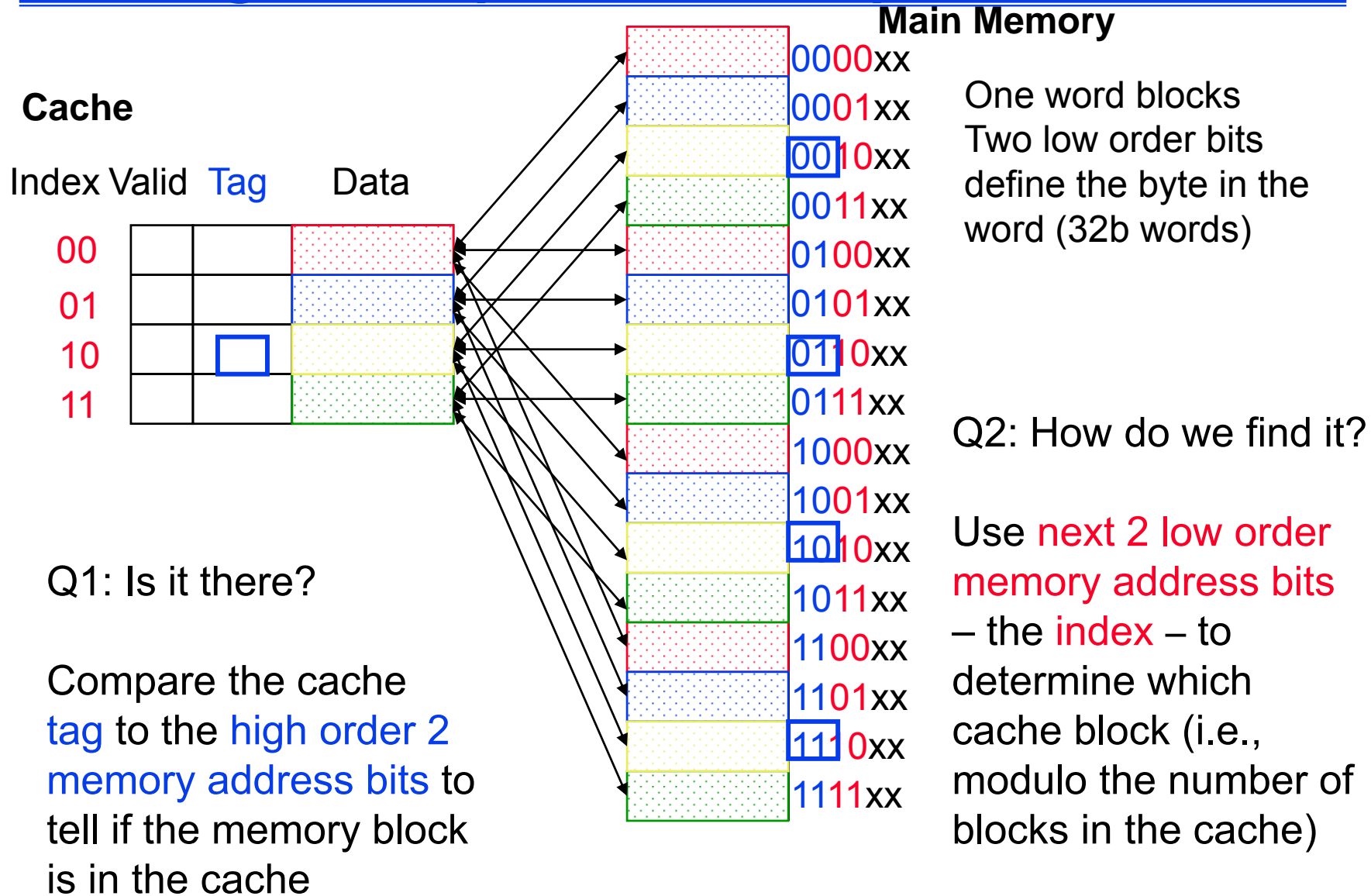1100xx
1101xx
1110xx
1111xx

One word blocks
Two low order bits
define the byte in the
word (32b words)

Q2: How do we find it?

Use next 2 low order
memory address bits
– the index – to
determine which
cache block (i.e.,
modulo the number of
blocks in the cache)

Q1: Is it there?

Compare the cache
tag to the high order 2
memory address bits to
tell if the memory block
is in the cache

(block address) modulo (# of blocks in the cache)

# Direct Mapped Cache

❑ Consider the main memory word reference string

Start with an empty cache - all
blocks initially marked as not valid

0   1   2   3   4   3   4   15

| 0 | | 1 | | 2 | | 3 | |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

| 4 | | 3 | | 4 | | 15 | |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

# Direct Mapped Cache

❑ Consider the main memory word reference string

Start with an empty cache - all blocks initially marked as not valid

0  1  2  3  4  3  4  15

**0** miss

| 00 | Mem(0) |
|----|--------|
|    |        |
|    |        |
|    |        |

**1** miss

| 00 | Mem(0) |
|----|--------|
| 00 | Mem(1) |
|    |        |
|    |        |

**2** miss

| 00 | Mem(0) |
|----|--------|
| 00 | Mem(1) |
| 00 | Mem(2) |
|    |        |

**3** miss

| 00 | Mem(0) |
|----|--------|
| 00 | Mem(1) |
| 00 | Mem(2) |
| 00 | Mem(3) |

**4** miss

01      4

| 00 | Mem(0) |
|----|--------|
| 00 | Mem(1) |
| 00 | Mem(2) |
| 00 | Mem(3) |

**3** hit

| 01 | Mem(4) |
|----|--------|
| 00 | Mem(1) |
| 00 | Mem(2) |
| 00 | Mem(3) |

**4** hit

| 01 | Mem(4) |
|----|--------|
| 00 | Mem(1) |
| 00 | Mem(2) |
| 00 | Mem(3) |

**15** miss

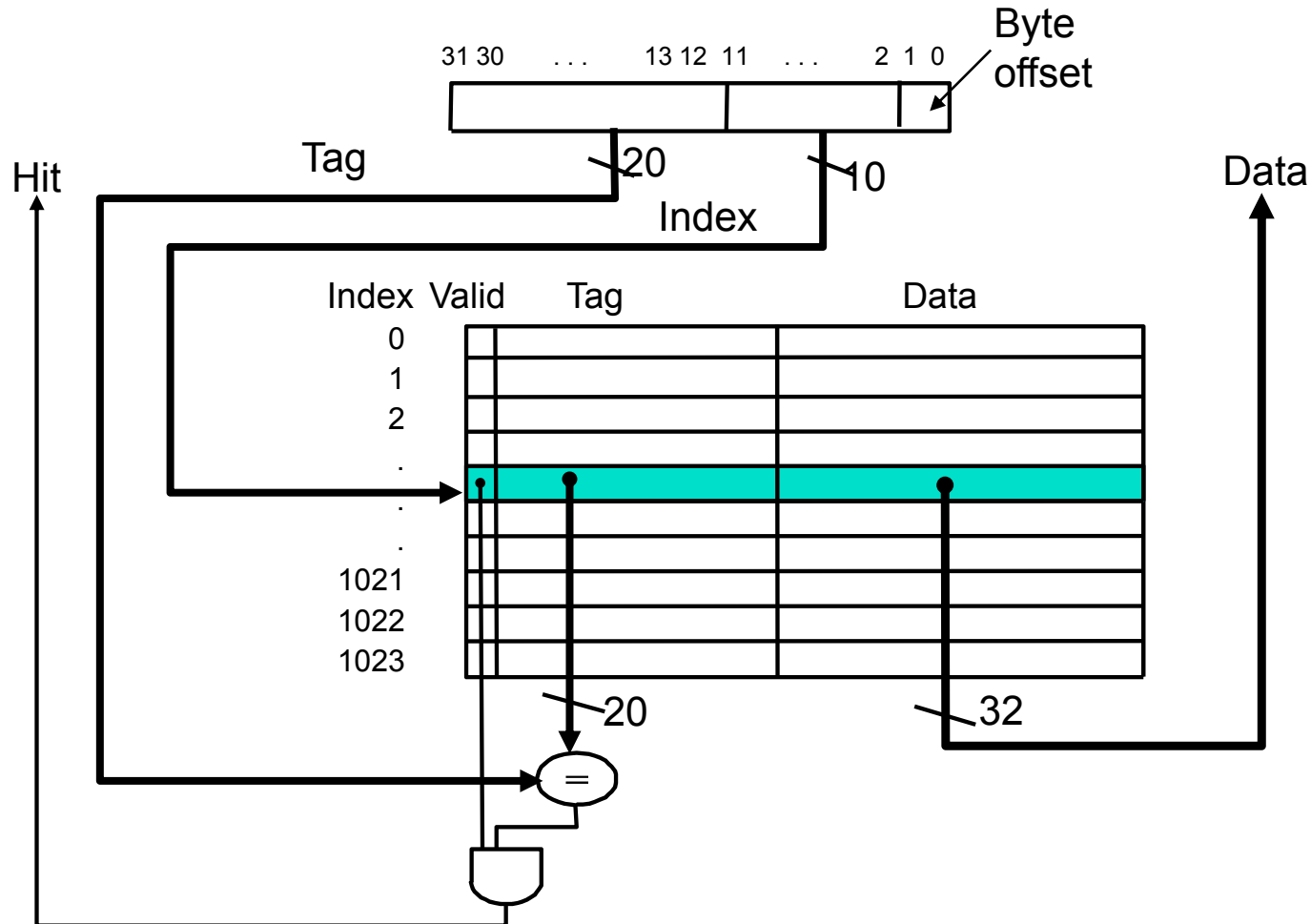| 01 | Mem(4) |
|----|--------|
| 00 | Mem(1) |
| 00 | Mem(2) |
| 00 | Mem(3) |

11      15

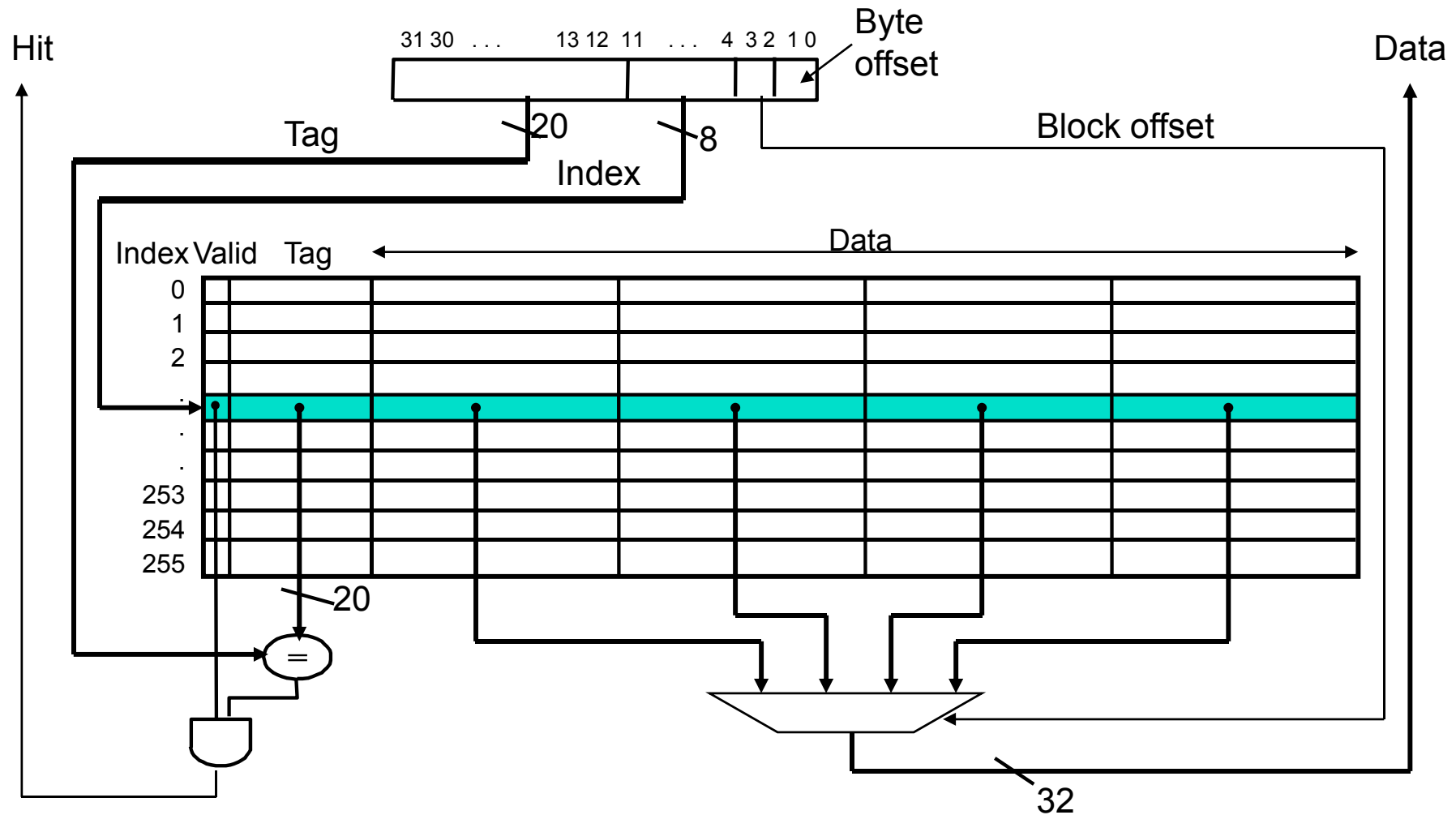● 8 requests, 6 misses

# MIPS Direct Mapped Cache Example

❑ One word blocks, cache size = 1K words (or 4KB)



*What kind of locality are we taking advantage of?*

# Multiword Block Direct Mapped Cache

❑ Four words/block, cache size = 1K words



*What kind of locality are we taking advantage of?*

# Taking Advantage of Spatial Locality

❑ Let cache block hold more than one word

Start with an empty cache - all
blocks initially marked as not valid

0  1  2  3  4  3  4  15

**0**

|  |  |  |
|--|--|--|
|  |  |  |

**1**

|  |  |  |
|--|--|--|
|  |  |  |

**2**

|  |  |  |
|--|--|--|
|  |  |  |

**3**

|  |  |  |
|--|--|--|
|  |  |  |

**4**

|  |  |  |
|--|--|--|
|  |  |  |

**3**

|  |  |  |
|--|--|--|
|  |  |  |

**4**

|  |  |  |
|--|--|--|
|  |  |  |

**15**

|  |  |  |
|--|--|--|
|  |  |  |

# Taking Advantage of Spatial Locality

❑ Let cache block hold more than one word

Start with an empty cache - all
blocks initially marked as not valid

0  1  2  3  4  3  4  15

**0** miss

| 00 | Mem(1) | Mem(0) |
|----|--------|--------|
|    |        |        |

**1** hit

| 00 | Mem(1) | Mem(0) |
|----|--------|--------|
|    |        |        |

**2** miss

| 00 | Mem(1) | Mem(0) |
|----|--------|--------|
| 00 | Mem(3) | Mem(2) |

**3** hit

| 00 | Mem(1) | Mem(0) |
|----|--------|--------|
| 00 | Mem(3) | Mem(2) |

**4** miss

01                    5          4
| 00 | Mem(1) | Mem(0) |
|----|--------|--------|
| 00 | Mem(3) | Mem(2) |

**3** hit

| 01 | Mem(5) | Mem(4) |
|----|--------|--------|
| 00 | Mem(3) | Mem(2) |

**4** hit

| 01 | Mem(5) | Mem(4) |
|----|--------|--------|
| 00 | Mem(3) | Mem(2) |

**15** miss

11                    15         14
| 01 | Mem(5) | Mem(4) |
|----|--------|--------|
| 00 | Mem(3) | Mem(2) |

● 8 requests, 4 misses

# Miss Rate vs Block Size vs Cache Size



**Block size (bytes)**

Legend:
- 8 KB
- 16 KB
- 64 KB
- 256 KB

❑ Miss rate goes up if the block size becomes a significant fraction of the cache size because the number of blocks that can be held in the same size cache is smaller (increasing capacity misses)

# Cache Field Sizes

❑ The number of bits in a cache includes both the storage for data and for the tags

  ● 32-bit byte address

  ● For a direct mapped cache with $2^n$ blocks, $n$ bits are used for the index

  ● For a block size of $2^m$ words ($2^{m+2}$ bytes), $m$ bits are used to address the word within the block and 2 bits are used to address the byte within the word

❑ What is the size of the tag field?

❑ The total number of bits in a direct-mapped cache is then

$2^n$ x (block size + tag field size + valid field size)

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

2^10 x (4x32 +18 + 1) = 2^10 x 147 = 147Kbits

# Handling Cache Hits

❑ Read hits (I$ and D$)

- this is what we want!

❑ Write hits (D$ only)

- require the cache and memory to be consistent
  - always write the data into both the cache block and the next level in the memory hierarchy (write-through)
  - writes run at the speed of the next level in the memory hierarchy – so slow! – or can use a write buffer and stall only if the write buffer is full
- allow cache and memory to be inconsistent
  - write the data only into the cache block (write-back the cache block to the next level in the memory hierarchy when that cache block is "evicted")
  - need a dirty bit for each data cache block to tell if it needs to be written back to memory when it is evicted – can use a write buffer to help "buffer" write-backs of dirty blocks

# Sources of Cache Misses

❑ Compulsory (cold start or process migration, first reference):

- First access to a block, "cold" fact of life, not a whole lot you can do about it.  If you are going to run "millions" of instruction, compulsory misses are insignificant

- Solution: increase block size (increases miss penalty; very large blocks could increase miss rate)

❑ Capacity:

- Cache cannot contain all blocks accessed by the program

- Solution: increase cache size (may increase access time)

❑ Conflict (collision):

- Multiple memory locations mapped to the same cache location

- Solution 1: increase cache size

- Solution 2: increase associativity (stay tuned) (may increase access time)

# Handling Cache Misses (Single Word Blocks)

❑ Read misses (I$ and D$)

- stall the pipeline, fetch the block from the next level in the memory hierarchy, install it in the cache and send the requested word to the processor, then let the pipeline resume

❑ Write misses (D$ only)

1. stall the pipeline, fetch the block from next level in the memory hierarchy, install it in the cache (which may involve having to evict a dirty block if using a write-back cache), write the word from the processor to the cache, then let the pipeline resume

or

2. Write allocate – just write the word into the cache updating both the tag and data, no need to check for cache hit, no need to stall

or

3. No-write allocate – skip the cache write (but must invalidate that cache block since it will now hold stale data) and just write the word to the write buffer (and eventually to the next memory level), no need to stall if the write buffer isn't full

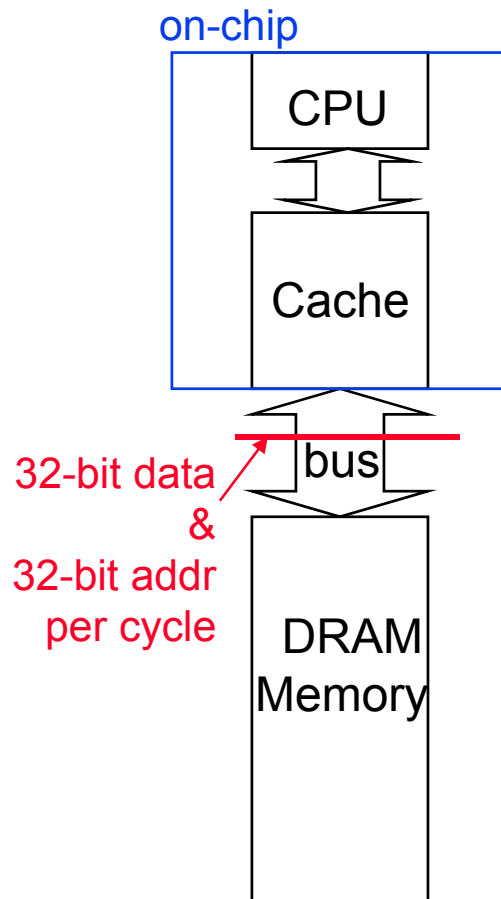# Multiword Block Considerations

❑ Read misses (I$ and D$)

- Processed the same as for single word blocks – a miss returns the entire block from memory

- Miss penalty grows as block size grows
    - Early restart – processor resumes execution as soon as the requested word of the block is returned
    - Requested word first – requested word is transferred from the memory to the cache (and processor) first

- Nonblocking cache – allows the processor to continue to access the cache while the cache is handling an earlier miss

❑ Write misses (D$)

- If using write allocate must *first* fetch the block from memory and then write the word to the block (or could end up with a "garbled" block in the cache (e.g., for 4 word blocks, a new tag, one word of data from the new block, and three words of data from the old block)

# Memory Systems that Support Caches

❑ The off-chip interconnect and memory architecture can affect overall system performance in dramatic ways

on-chip

```
┌─────────────────┐
│   ┌─────────┐   │
│   │   CPU   │   │
│   └─────────┘   │
│       ⬍         │
│   ┌─────────┐   │
│   │  Cache  │   │
│   └─────────┘   │
└───────⬍─────────┘
      ═══⬍═══
   32-bit data  ┌bus┐
        &       └───┘
   32-bit addr    ⬍
   per cycle  ┌─────────┐
             │  DRAM   │
             │ Memory  │
             └─────────┘
```

One word wide organization (one word wide bus and one word wide memory)

❑ Assume
1. 1 memory bus clock cycle to send the addr
2. 15 memory bus clock cycles to get the 1st word in the block from DRAM (row cycle time), 5 memory bus clock cycles for 2nd, 3rd, 4th words (column access time)
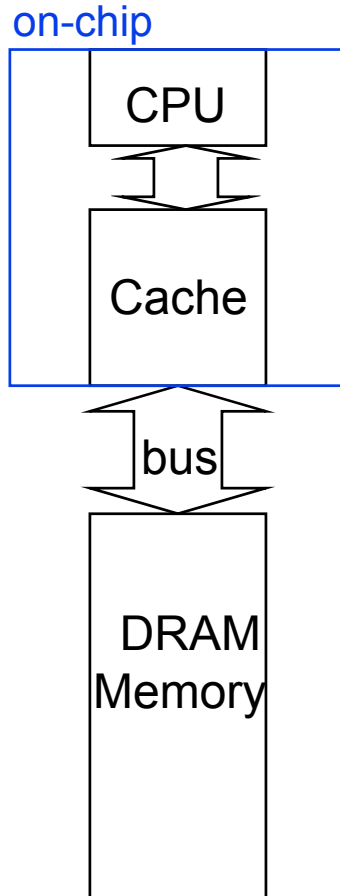3. 1 memory bus clock cycle to return a word of data

❑ Memory-Bus to Cache bandwidth
● number of bytes accessed from memory and transferred to cache/CPU per memory bus clock cycle

# Review: (DDR) SDRAM Operation

□ After a row is read into the SRAM register

- Input CAS as the starting "burst" address along with a burst length

- Transfers a burst of data (ideally a cache block) from a series of sequential addr's within that row

  - The memory bus clock controls transfer of successive words in the burst

**Column Address** → □+1

N cols

DRAM

N rows

Row Address

**N x M SRAM**

M bit planes

**M-bit Output**

**Cycle Time**

**1st M-bit Access**   **2nd M-bit**   **3rd M-bit**   **4th M-bit**

RAS

CAS

Row Address    Col Address                     Row Add

# One Word Wide Bus, One Word Blocks

on-chip

CPU

Cache

bus

DRAM Memory

❑ If the block size is one word, then for a memory access due to a cache miss, the pipeline will have to stall for the number of cycles required to return one data word from memory

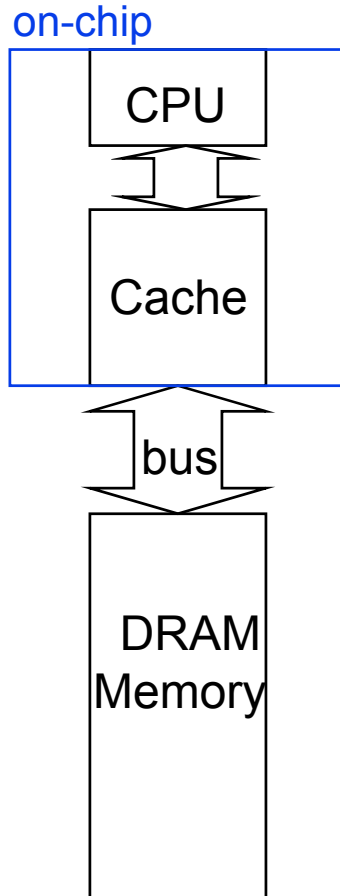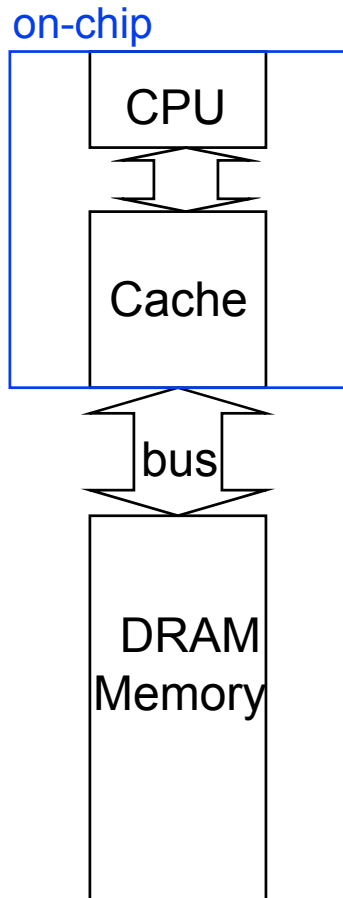cycle to send address

cycles to read DRAM

cycle to return data

——— total clock cycles miss penalty

❑ Number of bytes transferred per clock cycle (bandwidth) for a single miss is

bytes per memory bus clock cycle

# One Word Wide Bus, One Word Blocks

on-chip

CPU

Cache

bus

DRAM Memory

❑ If the block size is one word, then for a memory access due to a cache miss, the pipeline will have to stall for the number of cycles required to return one data word from memory

| | |
|---|---|
| 1 | memory bus clock cycle to send address |
| 15 | memory bus clock cycles to read DRAM |
| 1 | memory bus clock cycle to return data |
| 17 | total clock cycles miss penalty |

❑ Number of bytes transferred per clock cycle (bandwidth) for a single miss is

4/17 = 0.235  bytes per memory bus clock cycle

# One Word Wide Bus, Four Word Blocks

on-chip

CPU

Cache

bus

DRAM
Memory

❑ What if the block size is four words and each word is in a different DRAM row?

cycle to send 1st address
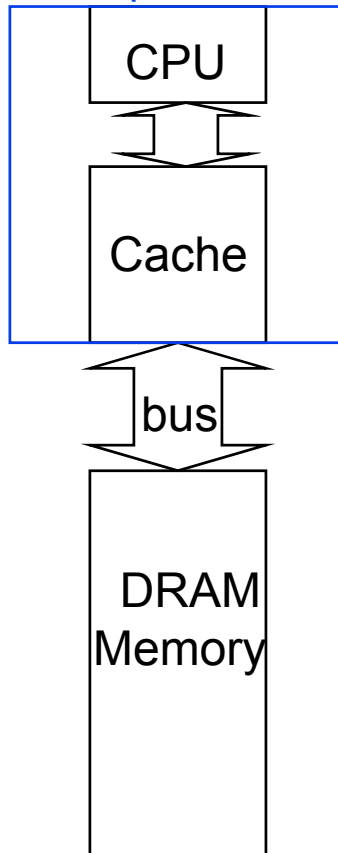
cycles to read DRAM

cycles to return last data word
___
total clock cycles miss penalty

❑ Number of bytes transferred per clock cycle (bandwidth) for a single miss is
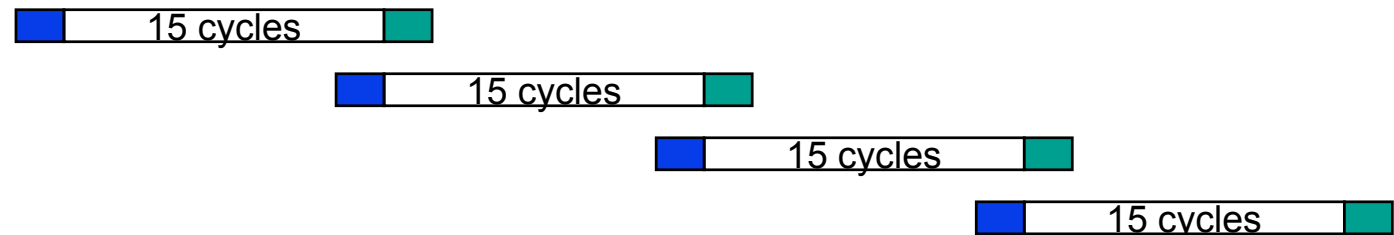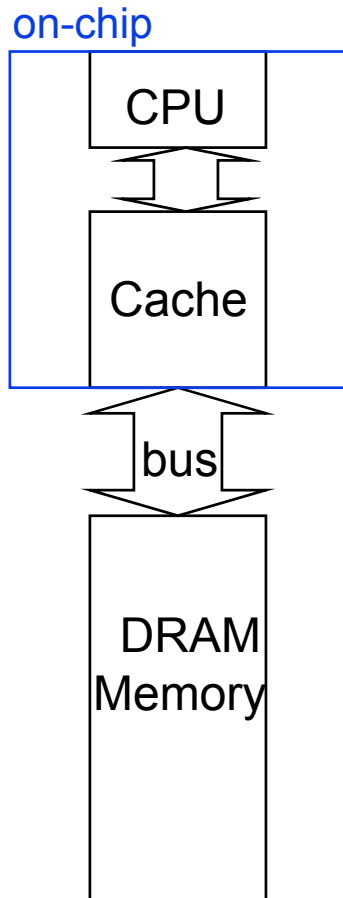
bytes per clock

# One Word Wide Bus, Four Word Blocks

❑ What if the block size is four words and each word is in a different DRAM row?

|   |   |
|---|---|
| 1 | cycle to send 1st address |
| 4 x 15 = 60 | cycles to read DRAM |
| 1 | cycles to return last data word |
| 62 | total clock cycles miss penalty |

on-chip

CPU

Cache

bus

DRAM
Memory

15 cycles

15 cycles

15 cycles

15 cycles

❑ Number of bytes transferred per clock cycle (bandwidth) for a single miss is

(4 x 4)/62 = 0.258  bytes per clock

# One Word Wide Bus, Four Word Blocks

on-chip



❑ What if the block size is four words and all words are in the same DRAM row?

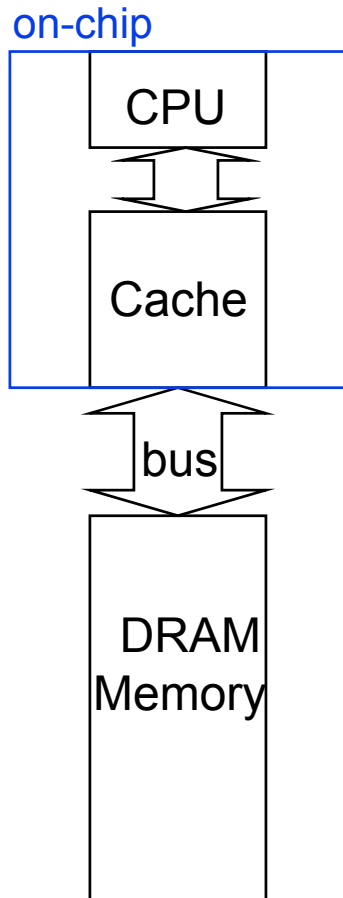cycle to send 1st address

cycles to read DRAM

cycles to return last data word
___
total clock cycles miss penalty

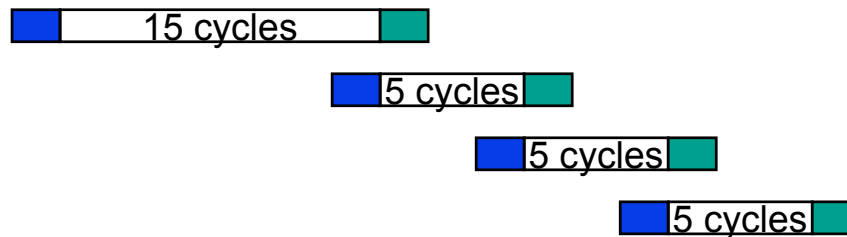❑ Number of bytes transferred per clock cycle (bandwidth) for a single miss is

bytes per clock

# One Word Wide Bus, Four Word Blocks

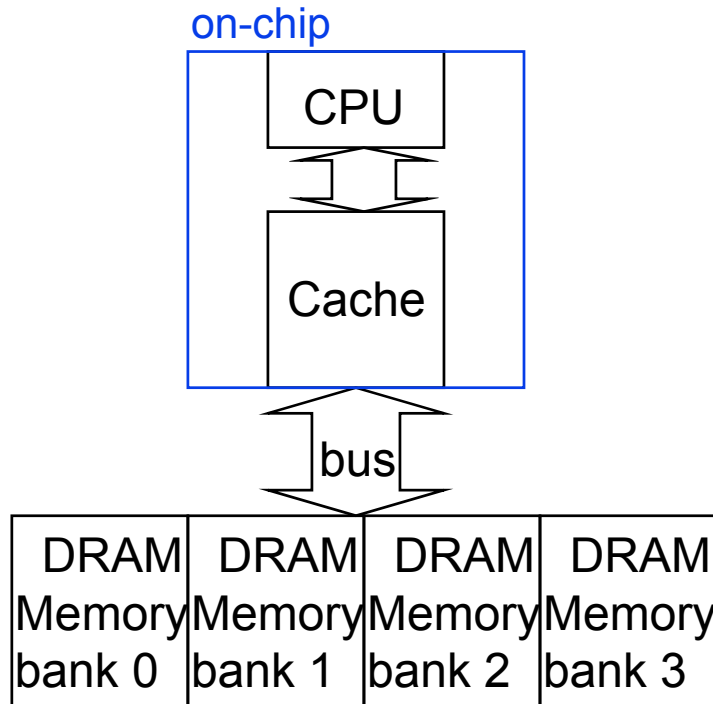❏ What if the block size is four words and all words are in the same DRAM row?

| | |
|---|---|
| 1 | cycle to send 1st address |
| 15 + 3*5 = 30 | cycles to read DRAM |
| 1 | cycles to return last data word |
| 32 | total clock cycles miss penalty |



❏ Number of bytes transferred per clock cycle (bandwidth) for a single miss is

(4 x 4)/32 = 0.5    bytes per clock

# Interleaved Memory, One Word Wide Bus

❑ For a block size of four words

on-chip

```
        CPU

       Cache

        bus
```

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

        cycle to send 1st address

        cycles to read DRAM banks

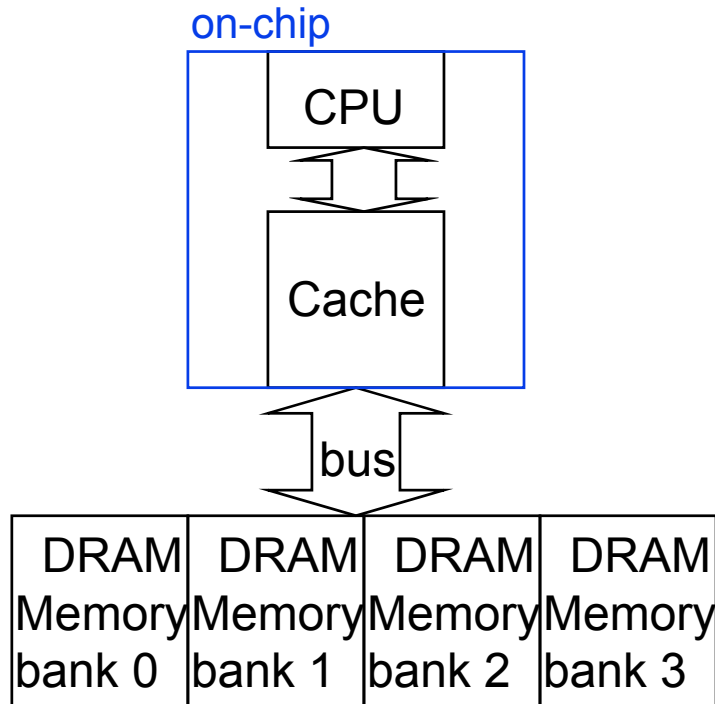    ____  cycles to return last data word

        total clock cycles miss penalty

❑ Number of bytes transferred per clock cycle (bandwidth) for a single miss is

        bytes per clock

# Interleaved Memory, One Word Wide Bus

❏ For a block size of four words

| | |
|---|---|
| 1 | cycle to send 1st address |
| 15 | cycles to read DRAM banks |
| 4*1 = 4 | cycles to return last data word |
| 20 | total clock cycles miss penalty |

on-chip

CPU

Cache

bus

| DRAM Memory bank 0 | DRAM Memory bank 1 | DRAM Memory bank 2 | DRAM Memory bank 3 |
|---|---|---|---|

15 cycles

15 cycles

15 cycles

15 cycles

❏ Number of bytes transferred per clock cycle (bandwidth) for a single miss is

(4 x 4)/20 = 0.8    bytes per clock

# DRAM Memory System Summary

❑ Its important to match the cache characteristics

  ● caches access one block at a time (usually more than one word)

❑ with the DRAM characteristics

  ● use DRAMs that support fast multiple word accesses, preferably ones that match the block size of the cache

❑ with the memory-bus characteristics

  ● make sure the memory-bus can support the DRAM access rates and patterns

  ● with the goal of increasing the Memory-Bus to Cache bandwidth

# Reducing Cache Miss Rates #1

1. Allow more flexible block placement

- In a direct mapped cache a memory block maps to exactly one cache block

- At the other extreme, could allow a memory block to be mapped to *any* cache block – fully associative cache

- A compromise is to divide the cache into sets each of which consists of n "ways" (n-way set associative).  A memory block maps to a unique set (specified by the index field) and can be placed in any way of that set (so there are n choices)

  (block address) modulo (# sets in the cache)

# Another Reference String Mapping

❑ Consider the main memory word reference string

Start with an empty cache - all
blocks initially marked as not valid

0   4   0   4   0   4   0   4

**0**

| | |
|---|---|
| | |
| | |
| | |

**4**

| | |
|---|---|
| | |
| | |
| | |

**0**

| | |
|---|---|
| | |
| | |
| | |

**4**

| | |
|---|---|
| | |
| | |
| | |

**0**

| | |
|---|---|
| | |
| | |
| | |

**4**

| | |
|---|---|
| | |
| | |
| | |

**0**

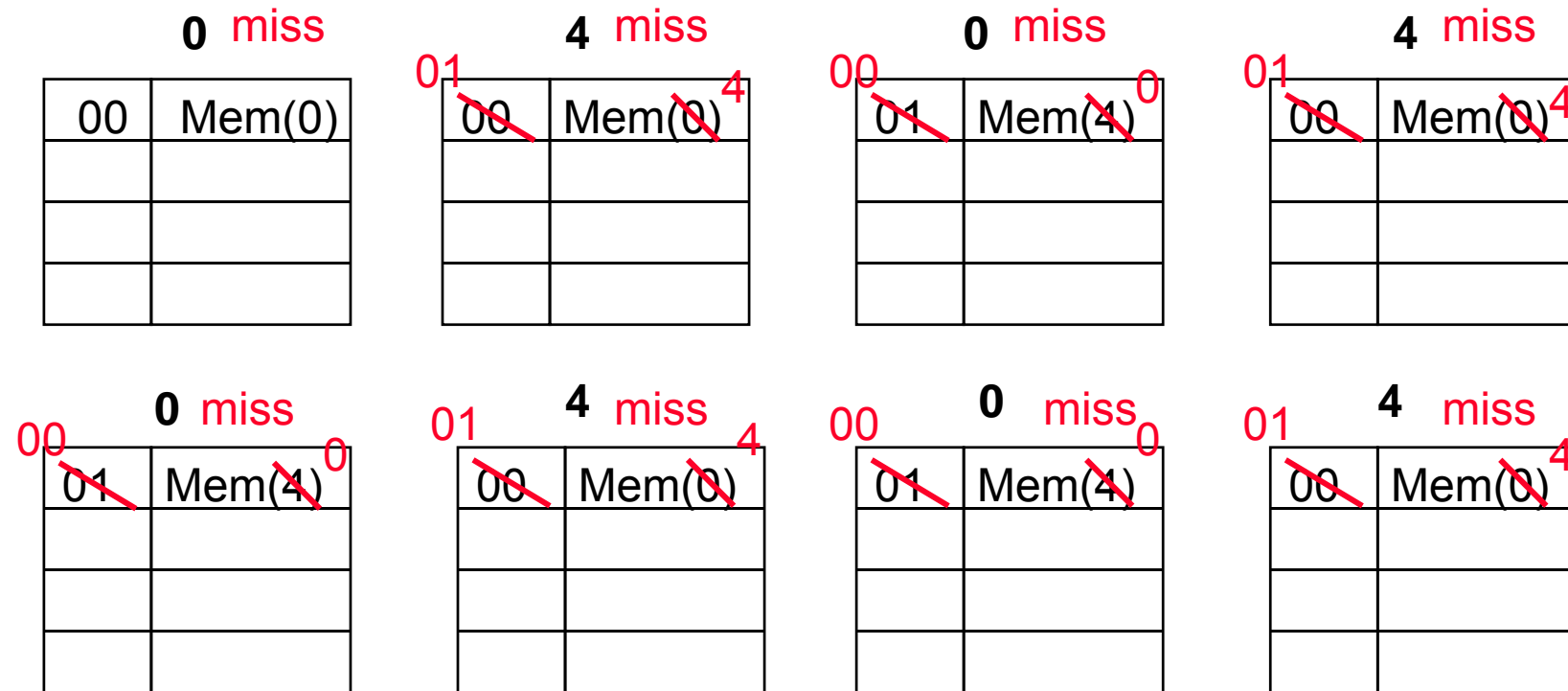| | |
|---|---|
| | |
| | |
| | |

**4**

| | |
|---|---|
| | |
| | |
| | |

# Another Reference String Mapping

❑ Consider the main memory word reference string

Start with an empty cache - all
blocks initially marked as not valid

0   4   0   4   0   4   0   4



● 8 requests, 8 misses

❑ Ping pong effect due to conflict misses - two memory
locations that map into the same cache block

# Set Associative Cache Example

**Main Memory**

**Cache**

| Way | Set | V | Tag | Data |
|-----|-----|---|-----|------|
| 0 | 0 | | | |
| | 1 | | | |
| 1 | 0 | | | |
| | 1 | | | |

0000xx
00001xx
0010xx
00111xx
0100xx
01001xx
0110xx
01111xx
1000xx
10001xx
1010xx
10111xx
1100xx
11001xx
1110xx
11111xx

One word blocks
Two low order bits define the byte in the word (32b words)

Q2: How do we find it?

Use next 1 low order memory address bit to determine which cache set (i.e., modulo the number of sets in the cache)

Q1: Is it there?

Compare *all* the cache tags in the set to the high order 3 memory address bits to tell if the memory block is in the cache

# Another Reference String Mapping

❑ Consider the main memory word reference string

Start with an empty cache - all
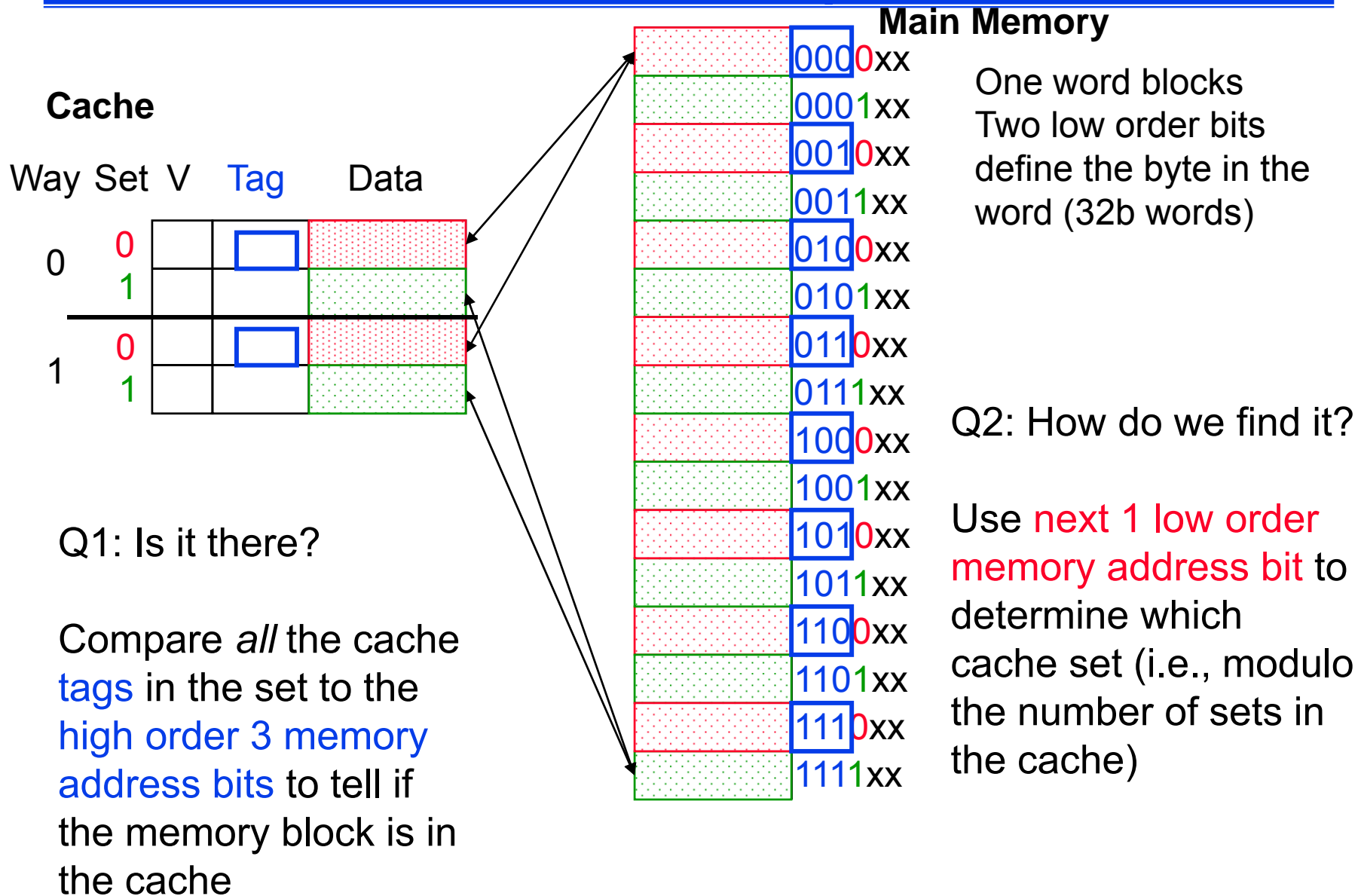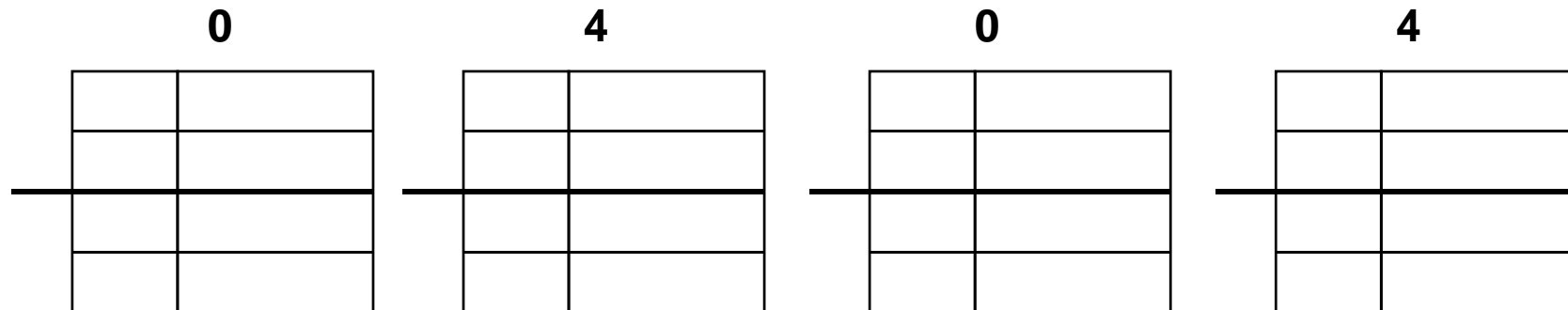blocks initially marked as not valid

0   4   0   4   0   4   0   4

| 0 | | 4 | | 0 | | 4 | |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

# Another Reference String Mapping

❑ Consider the main memory word reference string

Start with an empty cache - all
blocks initially marked as not valid
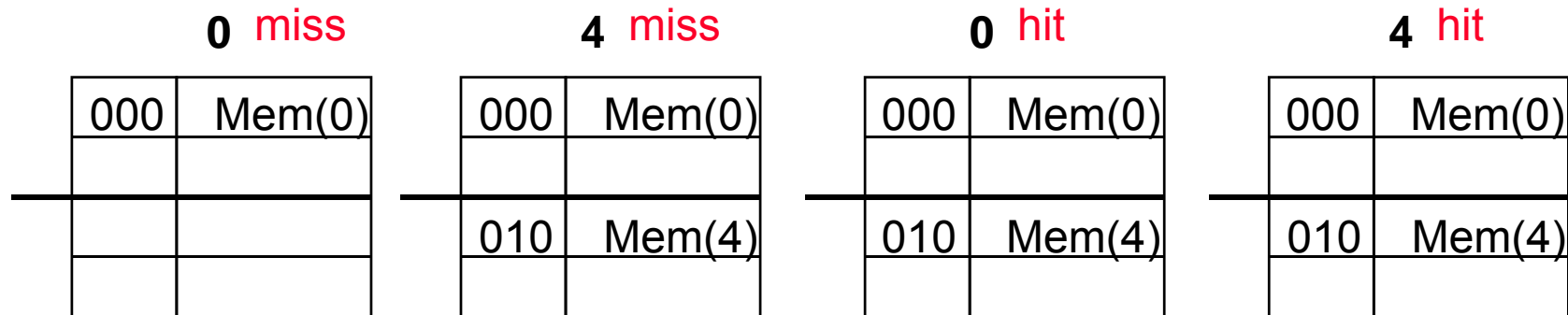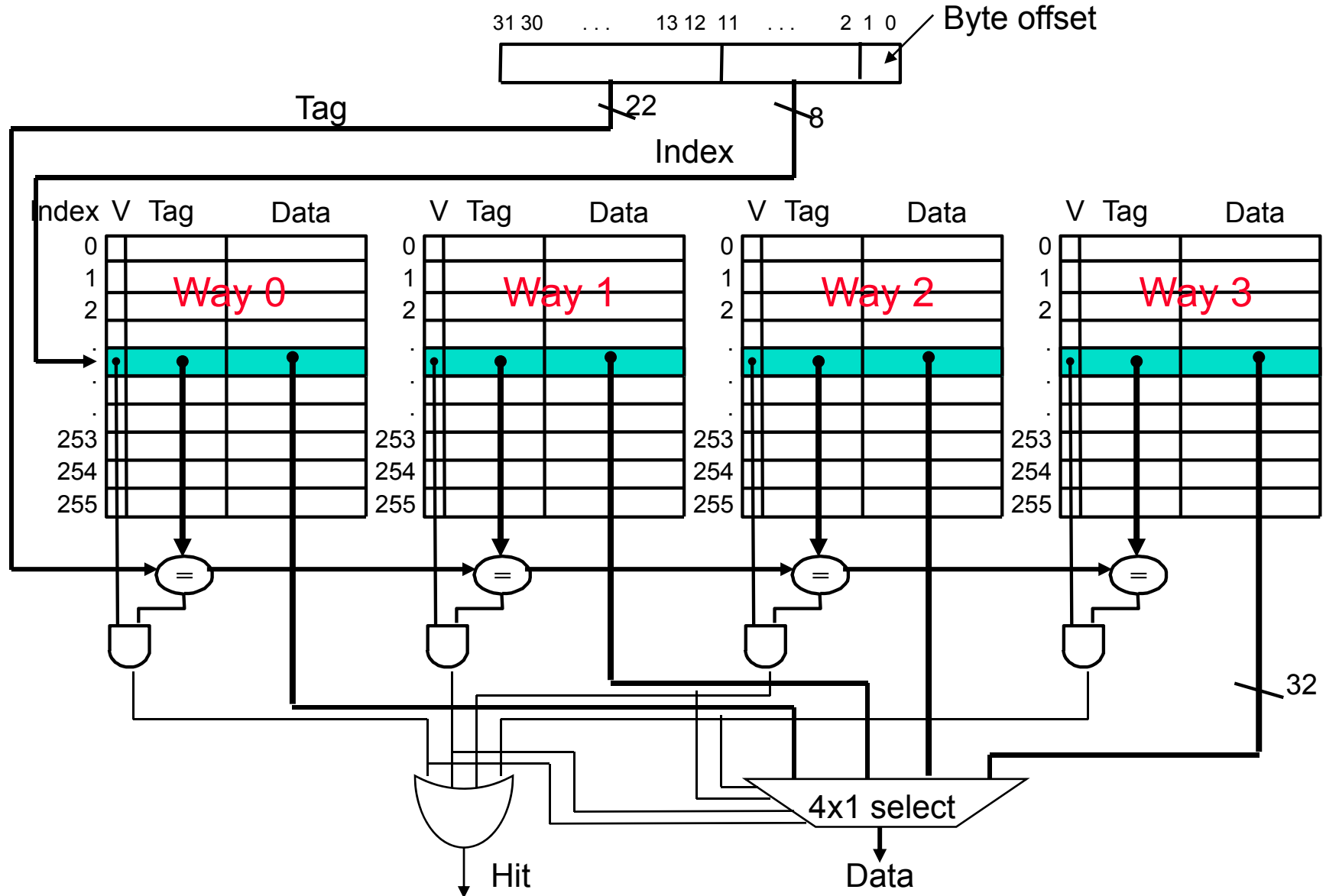
0   4   0   4   0   4   0   4

**0** miss

| 000 | Mem(0) |
|-----|--------|
|     |        |
|     |        |
|     |        |

**4** miss

| 000 | Mem(0) |
|-----|--------|
|     |        |
| 010 | Mem(4) |
|     |        |

**0** hit

| 000 | Mem(0) |
|-----|--------|
|     |        |
| 010 | Mem(4) |
|     |        |

**4** hit

| 000 | Mem(0) |
|-----|--------|
|     |        |
| 010 | Mem(4) |
|     |        |

● 8 requests, 2 misses

❑ Solves the ping pong effect in a direct mapped cache due to conflict misses since now two memory locations that map into the same cache set can co-exist!

# Four-Way Set Associative Cache

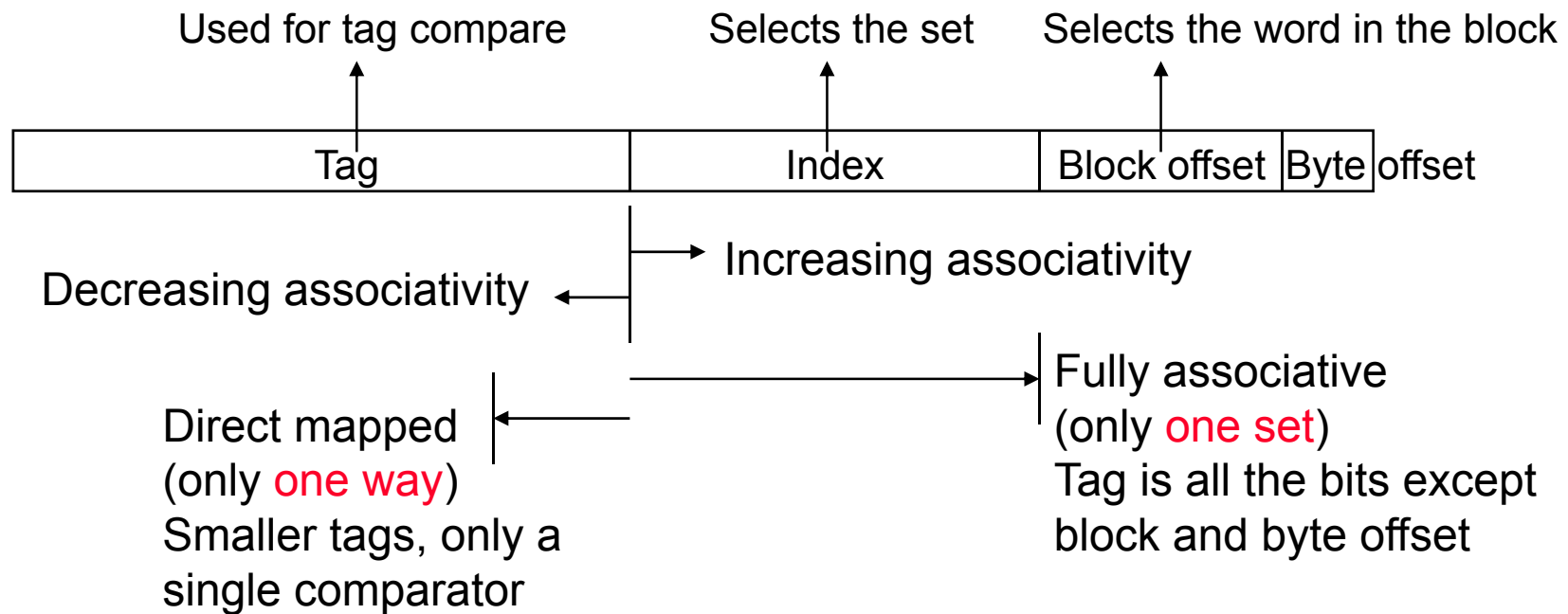❑ $2^8 = 256$ sets each with four ways (each with one block)

# Range of Set Associative Caches

❑ For a fixed size cache, each increase by a factor of two in associativity doubles the number of blocks per set (i.e., the number or ways) and halves the number of sets – decreases the size of the index by 1 bit and increases the size of the tag by 1 bit

| Tag | Index | Block offset | Byte | offset |
|-----|-------|--------------|------|--------|

# Range of Set Associative Caches

❑ For a fixed size cache, each increase by a factor of two in associativity doubles the number of blocks per set (i.e., the number or ways) and halves the number of sets – decreases the size of the index by 1 bit and increases the size of the tag by 1 bit
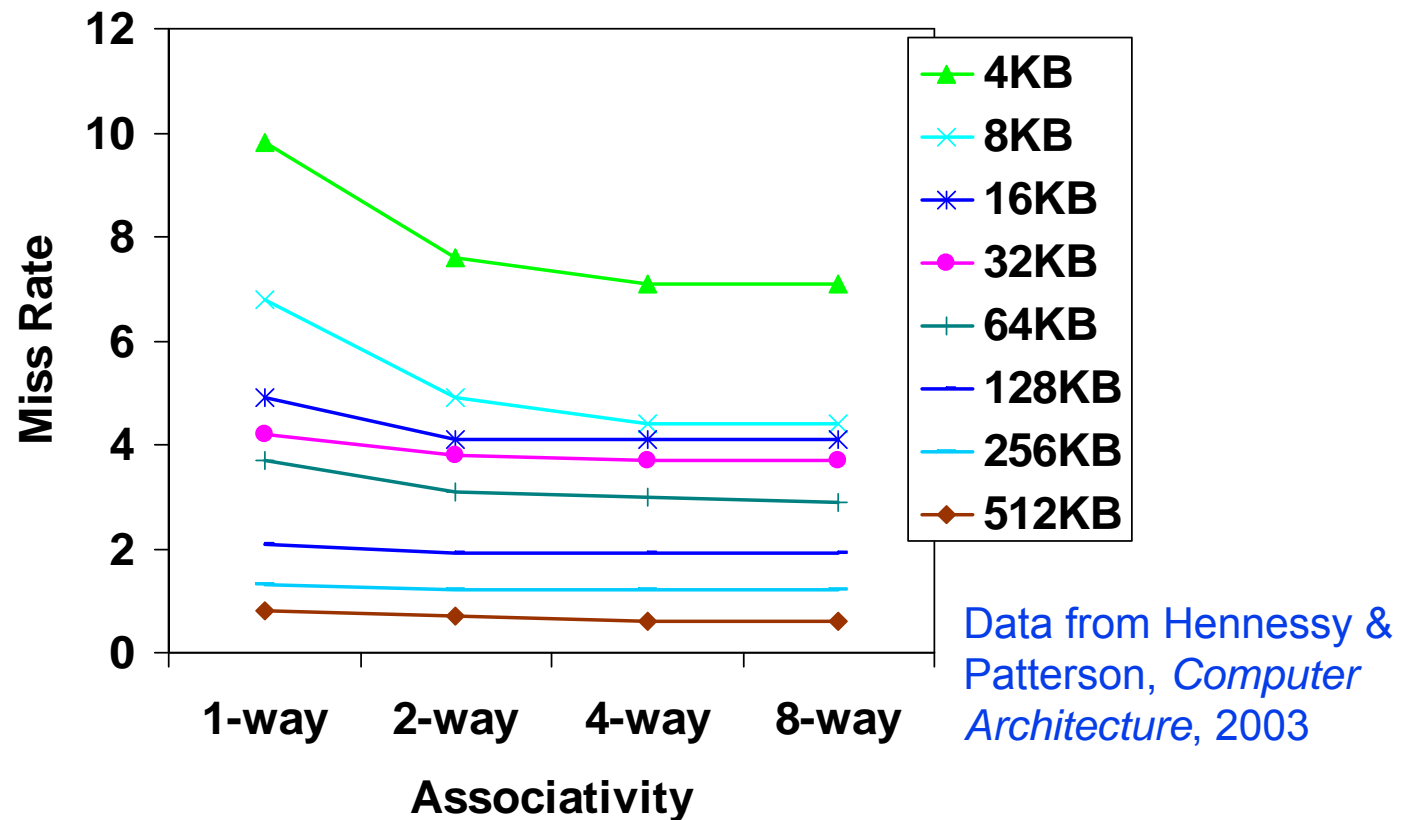
Used for tag compare      Selects the set      Selects the word in the block

| Tag | Index | Block offset | Byte offset |

Increasing associativity →

← Decreasing associativity

Direct mapped
(only one way)
Smaller tags, only a
single comparator

Fully associative
(only one set)
Tag is all the bits except
block and byte offset

# Costs of Set Associative Caches

- ❑ When a miss occurs, which way's block do we pick for replacement?
  - ● Least Recently Used (LRU): the block replaced is the one that has been unused for the longest time
    - Must have hardware to keep track of when each way's block was used relative to the other blocks in the set
    - For 2-way set associative, takes one bit per set → set the bit when a block is referenced (and reset the other way's bit)
- ❑ N-way set associative cache costs
  - ● N comparators (delay and area)
  - ● MUX delay (set selection) before data is available
  - ● Data available after set selection (and Hit/Miss decision). In a direct mapped cache, the cache block is available before the Hit/Miss decision
    - So its not possible to just assume a hit and continue and recover later if it was a miss

# Benefits of Set Associative Caches

❑ The choice of direct mapped or set associative depends on the cost of a miss versus the cost of implementation



Data from Hennessy & Patterson, *Computer Architecture*, 2003

❑ Largest gains are in going from direct mapped to 2-way (20%+ reduction in miss rate)

# Reducing Cache Miss Rates #2

2. Use multiple levels of caches

❑ With advancing technology have more than enough room on the die for bigger L1 caches *or* for a second level of caches – normally a unified L2 cache (i.e., it holds both instructions and data) and in some cases even a unified L3 cache

❑ For our example, $CPI_{ideal}$ of 2, 100 cycle miss penalty (to main memory) and a 25 cycle miss penalty (to UL2$), 36% load/stores, a 2% (4%) L1 I$ (D$) miss rate, add a 0.5% UL2$ miss rate

$$CPI_{stalls} = 2 + .02×25 + .36×.04×25 + .005×100 +$$
$$.36×.005×100 = 3.54$$

(as compared to 5.44 with no L2$)

# Multilevel Cache Design Considerations

❑ Design considerations for L1 and L2 caches are very different

- Primary cache should focus on minimizing hit time in support of a shorter clock cycle
  - Smaller with smaller block sizes

- Secondary cache(s) should focus on reducing miss rate to reduce the penalty of long main memory access times
  - Larger with larger block sizes
  - Higher levels of associativity

❑ The miss penalty of the L1 cache is significantly reduced by the presence of an L2 cache – so it can be smaller (i.e., faster) but have a higher miss rate

❑ For the L2 cache, hit time is less important than miss rate

- The L2$ hit time determines L1$'s miss penalty
- L2$ local miss rate >> than the global miss rate

# Two Machines' Cache Parameters

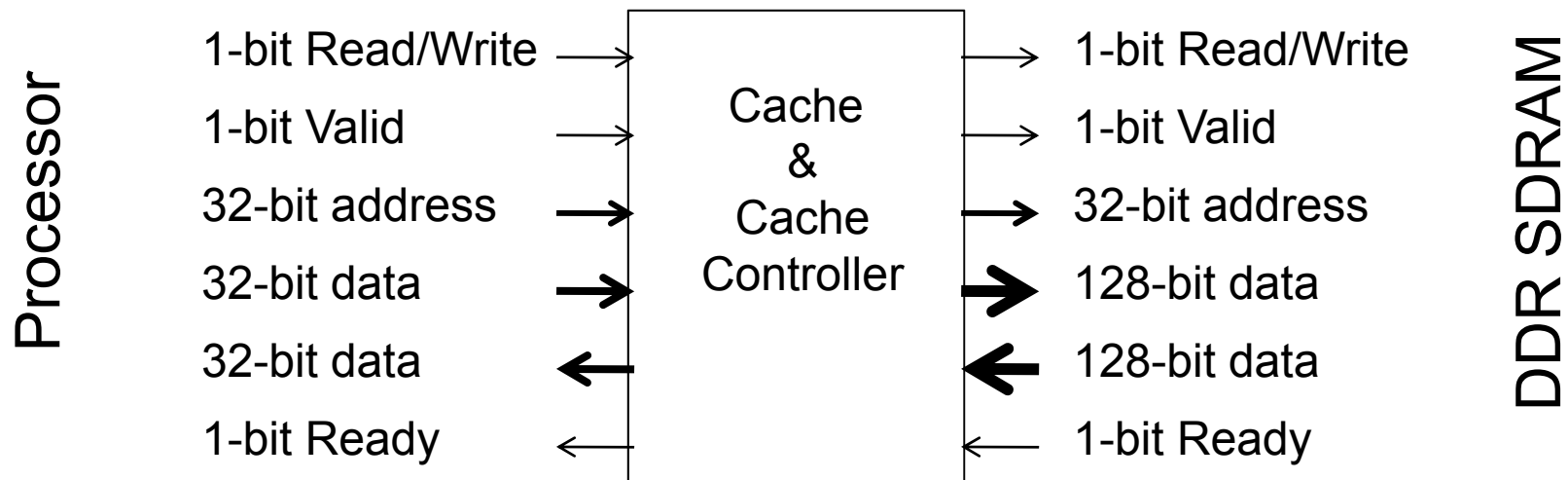| | Intel Nehalem | AMD Barcelona |
|---|---|---|
| L1 cache organization & size | Split I$ and D$; 32KB for each per core; 64B blocks | Split I$ and D$; 64KB for each per core; 64B blocks |
| L1 associativity | 4-way (I), 8-way (D) set assoc.; ~LRU replacement | 2-way set assoc.; LRU replacement |
| L1 write policy | write-back, write-allocate | write-back, write-allocate |
| L2 cache organization & size | Unified; 256MB (0.25MB) per core; 64B blocks | Unified; 512KB (0.5MB) per core; 64B blocks |
| L2 associativity | 8-way set assoc.; ~LRU | 16-way set assoc.; ~LRU |
| L2 write policy | write-back | write-back |
| L2 write policy | write-back, write-allocate | write-back, write-allocate |
| L3 cache organization & size | Unified; 8192KB (8MB) shared by cores; 64B blocks | Unified; 2048KB (2MB) shared by cores; 64B blocks |
| L3 associativity | 16-way set assoc. | 32-way set assoc.; evict block shared by fewest cores |
| L3 write policy | write-back, write-allocate | write-back; write-allocate |

# Two Machines' Cache Parameters

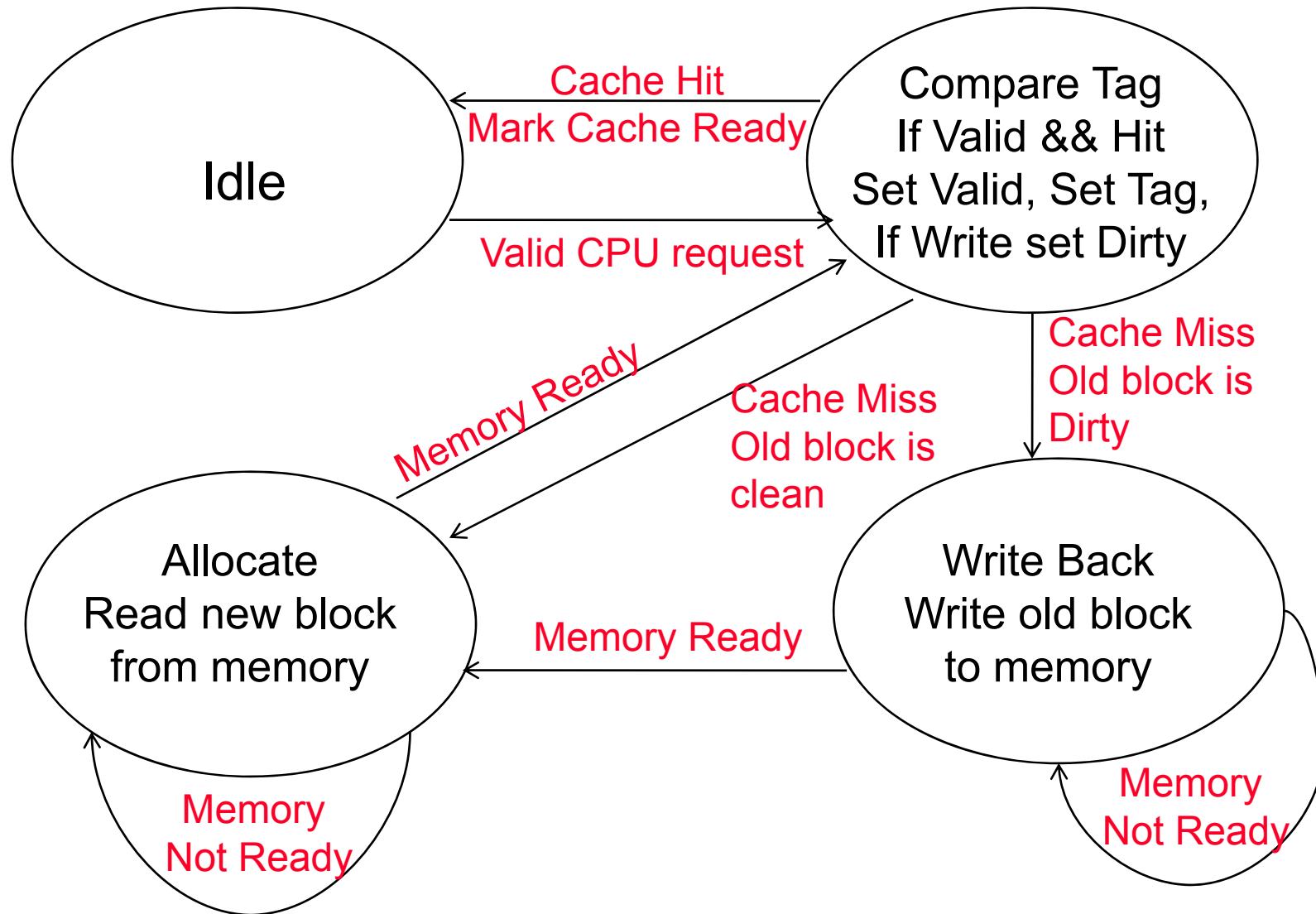| | Intel P4 | AMD Opteron |
|---|---|---|
| L1 organization | Split I$ and D$ | Split I$ and D$ |
| L1 cache size | 8KB for D$, 96KB for trace cache (~I$) | 64KB for each of I$ and D$ |
| L1 block size | 64 bytes | 64 bytes |
| L1 associativity | 4-way set assoc. | 2-way set assoc. |
| L1 replacement | ~ LRU | LRU |
| L1 write policy | write-through | write-back |
| L2 organization | Unified | Unified |
| L2 cache size | 512KB | 1024KB (1MB) |
| L2 block size | 128 bytes | 64 bytes |
| L2 associativity | 8-way set assoc. | 16-way set assoc. |
| L2 replacement | ~LRU | ~LRU |
| L2 write policy | write-back | write-back |

# FSM Cache Controller

❑ Key characteristics for a simple L1 cache

- Direct mapped

- Write-back using write-allocate

- Block size of 4 32-bit words (so 16B); Cache size of 16KB (so 1024 blocks)

- 18-bit tags, 10-bit index, 2-bit block offset, 2-bit byte offset, dirty bit, valid bit, LRU bits (if set associative)
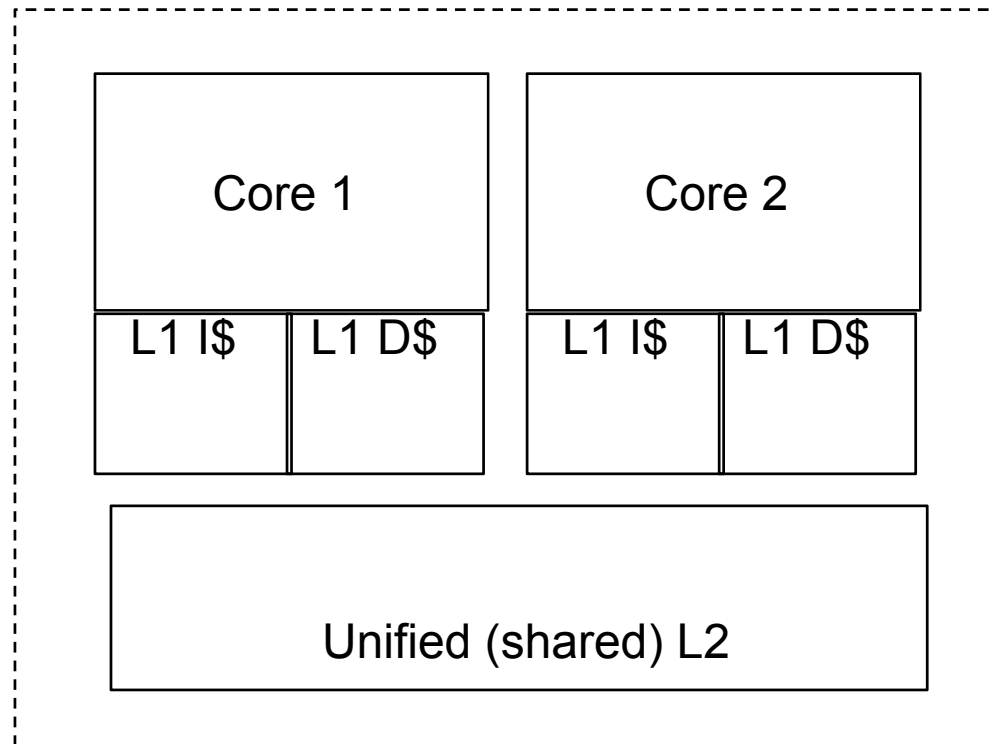
| Processor | | Cache & Cache Controller | | DDR SDRAM |
|---|---|---|---|---|
| 1-bit Read/Write | → | | → | 1-bit Read/Write |
| 1-bit Valid | → | | → | 1-bit Valid |
| 32-bit address | → | | → | 32-bit address |
| 32-bit data | → | | → | 128-bit data |
| 32-bit data | ← | | ← | 128-bit data |
| 1-bit Ready | ← | | ← | 1-bit Ready |

# Four State Cache Controller

# Cache Coherence in Multicores

❑ In multicore processors the cores *share* a common physical address space, causing a cache coherence problem

```
┌─────────────────────────────────────────────────┐
│  ┌──────────────────┐   ┌──────────────────┐     │
│  │                  │   │                  │     │
│  │     Core 1       │   │     Core 2       │     │
│  │                  │   │                  │     │
│  ├────────┬─────────┤   ├────────┬─────────┤     │
│  │        │         │   │        │         │     │
│  │ L1 I$  │ L1 D$   │   │ L1 I$  │ L1 D$   │     │
│  │        │         │   │        │         │     │
│  └────────┴─────────┘   └────────┴─────────┘     │
│  ┌───────────────────────────────────────┐       │
│  │                                        │       │
│  │          Unified (shared) L2           │       │
│  │                                        │       │
│  └───────────────────────────────────────┘       │
└─────────────────────────────────────────────────┘
```
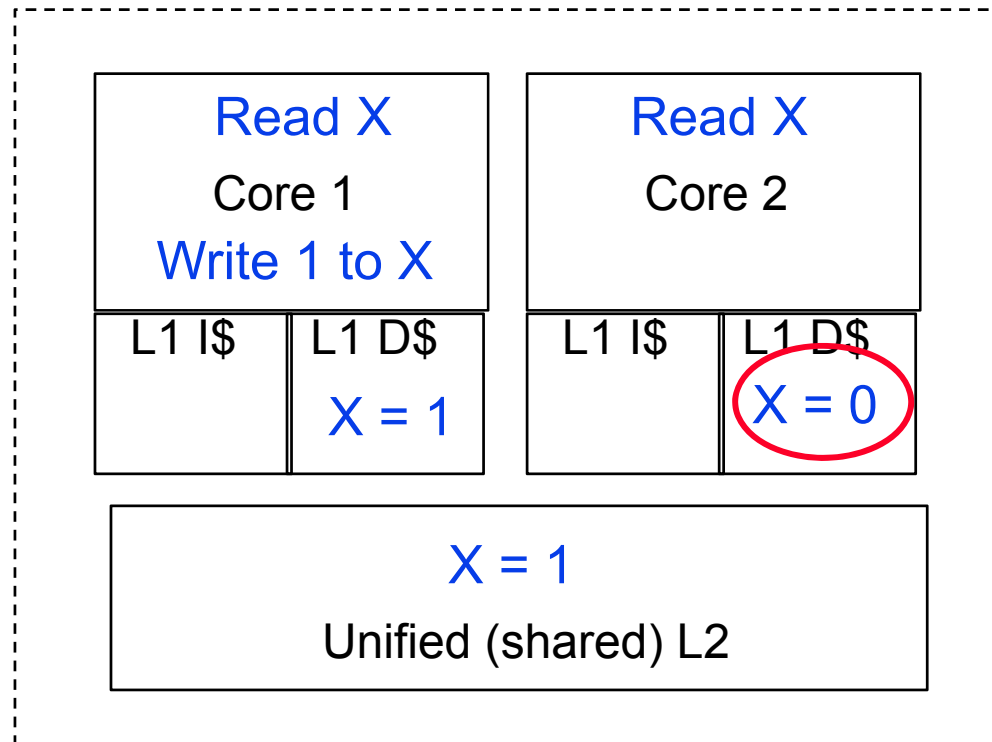
# Cache Coherence in Multicores

❑ In multicore processors the cores *share* a common physical address space, causing a cache coherence problem

# A Coherent Memory System

❑ Any read of a data item should return the most recently written value of the data item

- Coherence – defines what values can be returned by a read
  - Writes to the same location are serialized (two writes to the same location must be seen in the same order by all cores)
- Consistency – determines when a written value will be returned by a read

❑ To enforce coherence, caches must provide

- Replication of shared data items in multiple cores' caches
  - Replication reduces both latency and contention for a read shared data item
- Migration of shared data items to a core's local cache
  - Migration reduced the latency of the access the data and the bandwidth demand on the shared memory (L2 in our example)

# Cache Coherence Protocols

❑ Need a hardware protocol to ensure cache coherence the most popular of which is snooping

- The cache controllers monitor (snoop) on the broadcast medium (e.g., bus) with duplicate address tag hardware (so they don't interfere with core's access to the cache) to determine if their cache has a copy of a block that is requested

❑ Write invalidate protocol – writes require exclusive access and invalidate *all* other copies

- Exclusive access ensure that no other readable or writable copies of an item exists

❑ If two processors attempt to write the same data at the same time, one of them wins the race causing the other core's copy to be invalidated. For the other core to complete, it must obtain a new copy of the data which must now contain the updated value – thus enforcing write serialization

# Handling Writes

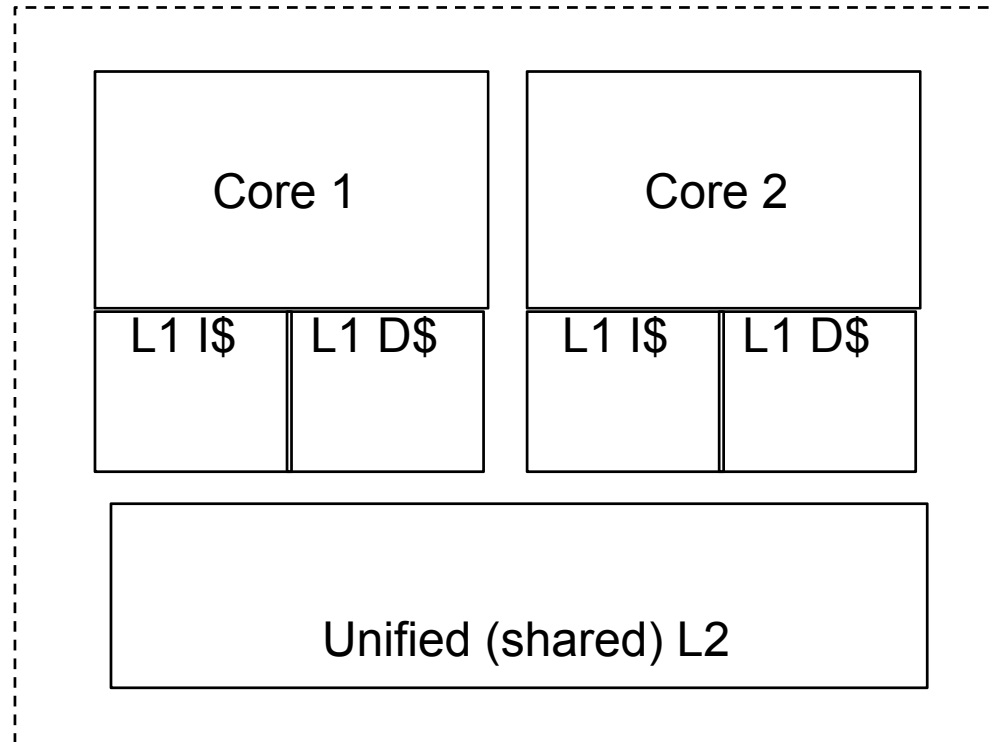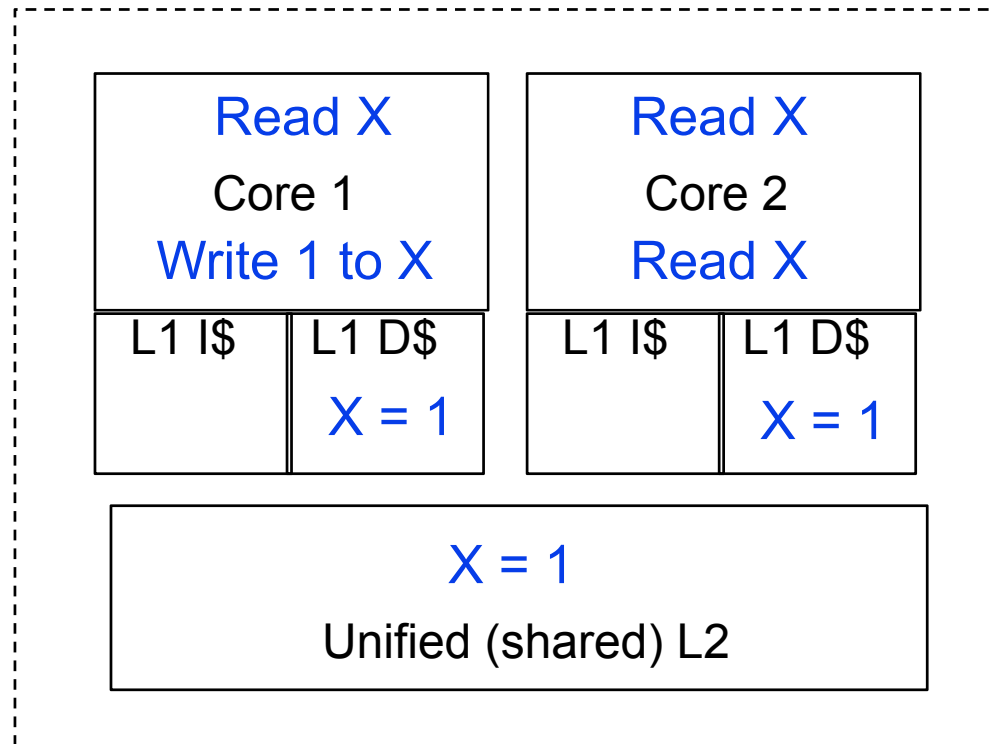Ensuring that all other processors sharing data are informed of writes can be handled two ways:

1.  Write-update (write-broadcast) – writing processor broadcasts new data over the bus, all copies are updated

    - All writes go to the bus $\rightarrow$ higher bus traffic

    - Since new values appear in caches sooner, can reduce latency

2.  Write-invalidate – writing processor issues invalidation signal on bus, cache snoops check to see if they have a copy of the data, if so they invalidate their cache block containing the word (this allows multiple readers but only one writer)

    - Uses the bus only on the first write $\rightarrow$ lower bus traffic, so better use of bus bandwidth
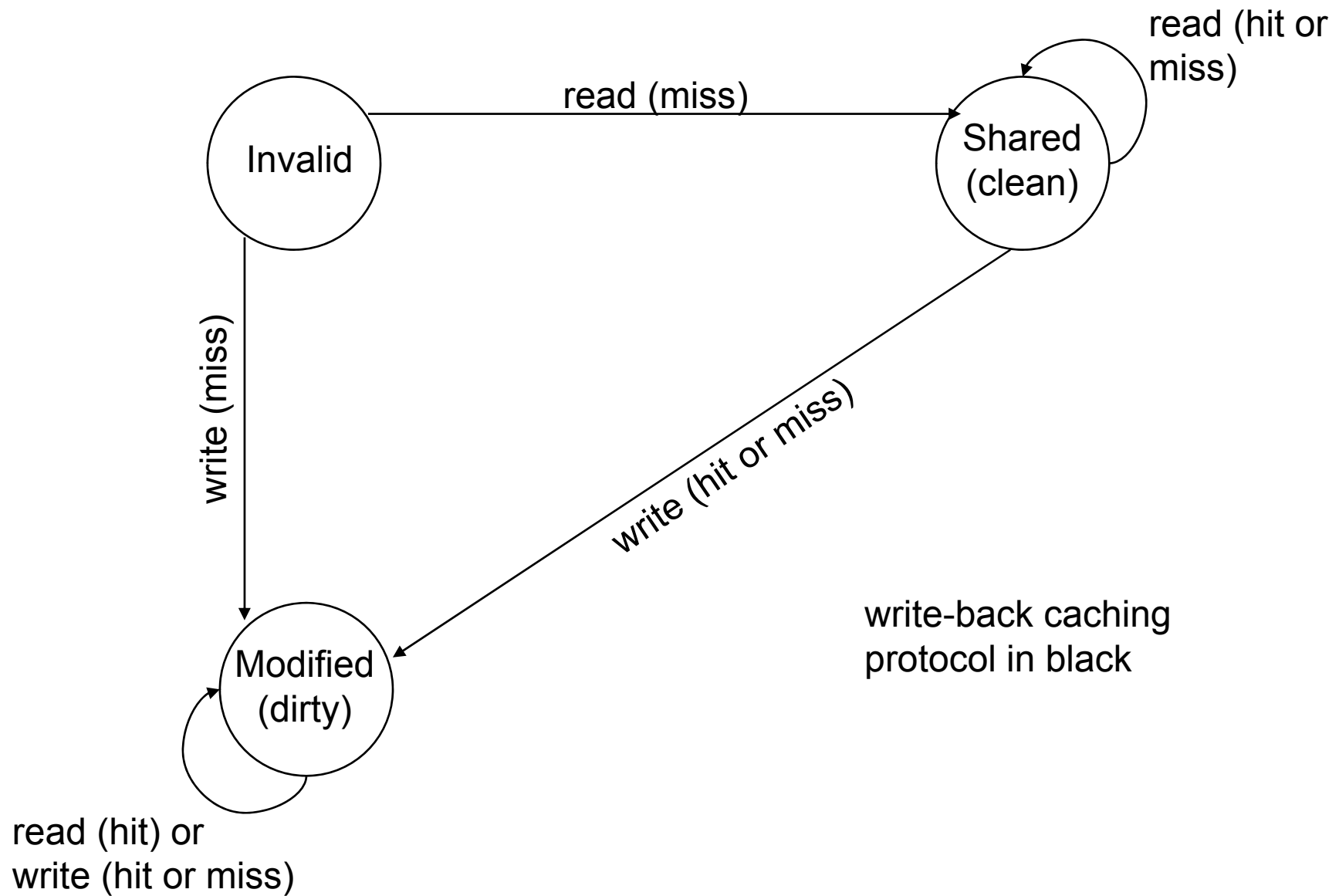
# Example of Snooping Invalidation

# Example of Snooping Invalidation



❑ When the second miss by Core 2 occurs, Core 1 responds with the value canceling the response from the L2 cache (and also updating the L2 copy)
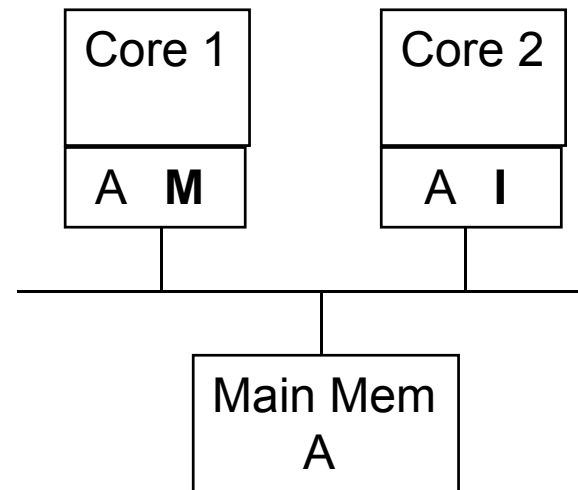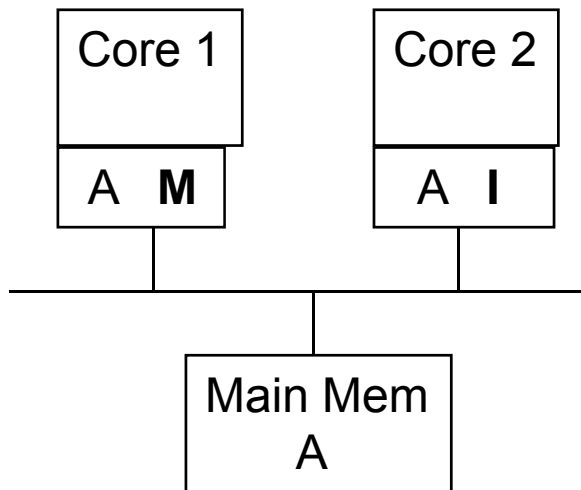
# A Write-Invalidate CC Protocol



read (hit or miss)

Invalid

Shared (clean)

read (miss)

write (miss)

write (hit or miss)

Modified (dirty)
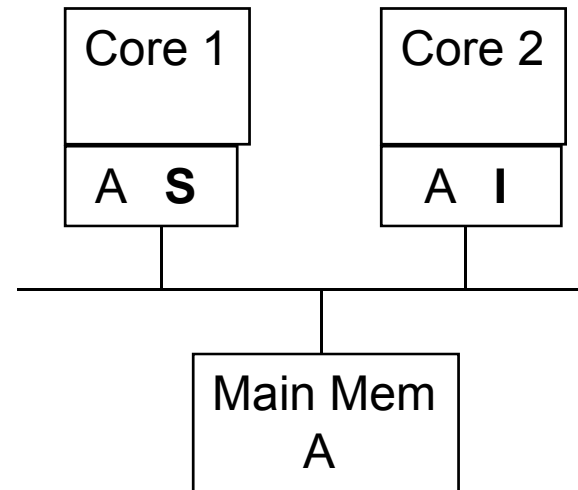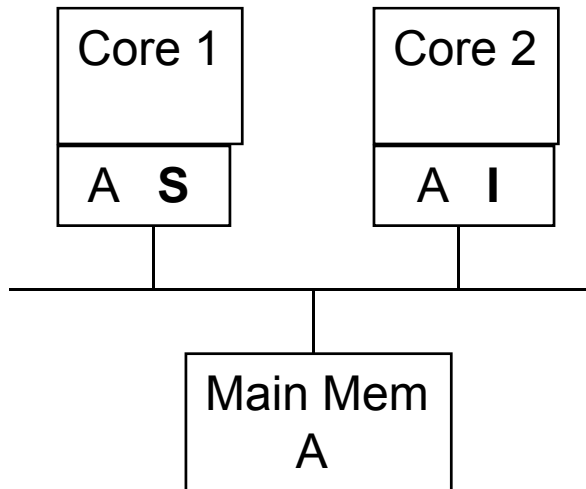
read (hit) or
write (hit or miss)

write-back caching
protocol in black

# A Write-Invalidate CC Protocol
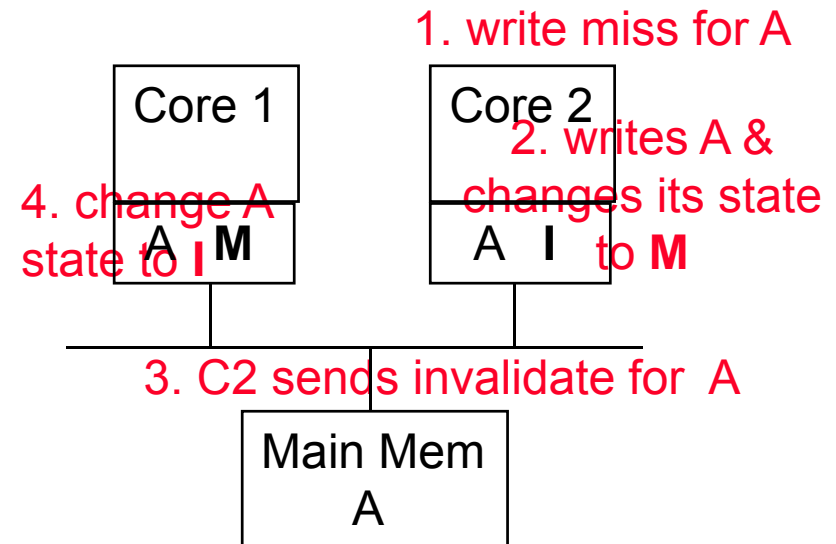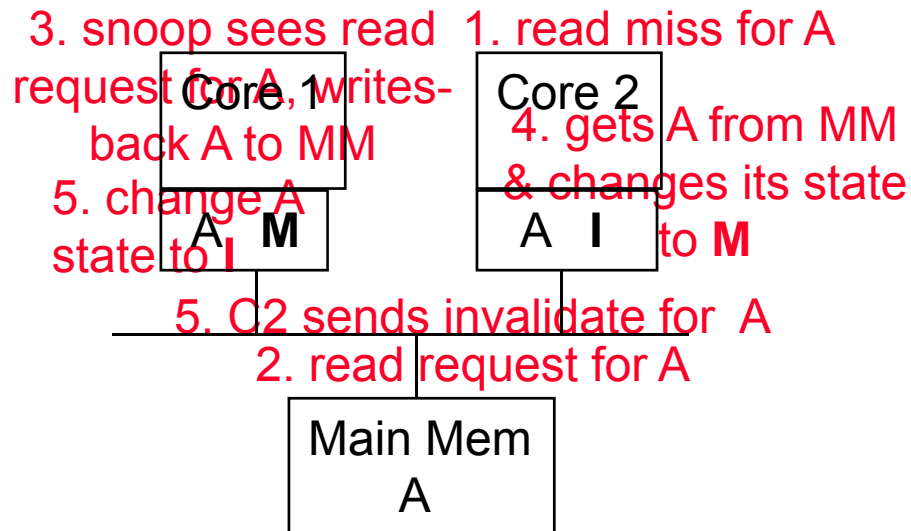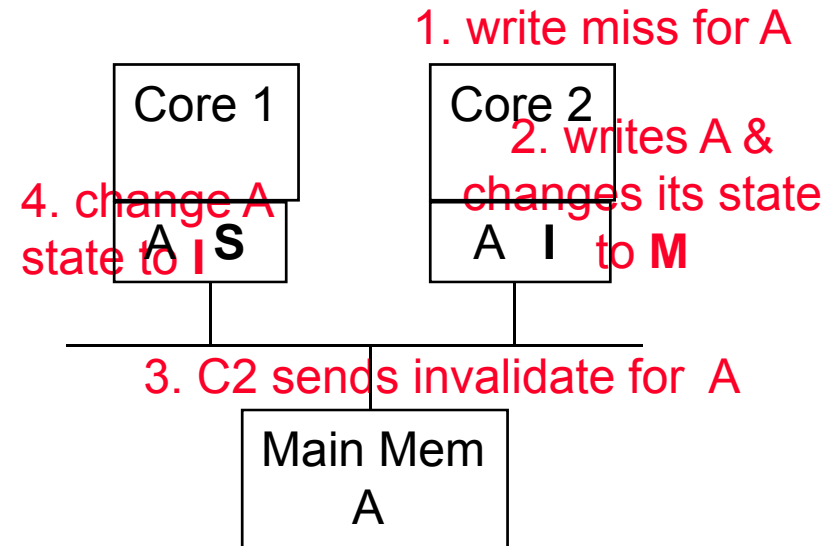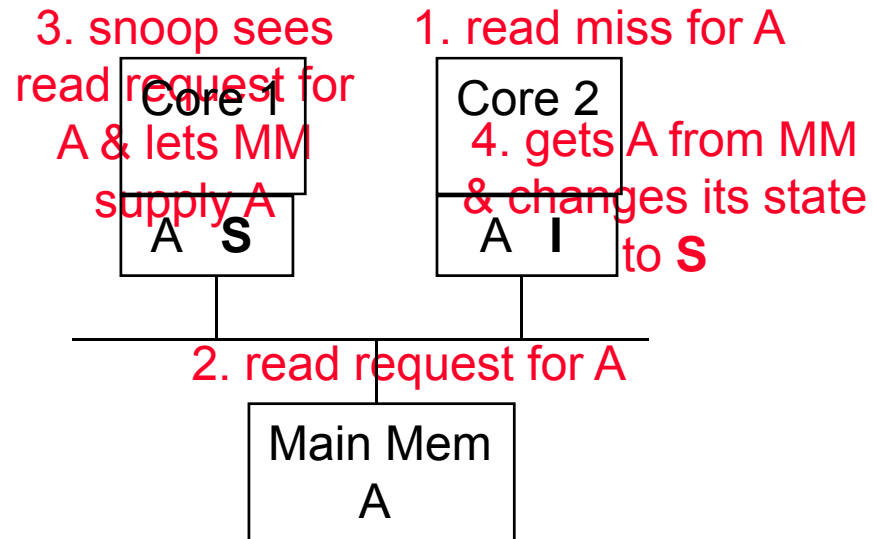
# Write-Invalidate CC Examples

- I = invalid (many), S = shared (many), M = modified (only one)

# Write-Invalidate CC Examples

- I = invalid (many), S = shared (many), M = modified (only one)

**3. snoop sees read request for A & lets MM supply A**

**1. read miss for A**

**1. write miss for A**

| Core 1 | Core 2 |
|--------|--------|
| A  **S** | A  **I** |

4. gets A from MM & changes its state to **S**

| Core 1 | Core 2 |
|--------|--------|
| A  **S** | A  **I** |

4. change A state to **I**

2. writes A & changes its state to **M**

**2. read request for A**

**Main Mem**
A

**3. C2 sends invalidate for A**

**Main Mem**
A

**3. snoop sees read request for A, writes-back A to MM**

**5. change A state to I**

**1. read miss for A**

**1. write miss for A**

| Core 1 | Core 2 |
|--------|--------|
| A  **M** | A  **I** |

4. gets A from MM & changes its state to **M**

| Core 1 | Core 2 |
|--------|--------|
| A  **M** | A  **I** |

4. change A state to **I**

2. writes A & changes its state to **M**

**5. C2 sends invalidate for A**

**2. read request for A**

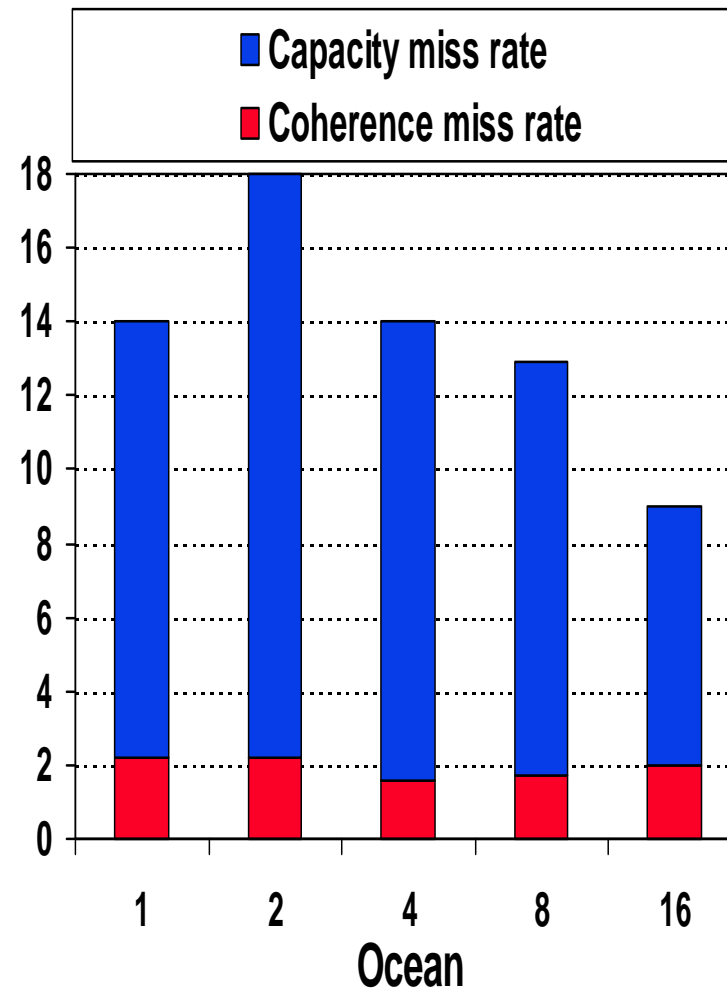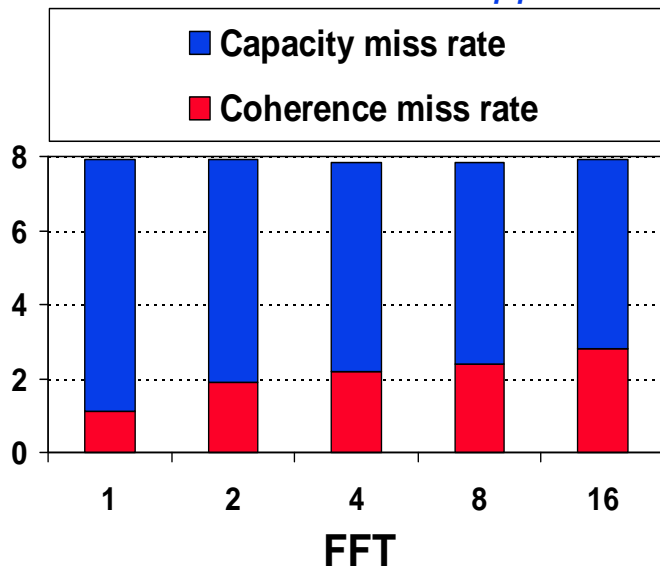**Main Mem**
A

**3. C2 sends invalidate for A**

**Main Mem**
A

# Data Miss Rates

❑ Shared data has lower spatial and temporal locality

● Share data misses often dominate cache behavior even though they may only be 10% to 40% of the data accesses
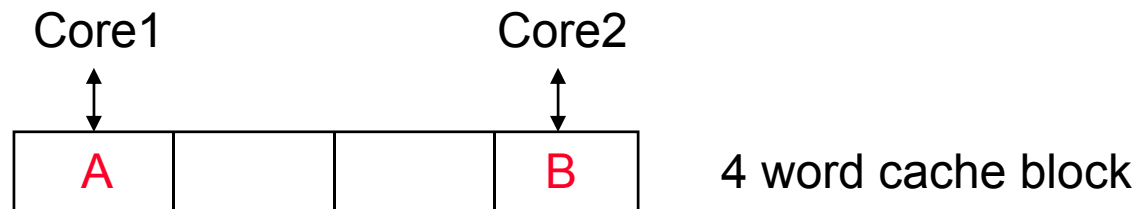
64KB 2-way set associative data cache with 32B blocks

Hennessy & Patterson, *Computer Architecture: A Quantitative Approach*

# Block Size Effects

❏ Writes to one word in a multi-word block mean that the full block is invalidated

❏ Multi-word blocks can also result in false sharing: when two cores are writing to two different variables that happen to fall in the same cache block

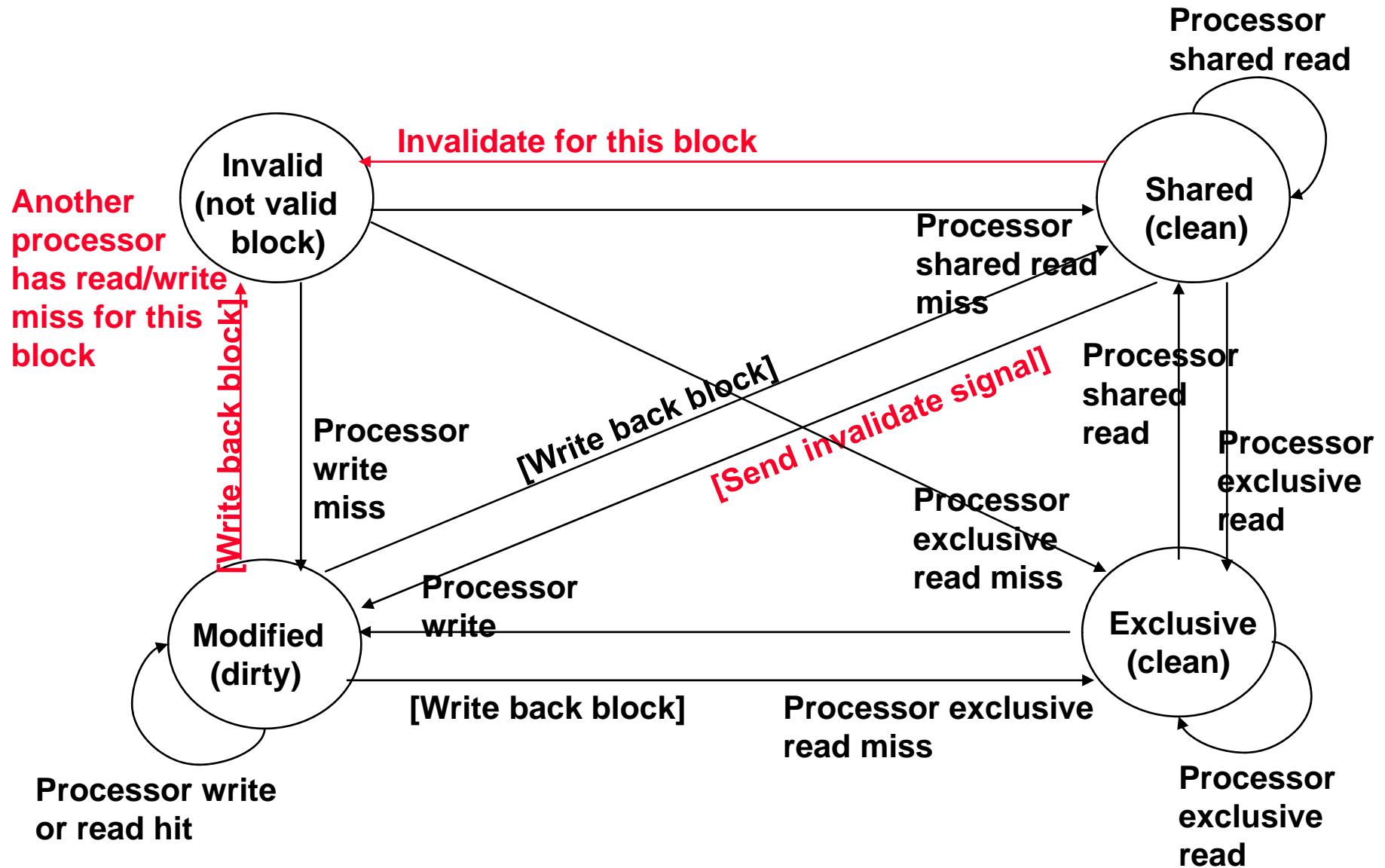- With write-invalidate false sharing increases cache miss rates

Core1                    Core2

| A | | | B | 4 word cache block

❏ Compilers can help reduce false sharing by allocating highly correlated data to the same cache block

# Other Coherence Protocols

❑ There are many variations on cache coherence protocols

❑ Another write-invalidate protocol used in the Pentium 4 (and many other processors) is MESI with four states:

- Modified – same

- Exclusive – only one copy of the shared data is allowed to be cached; memory has an up-to-date copy

  - Since there is only one copy of the block, write hits don't need to send invalidate signal

- Shared – multiple copies of the shared data may be cached (i.e., data permitted to be cached with more than one processor); memory has an up-to-date copy

- Invalid – same

# MESI Cache Coherency Protocol

# Summary:  Improving Cache Performance

## 0. Reduce the time to hit in the cache

- smaller cache

- direct mapped cache

- smaller blocks

- for writes

  - no write allocate – no "hit" on cache, just write to write buffer

  - write allocate – to avoid two cycles (first check for hit, then write) pipeline writes via a delayed write buffer to cache

## 1. Reduce the miss rate

- bigger cache

- more flexible placement (increase associativity)

- larger blocks (16 to 64 bytes typical)

- victim cache – small buffer holding most recently discarded blocks

# Summary: Improving Cache Performance

## 2. Reduce the miss penalty

- smaller blocks

- use a write buffer to hold dirty blocks being replaced so don't have to wait for the write to complete before reading

- check write buffer (and/or victim cache) on read miss – may get lucky

- for large blocks fetch critical word first

- use multiple cache levels – L2 cache not tied to CPU clock rate

- faster backing store/improved memory bandwidth
  - wider buses
  - memory interleaving, DDR SDRAMs

# Summary: The Cache Design Space

❑ Several interacting dimensions

- cache size

- block size

- associativity

- replacement policy

- write-through vs write-back

- write allocation

**Cache Size**

**Associativity**

**Block Size**

❑ The optimal choice is a compromise

- depends on access characteristics
  - workload
  - use (I-cache, D-cache, TLB)
- depends on technology / cost

**Bad**

**Good**   Factor A          Factor B

**Less**                          **More**

❑ Simplicity often wins

# Measuring Cache Performance

❏ Assuming cache hit costs are included as part of the normal CPU execution cycle, then

$$\text{CPU time} = IC \times CPI \times CC$$

$$= IC \times \underbrace{(CPI_{ideal} + \text{Memory-stall cycles})}_{CPI_{stall}} \times CC$$

❏ Memory-stall cycles come from cache misses (a sum of read-stalls and write-stalls)

$$\text{Read-stall cycles} = \text{reads/program} \times \text{read miss rate} \times \text{read miss penalty}$$

$$\text{Write-stall cycles} = (\text{writes/program} \times \text{write miss rate} \times \text{write miss penalty})$$
$$+ \text{ write buffer stalls}$$

❏ For write-through caches, we can simplify this to

$$\text{Memory-stall cycles} = \text{accesses/program} \times \text{miss rate} \times \text{miss penalty}$$

# Impacts of Cache Performance

❑ Relative cache penalty increases as processor performance improves (faster clock rate and/or lower CPI)

  ● The memory speed is unlikely to improve as fast as processor cycle time. When calculating $CPI_{stall}$, the cache miss penalty is measured in *processor* clock cycles needed to handle a miss

  ● The lower the $CPI_{ideal}$, the more pronounced the impact of stalls

❑ A processor with a $CPI_{ideal}$ of 2, a 100 cycle miss penalty, 36% load/store instr's, and 2% I$ and 4% D$ miss rates

$$\text{Memory-stall cycles} = 2\% \times 100 + 36\% \times 4\% \times 100 = 3.44$$

$$\text{So} \quad CPI_{stalls} = 2 + 3.44 = \mathbf{5.44}$$

more than twice the $CPI_{ideal}$ !

❑ What if the $CPI_{ideal}$ is reduced to 1? 0.5? 0.25?

❑ What if the D$ miss rate went up 1%? 2%?

❑ What if the processor clock rate is doubled (doubling the miss penalty)?

# Average Memory Access Time (AMAT)

❑ A larger cache will have a longer access time.  An increase in hit time will likely add another stage to the pipeline.  At some point the increase in hit time for a larger cache will overcome the improvement in hit rate leading to a decrease in performance.

❑ Average Memory Access Time (AMAT) is the average to access memory considering both hits and misses

AMAT =  Time for a hit  +  Miss rate x Miss penalty

❑ What is the AMAT for a processor with a 20 psec clock, a miss penalty of 50 clock cycles, a miss rate of 0.02 misses per instruction and a cache access time of 1 clock cycle?

# Virtual Memory

❑ Use main memory as a "cache" for secondary memory

- Allows efficient and <span style="color:red">safe</span> sharing of memory among multiple programs
- Provides the ability to easily run programs larger than the size of physical memory
- Simplifies loading a program for execution by providing for code relocation (i.e., the code can be loaded anywhere in main memory)

❑ What makes it work? – again the Principle of Locality

- A program is likely to access a relatively small portion of its address space during any period of time

❑ Each program is compiled into its own address space – a "virtual" address space

- During run-time each <span style="color:red">virtual</span> address must be translated to a <span style="color:red">physical</span> address (an address in main memory)

# Two Programs Sharing Physical Memory

❑ A program's address space is divided into pages (all one fixed size) or segments (variable sizes)

- The starting location of each page (either in main memory or in secondary memory) is contained in the program's page table

Program 1
virtual address space

main memory

Program 2
virtual address space

# Address Translation

❑ A virtual address is translated to a physical address by a combination of hardware and software

Virtual Address (VA)

31  30                          . . .                    12  11        . . .        0

| Virtual page number | Page offset |
|---|---|

Translation

| Physical page number | Page offset |
|---|---|

29                    . . .                    12  11                    0

Physical Address (PA)

❑ So each memory request *first* requires an address translation from the virtual space to the physical space

- A virtual memory miss (i.e., when the page is not in physical memory) is called a page fault

# Address Translation Mechanisms



Virtual page #    Offset

Physical page #

Offset

Page table register

Physical page base addr

V

1
1
1
1
1
1
0
1
0
1
0

**Page Table**
(in main memory)

**Main memory**

**Disk storage**

# Virtual Addressing with a Cache

❑ Thus it takes an *extra* memory access to translate a VA to a PA

```
            VA            PA              miss
┌──────┐      ┌────────┐     ┌────────┐        ┌────────┐
│      │─────▶│ Trans- │────▶│        │───────▶│  Main  │
│ CPU  │      │ lation │     │ Cache  │        │ Memory │
│      │◀──┐  └────────┘  ┌─▶│        │◀──┐    │        │
└──────┘   │         hit │  └────────┘   │    └────────┘
           │             │     │         │
           └─────────────┴─────┘         │
                  data      ▼            │
```

❑ This makes memory (cache) accesses very expensive (if every access was really *two* accesses)

❑ The hardware fix is to use a Translation Lookaside Buffer (TLB) – a small cache that keeps track of recently used address mappings to avoid having to do a page table lookup

# Making Address Translation Fast

Virtual page #



Page table register

**TLB**

Physical page base addr

V   Tag

| 1 | | |
| 1 | | |
| 1 | | |
| 0 | | |
| 1 | | |

Physical page base addr

V

| 1 | |
| 1 | |
| 1 | |
| 1 | |
| 1 | |
| 1 | |
| 0 | |
| 1 | |
| 0 | |
| 1 | |
| 0 | |

**Page Table**
(in physical memory)

**Main memory**

**Disk storage**

# Translation Lookaside Buffers (TLBs)

❑ Just like any other cache, the TLB can be organized as fully associative, set associative, or direct mapped

| V | Virtual Page # | Physical Page # | Dirty | Ref | Access |
|---|----------------|-----------------|-------|-----|--------|
|   |                |                 |       |     |        |

❑ TLB access time is typically smaller than cache access time (because TLBs are much smaller than caches)

● TLBs are typically not more than 512 entries even on high end machines

# A TLB in the Memory Hierarchy



❑ A TLB miss – is it a page fault or merely a TLB miss?

- If the page is loaded into main memory, then the TLB miss can be handled (in hardware or software) by loading the translation information from the page table into the TLB

  - Takes 10's of cycles to find and load the translation info into the TLB

- If the page is not in main memory, then it's a true page fault

  - Takes 1,000,000's of cycles to service a page fault

❑ TLB misses are much more frequent than true page faults

# TLB Event Combinations

| TLB | Page Table | Cache | Possible?  Under what circumstances? |
|---|---|---|---|
| Hit | Hit | Hit | |
| Hit | Hit | Miss | |
| Miss | Hit | Hit | |
| Miss | Hit | Miss | |
| Miss | Miss | Miss | |
| Hit | Miss | Miss/ Hit | |
| Miss | Miss | Hit | |

# TLB Event Combinations

| TLB | Page Table | Cache | Possible?  Under what circumstances? |
|---|---|---|---|
| Hit | Hit | Hit | Yes – what we want! |
| Hit | Hit | Miss | Yes – although the page table is not checked if the TLB hits |
| Miss | Hit | Hit | Yes – TLB miss, PA in page table |
| Miss | Hit | Miss | Yes – TLB miss, PA in page table, but data not in cache |
| Miss | Miss | Miss | Yes – page fault |
| Hit | Miss | Miss/ Hit | Impossible – TLB translation not possible if page is not present in memory |
| Miss | Miss | Hit | Impossible – data not allowed in cache if page is not in memory |

# Handling a TLB Miss

❑ Consider a TLB miss for a page that is present in memory (i.e., the Valid bit in the page table is set)

● A TLB miss (or a page fault exception) must be asserted by the end of the same clock cycle that the memory access occurs so that the next clock cycle will begin exception processing

| Register | Description |
|----------|-------------|
| EPC | Where to restart after exception |
| Cause | Cause of exception |
| BadVAddr | Address that caused exception |
| Index | Location in TLB to be read/written |
| Random | Pseudorandom location in TLB |
| EntryLo | Physical page address and flags |
| EntryHi | Virtual page address |
| Context | Page table address & page number |

# A MIPS Software TLB Miss Handler

- When a TLB miss occurs, the hardware saves the address that caused the miss in `BadVAddr` and transfers control to 8000 0000$_{hex}$, the location of the TLB miss handler

```
TLBmiss:
 mfc0   $k1, Context    #copy addr of PTE into $k1
 lw     $k1, 0($k1)     #put PTE into $k1
 mtc0   $k1, EntryLo    #put PTE into EntryLo
 tlbwr                  #put EntryLo into TLB
                        #     at Random
 eret                   #return from exception
```

- `tlbwr` copies from `EntryLo` into the TLB entry selected by the control register `Random`

- A TLB miss takes about a dozen clock cycles to handle

# Some Virtual Memory Design Parameters

|  | Paged VM | TLBs |
|---|---|---|
| Total size | 16,000 to 250,000 words | 16 to 512 entries |
| Total size (KB) | 250,000 to 1,000,000,000 | 0.25 to 16 |
| Block size (B) | 4000 to 64,000 | 4 to 8 |
| Hit time |  | 0.5 to 1 clock cycle |
| Miss penalty (clocks) | 10,000,000 to 100,000,000 | 10 to 100 |
| Miss rates | 0.00001% to 0.0001% | 0.01% to 1% |

# Two Machines' TLB Parameters

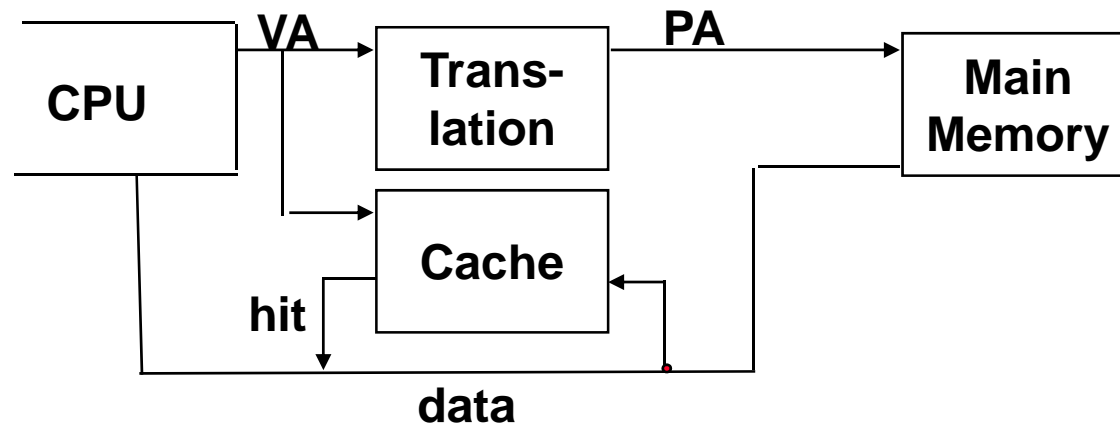| | Intel Nehalem | AMD Barcelona |
|---|---|---|
| Address sizes | 48 bits (vir); 44 bits (phy) | 48 bits (vir); 48 bits (phy) |
| Page size | 4KB | 4KB |
| TLB organization | L1 TLB for instructions and  L1 TLB for data per core; both are 4-way set assoc.; LRU<br><br>L1 ITLB has 128 entries, L2 DTLB has 64 entries<br><br>L2 TLB (unified) is 4-way set assoc.; LRU<br><br>L2 TLB has 512 entries<br><br>TLB misses handled in hardware | L1 TLB for instructions and L1 TLB for data per core; both are fully assoc.; LRU<br><br>L1 ITLB and DTLB each have 48 entries<br><br>L2 TLB for instructions and L2 TLB for data per core; each are 4-way set assoc.; round robin LRU<br><br>Both L2 TLBs have 512 entries<br><br>TLB misses handled in hardware |

# Two Machines' TLB Parameters

| | Intel P4 | AMD Opteron |
|---|---|---|
| TLB organization | 1 TLB for instructions and 1TLB for data | 2 TLBs for instructions and 2 TLBs for data |
| | Both 4-way set associative | Both L1 TLBs fully associative with ~LRU replacement |
| | Both use ~LRU replacement | Both L2 TLBs are 4-way set associative with round-robin LRU |
| | Both have 128 entries | Both L1 TLBs have 40 entries |
| | | Both L2 TLBs have 512 entries |
| | TLB misses handled in hardware | TBL misses handled in hardware |

# Why Not a Virtually Addressed Cache?

❑ A virtually addressed cache would only require address translation on cache misses

```
                VA              PA
CPU  ─────────┬──►  Trans-  ────────►  Main
              │      lation            Memory
              │                         ▲
              └──►   Cache  ◄───────────┘
       hit     ┌──────┘
              ▼
             data
```
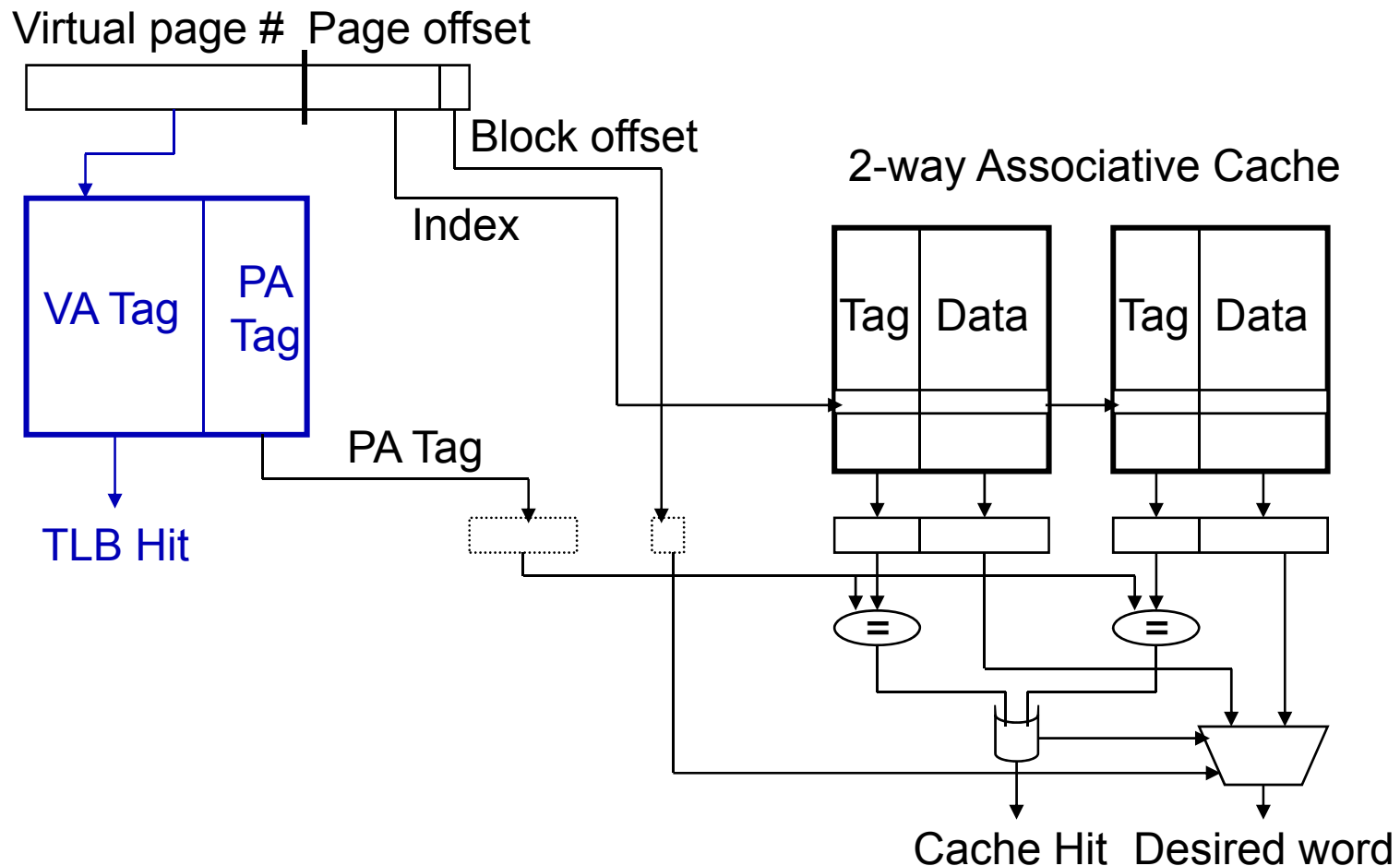
but

- Two programs which are sharing data will have two different virtual addresses for the same physical address – aliasing – so have two copies of the shared data in the cache and two entries in the TLB which would lead to coherence issues
  - Must update all cache entries with the same physical address or the memory becomes inconsistent

# Reducing Translation Time

❑ Can overlap the cache access with the TLB access

- Works when the high order bits of the VA are used to access the TLB while the low order bits are used as index into cache

Virtual page #  Page offset

Block offset

Index

2-way Associative Cache

VA Tag | PA Tag

Tag | Data      Tag | Data

TLB Hit

PA Tag

Cache Hit  Desired word

# The Hardware/Software Boundary

❑ What parts of the virtual to physical address translation is done by or assisted by the hardware?

- Translation Lookaside Buffer (TLB) that caches the recent translations
  - TLB access time is part of the cache hit time
  - May allot an extra stage in the pipeline for TLB access
- Page table storage, fault detection and updating
  - Page faults result in interrupts (precise) that are then handled by the OS
  - Hardware must support (i.e., update appropriately) Dirty and Reference bits (e.g., ~LRU) in the Page Tables
- Disk placement
  - Bootstrap (e.g., out of disk sector 0) so the system can service a limited number of page faults before the OS is even loaded

# Summary

- ❑ The Principle of Locality:
  - ● Program likely to access a relatively small portion of the address space at any instant of time.
    - - Temporal Locality: Locality in Time
    - - Spatial Locality: Locality in Space

- ❑ Caches, TLBs, Virtual Memory all understood by examining how they deal with the four questions
  1. Where can entry be placed?
  2. How is entry found?
  3. What entry is replaced on miss?
  4. How are writes handled?

- ❑ Page tables map virtual address to physical address
  - ● TLBs are important for fast translation

# SDR SDRAM
## (Single Data Rate SDRAM)
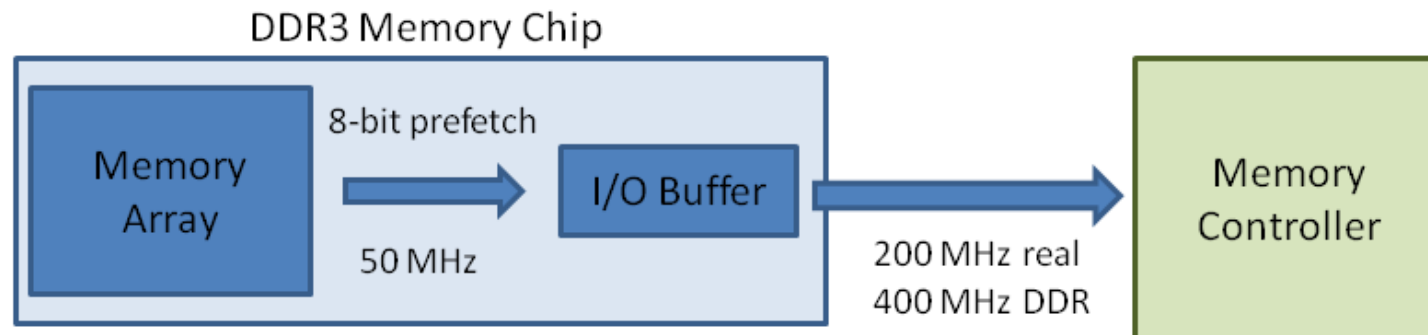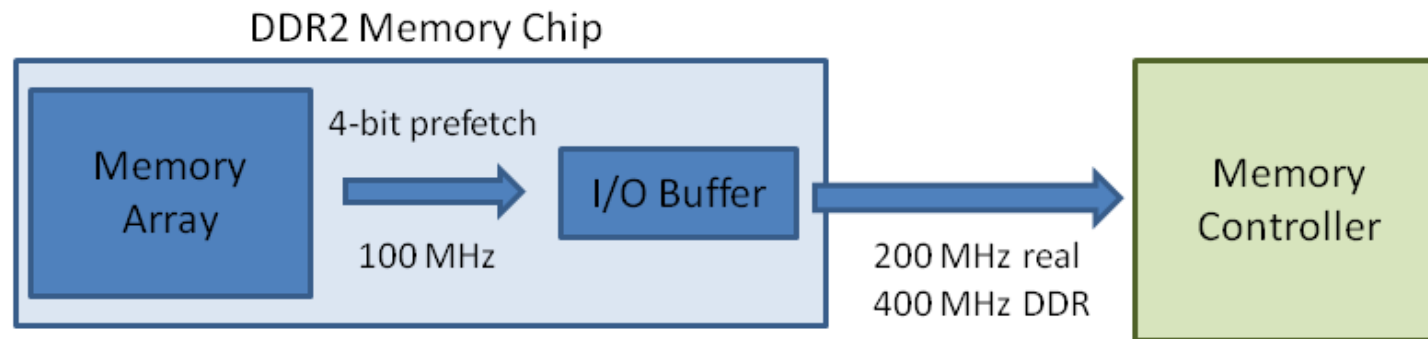## (Synchronous Dynamic Random Access Memory)

- SDRAM is designed to synchronize itself with the CPU timing

- The memory controller knows the exact clock cycle when the requested data will be ready, so the CPU no longer has to wait between memory accesses.

- For example, PC66 SDRAM runs at 66 MT/s, PC100 SDRAM runs at 100 MT/s, PC133 SDRAM runs at 133 MT/s, and so on.

- The I/O, internal and bus clocks are the same. For example, the I/O, internal clock and bus clock of PC133 are all 133 Mhz.

- Single Data Rate means that SDRAM can only read/write one time in a clock cycle. SDRAM have to wait for the completion of the previous command to be able to do another read/write.

# DDR SDRAM

DDR SDRAM (Double Data Rate SDRAM):
The next generation of SDRAM is DDR, which achieves greater bandwidth than the preceding single data rate SDRAM by transferring data on the rising and falling edges of the clock signal (double pumped). Effectively, it doubles the transfer rate without increasing the frequency of the clock. The transfer rate of DDR SDRAM is the double of SDR SDRAM without changing the internal clock. In DDR SDRAM, as the first generation of DDR memory, the prefetch buffer is 2n, which is the double of  SDR SDRAM. The transfer rate of DDR is between 266~400 MT/s. DDR266 and DDR400 are of this type.

## DDR Memory Chip

| | | |
|---|---|---|
| **Memory Array** | 2-bit prefetch → | **I/O Buffer** |
| | 200 MHz | |

200 MHz real
400 MHz DDR → **Memory Controller**

## DDR2 Memory Chip

| | | |
|---|---|---|
| **Memory Array** | 4-bit prefetch → | **I/O Buffer** |
| | 100 MHz | |

200 MHz real
400 MHz DDR → **Memory Controller**

## DDR3 Memory Chip

| | | |
|---|---|---|
| **Memory Array** | 8-bit prefetch → | **I/O Buffer** |
| | 50 MHz | |

200 MHz real
400 MHz DDR → **Memory Controller**

Prefetch

**1N** — Memory array (100-150MHz) ↔ I/O (100-150MHz) — SDR
100-150Mbps

**2N** — Memory array (100-200MHz), Memory array (100-200MHz) ↔ MUX ↔ I/O (100-200MHz) — DDR/LPDDR
200-400Mbps

**4N** — Memory array (100-266MHz) ×4 ↔ MUX ↔ I/O (200-533MHz) — DDR2/LPDDR2
400-1066Mbps

**8N** — Memory array (100-266MHz) ×4 ↔ MUX ↔ I/O (400-1066MHz) — DDR3/LPDDR3
800-2133Mbps

Bank group 0

**8N** — Memory array (100-266MHz) ×4 ↔ MUX ↔ MUX ↔ I/O (667-1600MHz) — DDR4
1333-3200Mbps

**8N** — Memory array (100-266MHz) ×4 ↔ MUX ↔ MUX

Bank group 1

| DDR SDRAM Standard | Internal rate (MHz) | Bus clock (MHz) | Prefetch | Data rate (MT/s) | Transfer rate (GB/s) | Voltage (V) |
|---|---|---|---|---|---|---|
| SDRAM | 100-166 | 100-166 | 1n | 100-166 | 0.8-1.3 | 3.3 |
| DDR | 133-200 | 133-200 | 2n | 266-400 | 2.1-3.2 | 2.5/2.6 |
| DDR2 | 133-200 | 266-400 | 4n | 533-800 | 4.2-6.4 | 1.8 |
| DDR3 | 133-200 | 533-800 | 8n | 1066-1600 | 8.5-14.9 | 1.35/1.5 |
| DDR4 | 133-200 | 1066-1600 | 8n | 2133-3200 | 17-21.3 | 1.2 |

# DDR2 (Double Data Rate 2 SDRAM)

Its primary benefit is the ability to operate the external data bus twice as fast as DDR SDRAM. This is achieved by improved bus signal. The prefetch buffer of DDR2 is 4 bit(double of DDR SDRAM). DDR2 memory is at the same internal clock speed (133~200MHz) as DDR,  but the transfer rate of DDR2 can reach 533~800 MT/s with the improved I/O bus signal. DDR2 533 and DDR2 800 memory types are on the market.

# DDR3 (Double Data Rate 3 SDRAM)

DDR3 memory reduces 40% power consumption compared to current DDR2 modules, allowing for lower operating currents and voltages (1.5 V, compared to DDR2's 1.8 V or DDR's 2.5 V). The transfer rate of DDR3 is 800~1600 MT/s. DDR3's prefetch buffer width is 8 bit, whereas DDR2's is 4 bit, and DDR's is 2 bit. DDR3 also adds two functions, such as ASR (Automatic Self-Refresh) and SRT (Self-Refresh Temperature). They can make the memory control the refresh rate according to the temperature variation.

# DDR4 (Double Data Rate 4 SDRAM)

DDR4 SDRAM provides the lower operating voltage (1.2V) and higher transfer rate. The transfer rate of DDR4 is 2133~3200 MT/s. DDR4 adds four new Bank Groups technology. Each bank group has the feature of singlehanded operation. DDR4 can process 4 data within a clock cycle, so DDR4's efficiency is better than DDR3 obviously. DDR4 also adds some functions, such as DBI
(Data Bus Inversion), CRC
(Cyclic Redundancy Check) and CA parity. They can enhance DDR4 memory's signal integrity, and improve the stability of data transmission/access.