# Gerarchia di Memoria

# Review:  Major Components of a Computer

**Processor**

Control

Datapath

**Memory**

**Devices**

Input

Output

Cache

Main Memory

Secondary Memory (Disk)

Moore's Law – The number of transistors on integrated circuit chips (1971-2018)

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are linked to Moore's law.

Our World in Data

μProc
55%/year
(2X/1.5yr)

DRAM
7%/year
(2X/10yrs)

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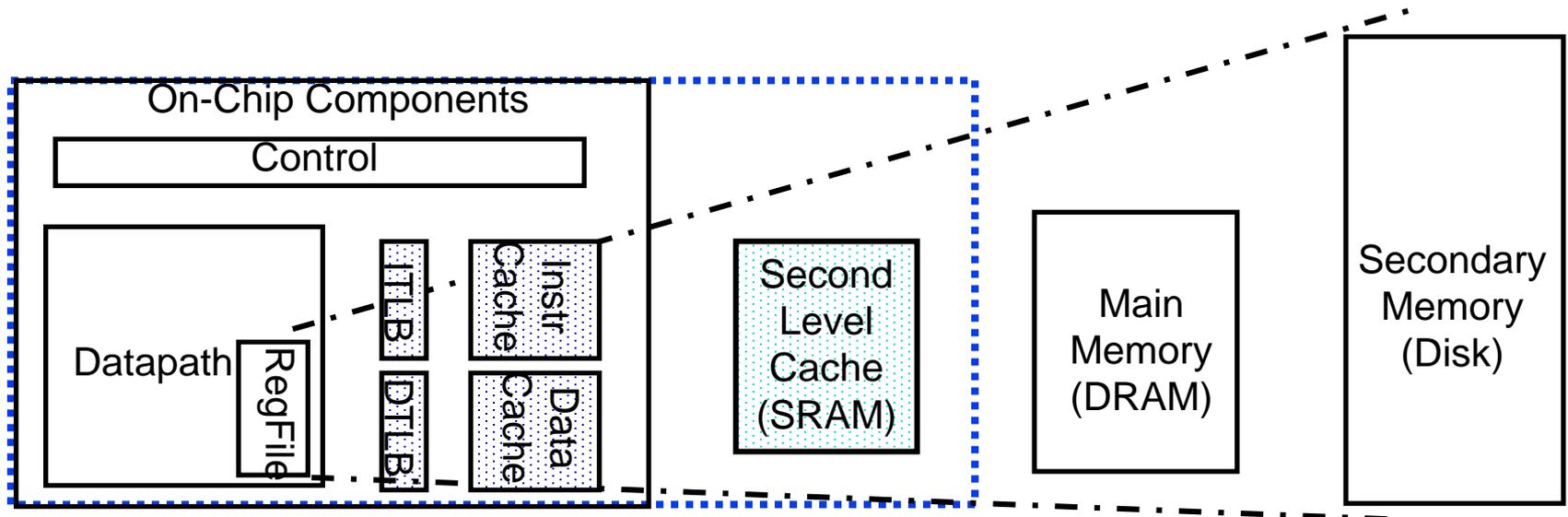
❑ Processor vs DRAM speed disparity continues to grow



❑ 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

Control

Datapath | RegFile | TLB | DTLB | Instr Cache | Data Cache | 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)
  - consumes 1% 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
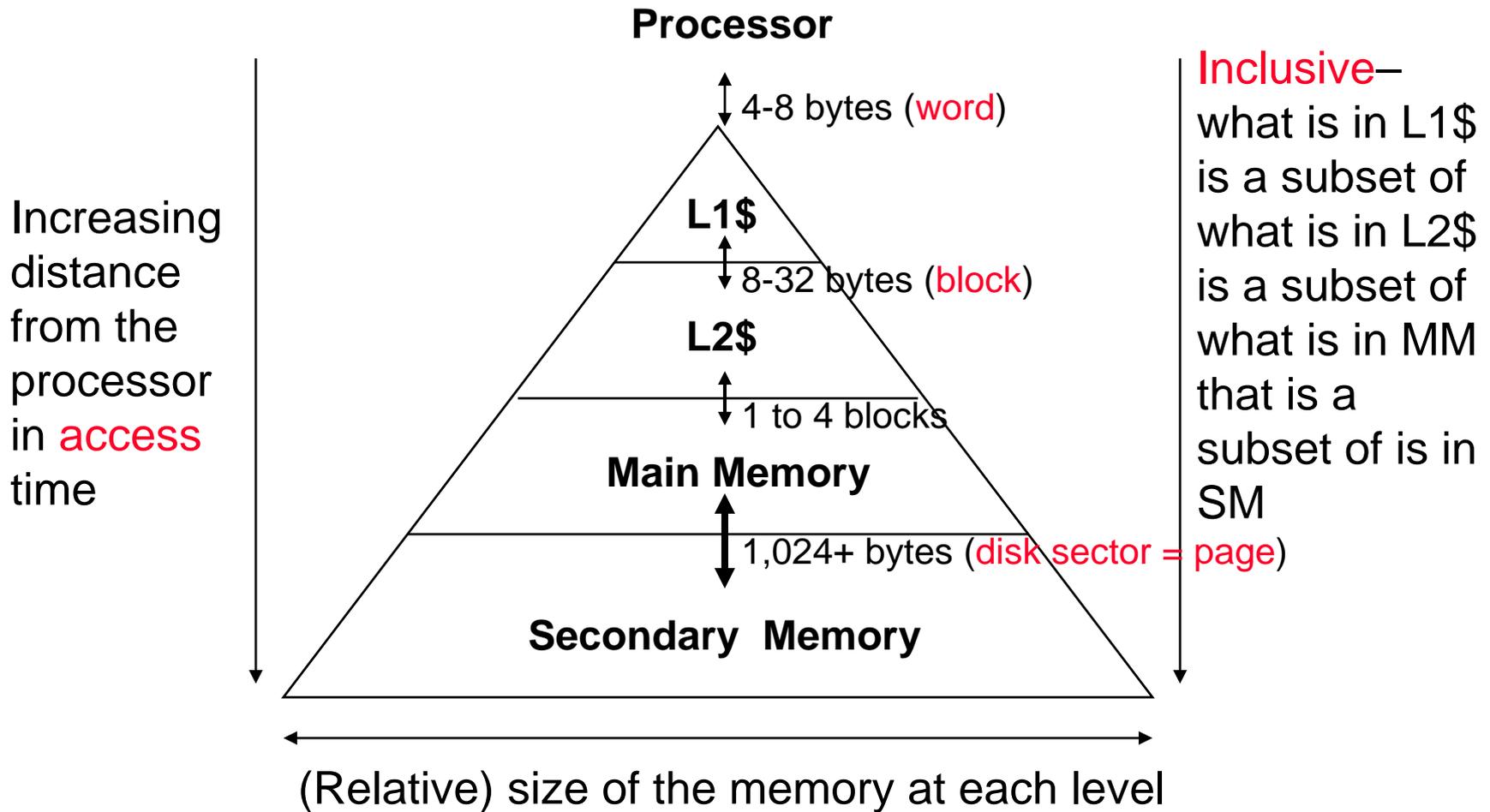
  ⇒ Move blocks consisting of contiguous words closer to the processor

# 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**

(Relative) size of the memory at each level

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

# How is the Hierarchy Managed?

❑ registers ↔ memory
  ❑ by compiler (programmer?)

❑ cache ↔ main memory
  ❑ by the cache controller hardware

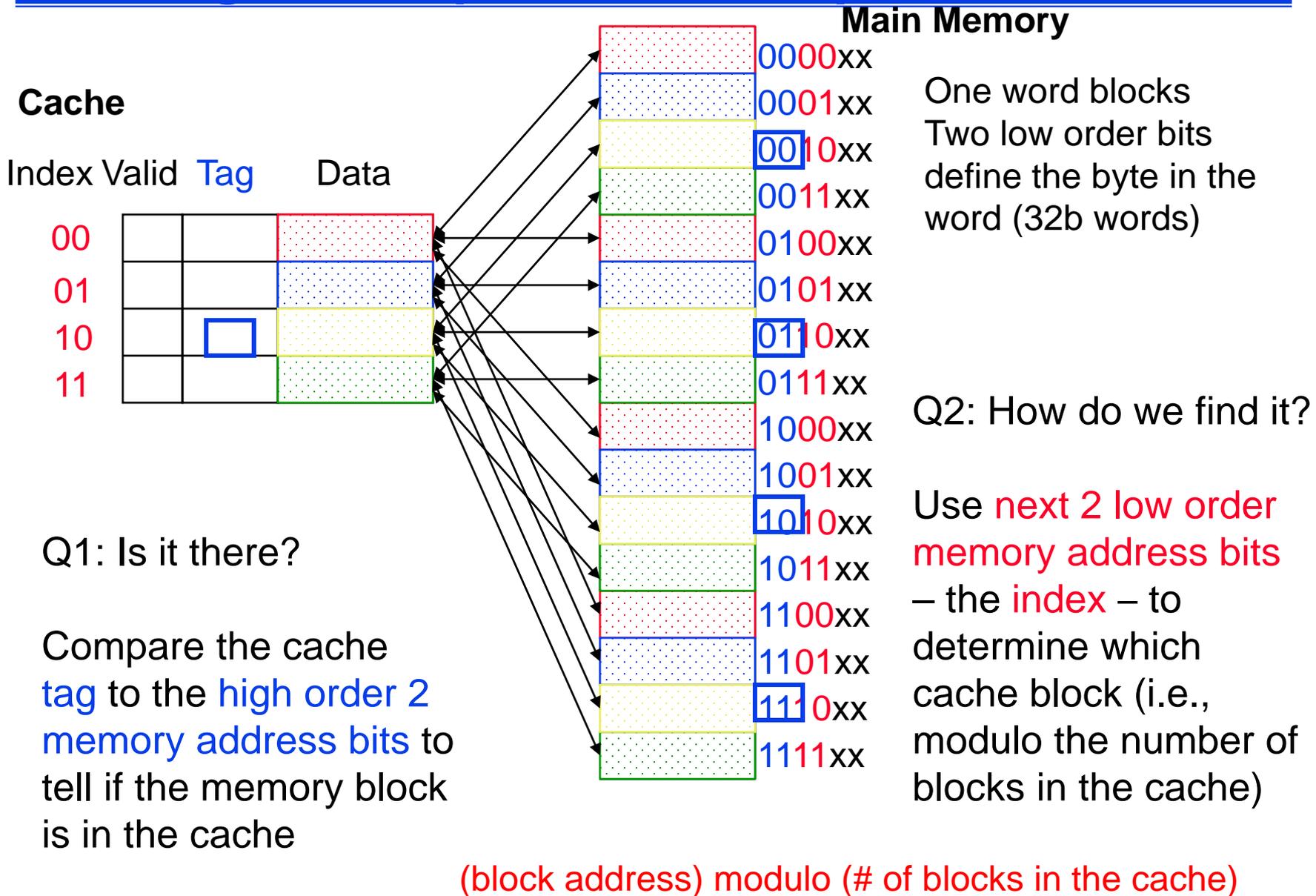❑ main memory ↔ disks
  ❑ by the operating system (virtual memory)
  ❑ virtual to physical address mapping assisted by the hardware (TLB)
  ❑ by the programmer (files)

# 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)

# 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** 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 ~~00~~ ~~Mem(1)~~ 5 ~~Mem(0)~~ 4

| 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 ~~01~~ Mem(5) 15 Mem(4) 14

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

▢ 8 requests, 4 misses

# Miss Rate vs Block Size vs Cache Size



**Block size (bytes)**

❑ 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

$$2^n \text{ x (block size + tag field size + valid field size)}$$

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

$$2\text{^}10 \text{ x } (4\text{x}32 + 18 + 1) = 2\text{^}10 \text{ x } 147 = 147\text{Kbits}$$

# 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

bus

32-bit data
&
32-bit addr
per cycle

DRAM
Memory

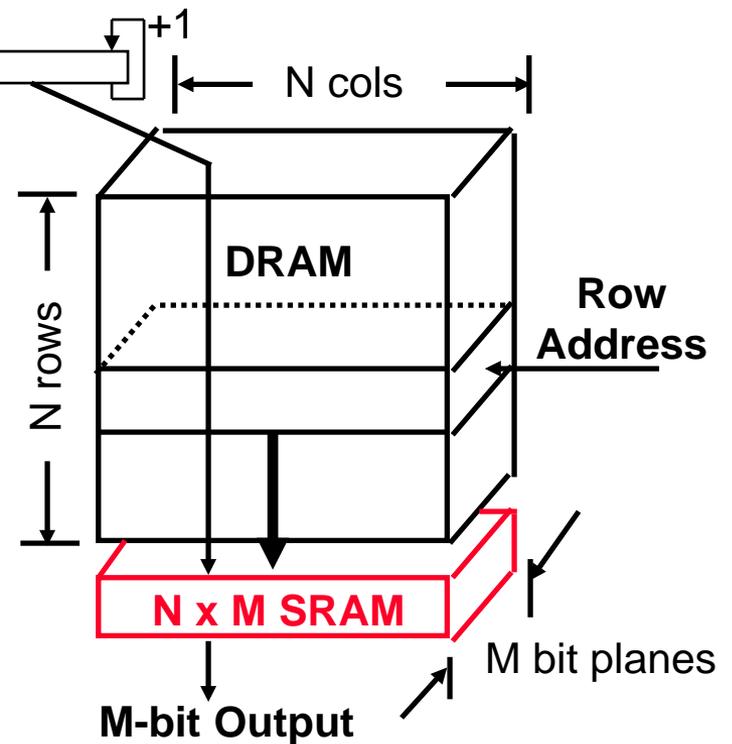One word wide organization (one word wide bus and one word wide memory)

❑ 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

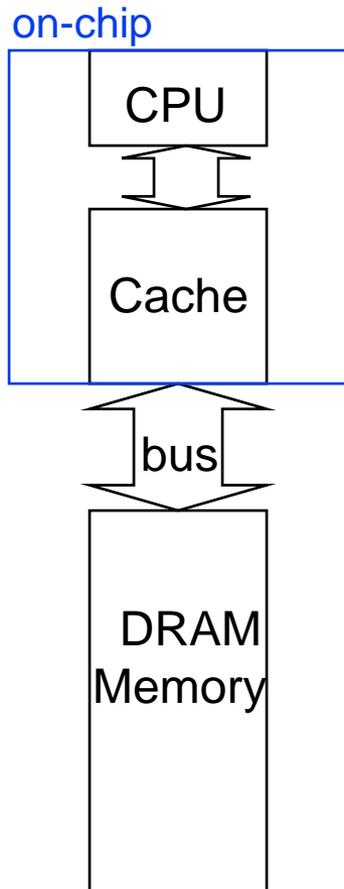  - The memory bus clock controls transfer of successive words in the burst

**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

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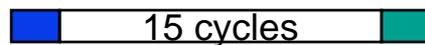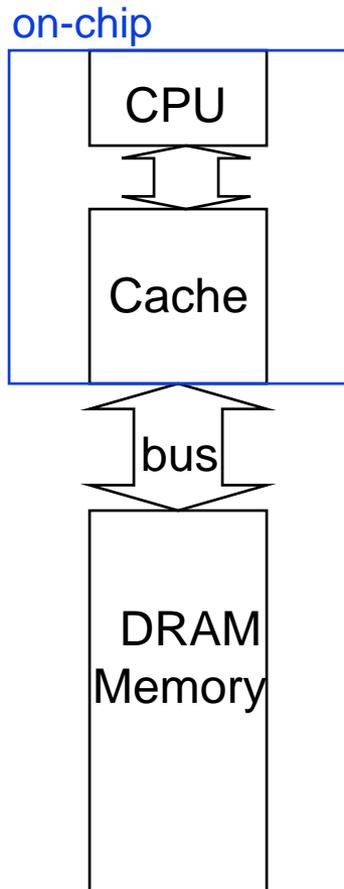
| | |
|---|---|
| 1 | cycle to send 1$^{st}$ address |
| 4 x 15 = 60 | cycles to read DRAM |
| 1 | cycles to return last data word |
| 62 | total clock cycles miss penalty |

15 cycles

15 cycles

15 cycles

15 cycles

❑ Number of bytes transferred per clock cycle (bandwidth) for a single miss is
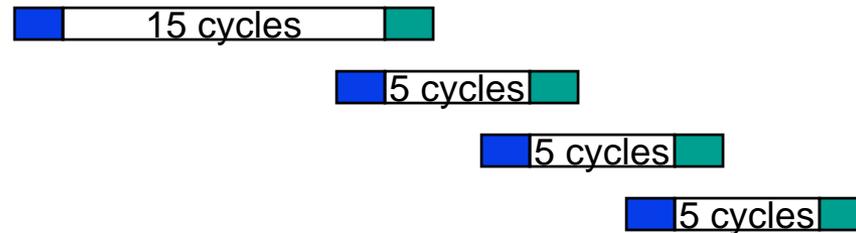
(4 x 4)/62 = 0.258  bytes per clock

# One Word Wide Bus, Four Word Blocks

on-chip

CPU

Cache

bus

DRAM
Memory

❑ 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 |

15 cycles

5 cycles

5 cycles
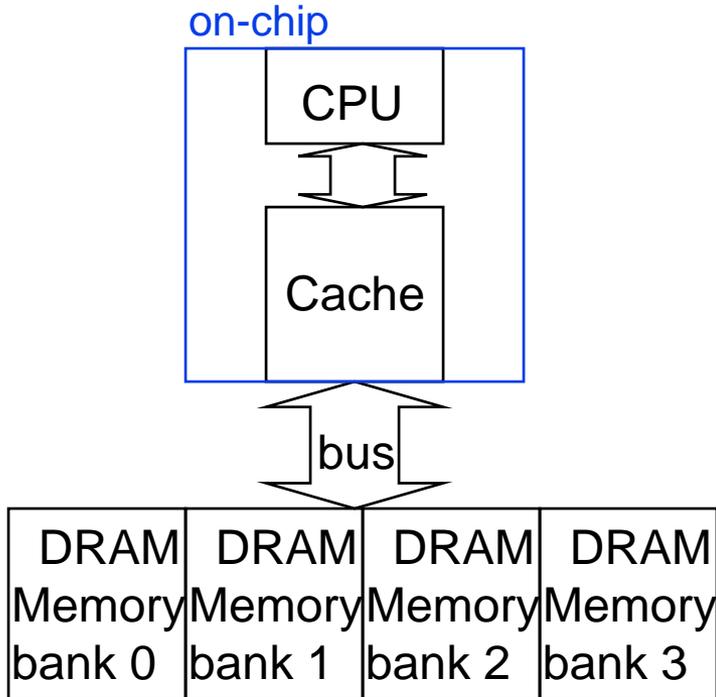
5 cycles

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

$(4 \times 4)/32 = 0.5$    bytes per clock

# Interleaved Memory, One Word Wide Bus

on-chip

CPU

Cache

bus

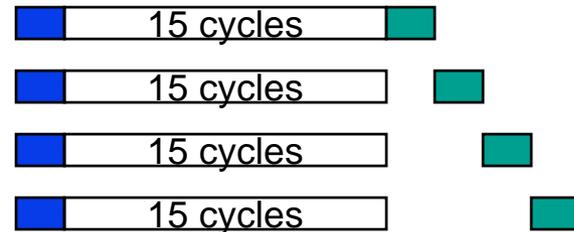| DRAM Memory bank 0 | DRAM Memory bank 1 | DRAM Memory bank 2 | DRAM Memory bank 3 |

❑ 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

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

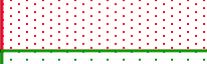Start with an empty cache - all blocks initially marked as not valid

0   4   0   4   0   4   0   4

**0** miss

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

**4** miss

| ~~01~~ ~~00~~ | Mem(~~0~~) 4 |
|----|--------|
|    |        |
|    |        |
|    |        |

**0** miss

| ~~00~~ ~~01~~ | Mem(4) 0 |
|----|--------|
|    |        |
|    |        |
|    |        |

**4** miss

| ~~01~~ ~~00~~ | Mem(~~0~~) 4 |
|----|--------|
|    |        |
|    |        |
|    |        |

**0** miss

| ~~00~~ ~~01~~ | Mem(4) 0 |
|----|--------|
|    |        |
|    |        |
|    |        |

**4** miss

| ~~01~~ ~~00~~ | Mem(0) 4 |
|----|--------|
|    |        |
|    |        |
|    |        |

**0** miss

| ~~00~~ ~~01~~ | Mem(4) 0 |
|----|--------|
|    |        |
|    |        |
|    |        |

**4** miss

| ~~01~~ ~~00~~ | Mem(0) 4 |
|----|--------|
|    |        |
|    |        |
|    |        |

▢ 8 requests, 8 misses

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

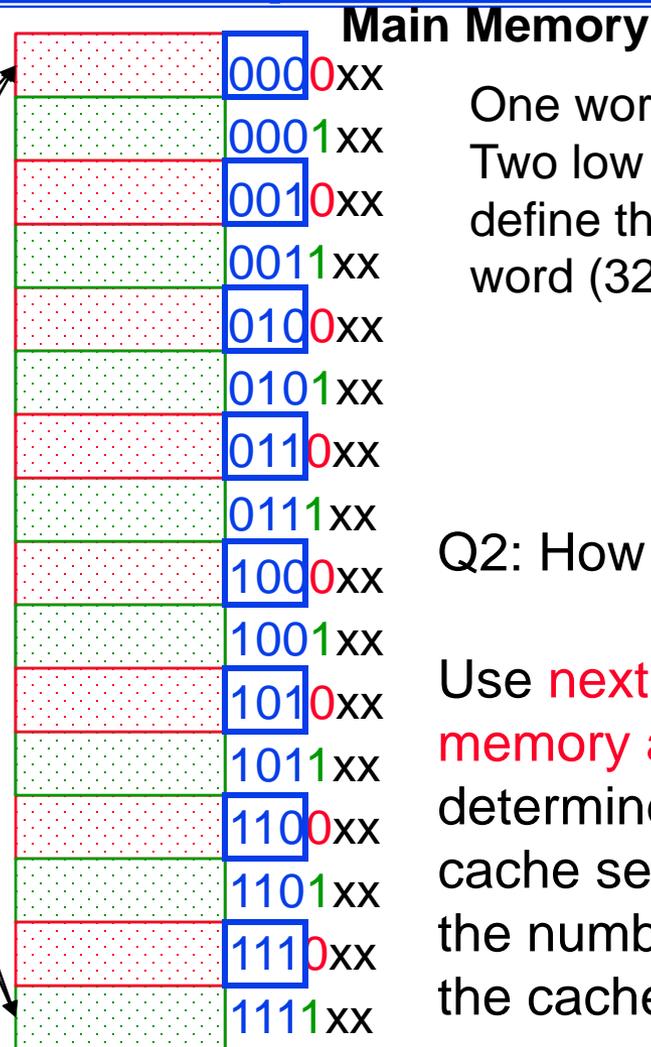# Set Associative Cache Example

**Cache**

Way  Set  V    Tag        Data

0        0
         1

1        0
         1

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

**Main Memory**

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 1 low order
memory address bit to
determine which
cache set (i.e., modulo
the number of 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** 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!
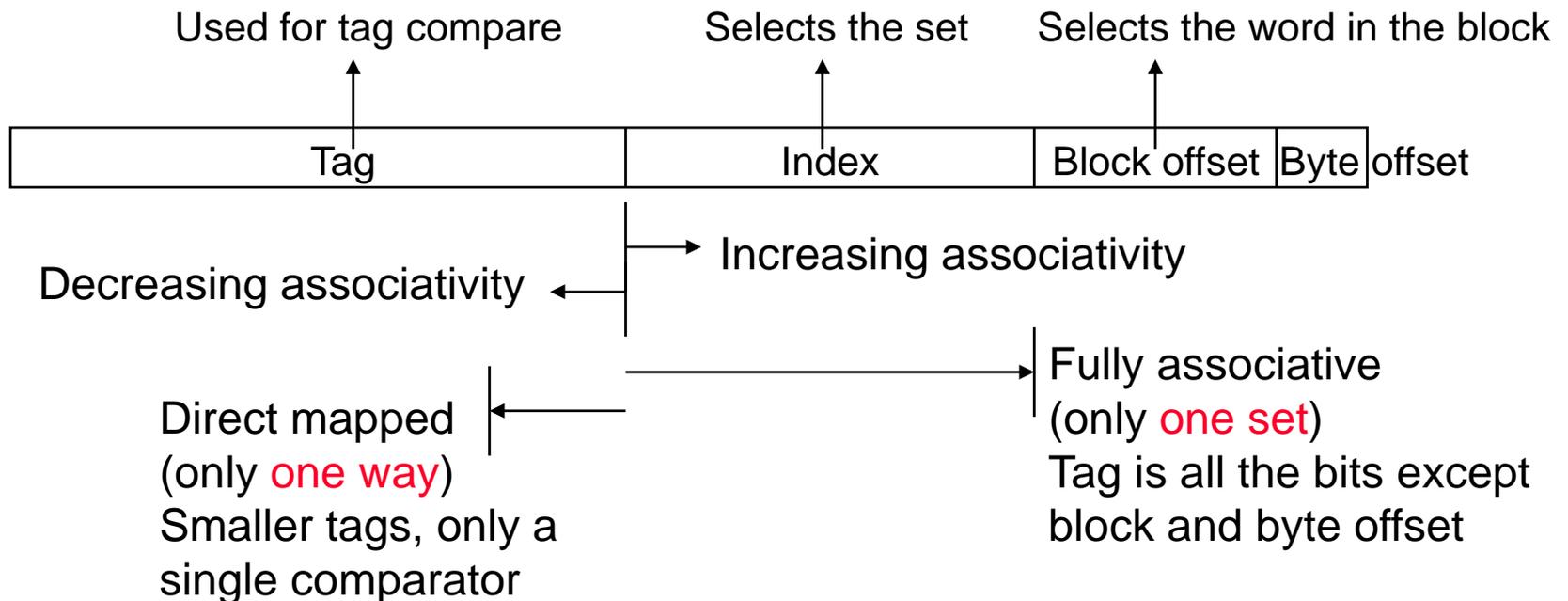
# Four-Way Set Associative Cache

❑ $2^8 = 256$ 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
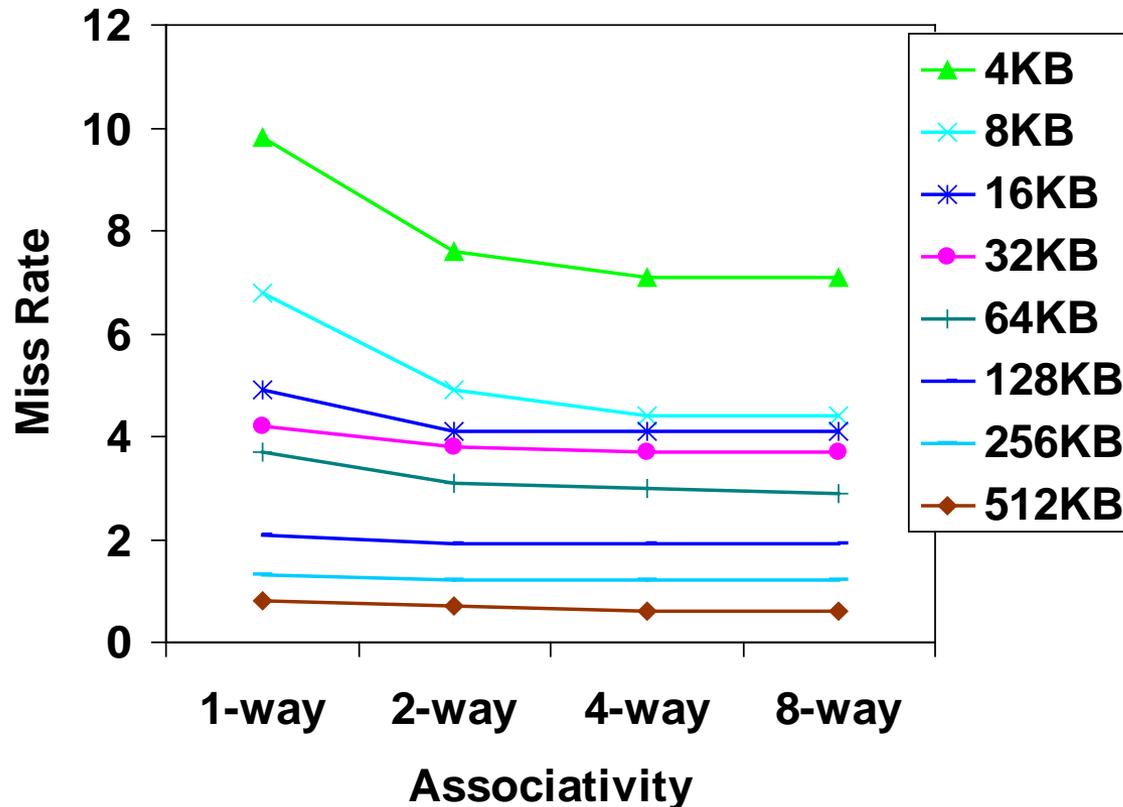
Used for tag compare | Selects the set | Selects the word in the block

| 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



❑ 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
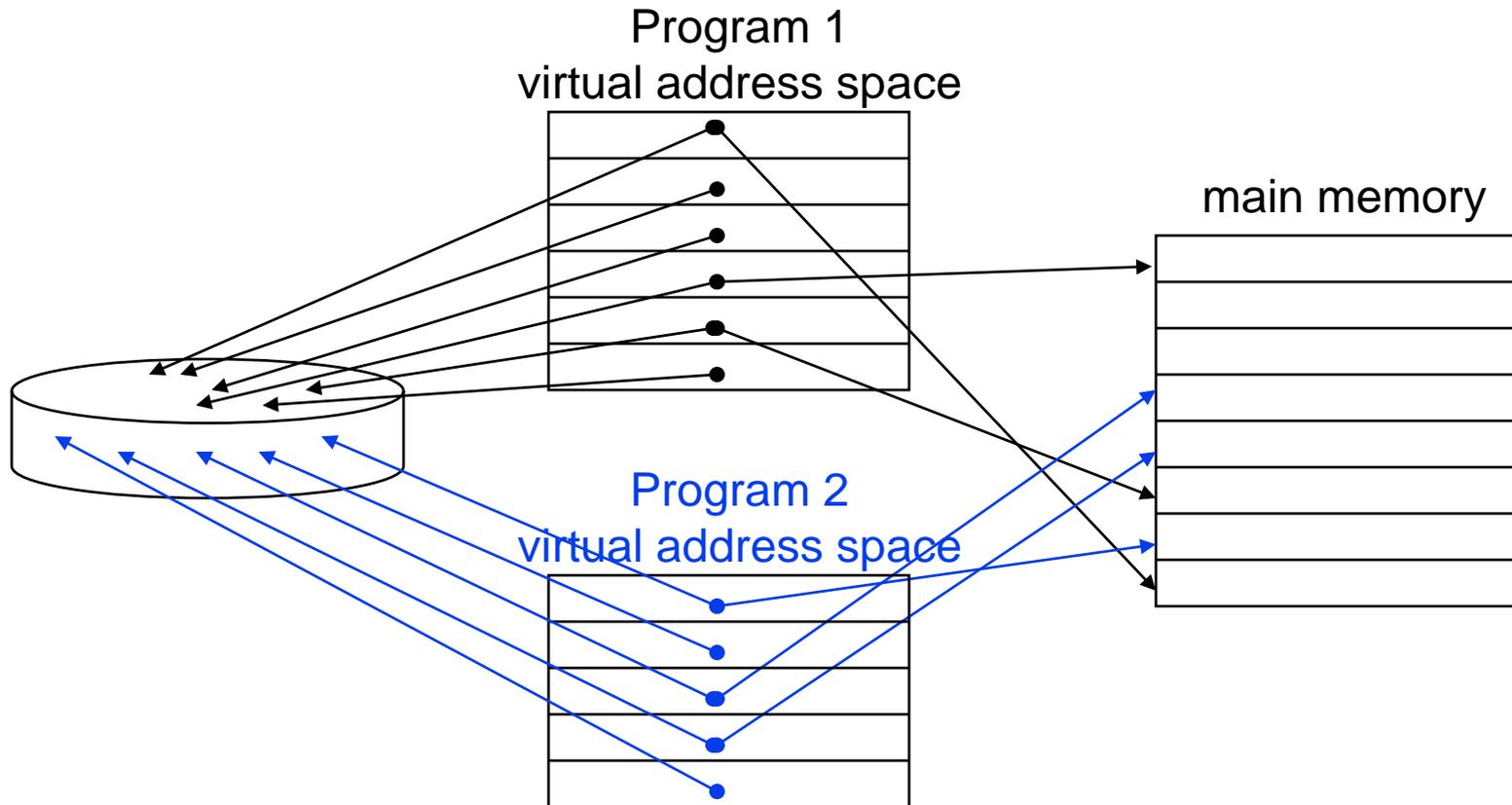
# Two Machines' Cache Parameters

| | 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 |

# 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)
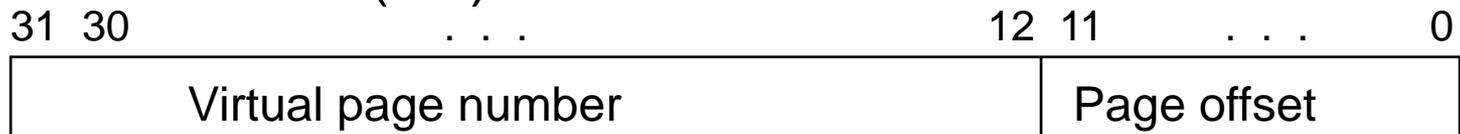
# Two Programs Sharing Physical Memory

❑ 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

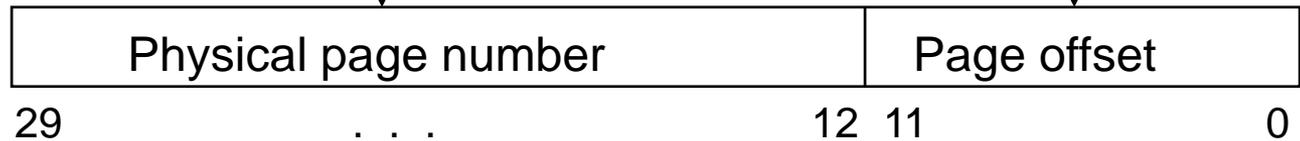Program 1
virtual address space

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)

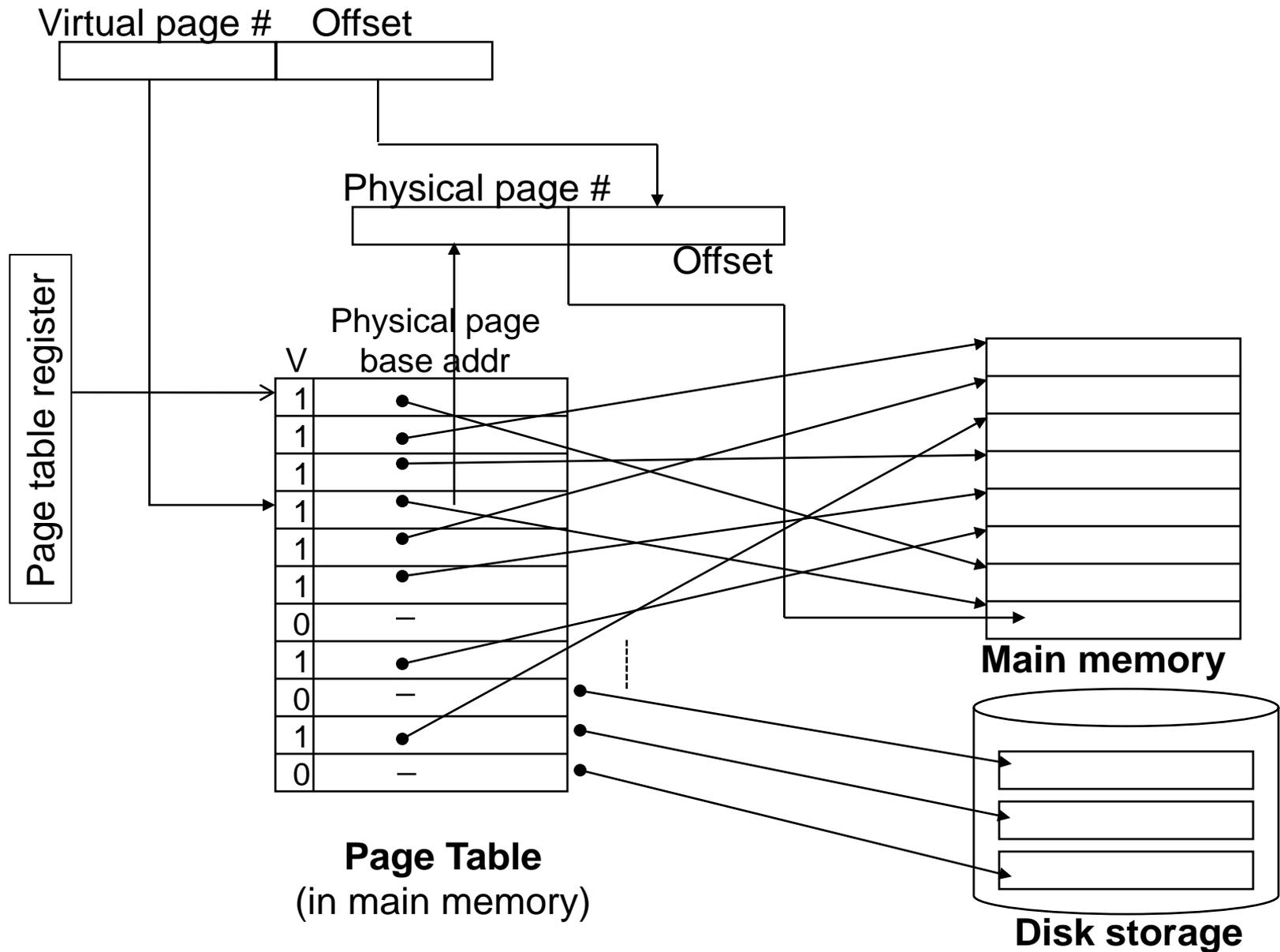| 31 30 . . . 12 | 11 . . . 0 |
|---|---|
| Virtual page number | Page offset |

Translation

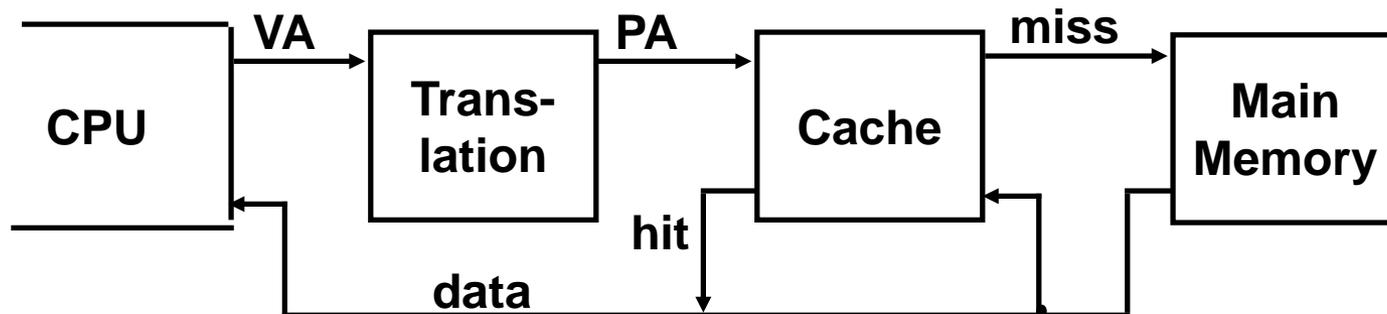| Physical page number | Page offset |
|---|---|
| 29 . . . 12 | 11 0 |

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
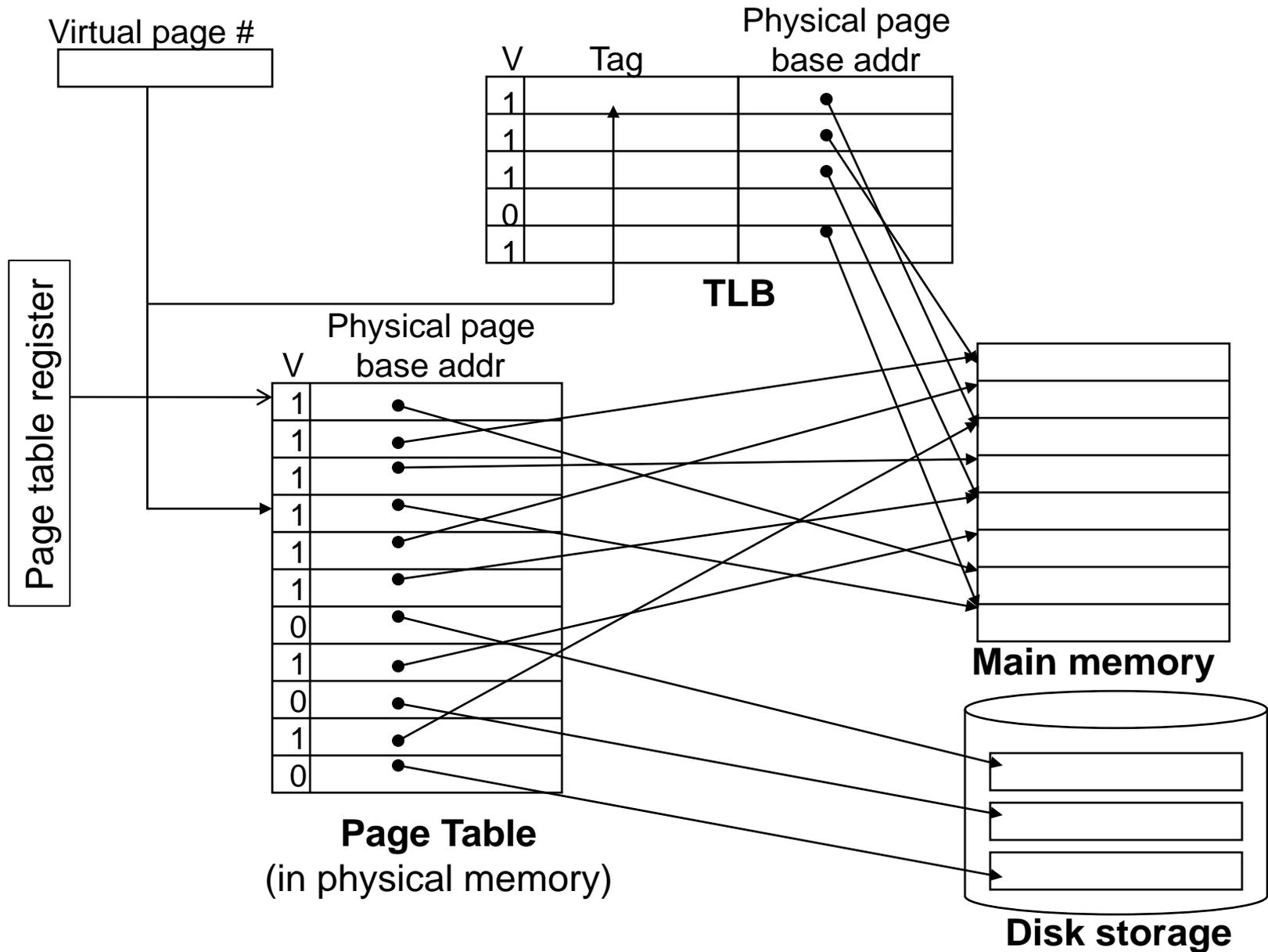
# Address Translation Mechanisms



**Page Table**
(in main memory)

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    –

**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
CPU  ──────►  Trans-  ──────►  Cache  ──────►  Main
             lation                            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

V Tag | Physical page base addr

**TLB**

V | Physical page base addr

**Main memory**

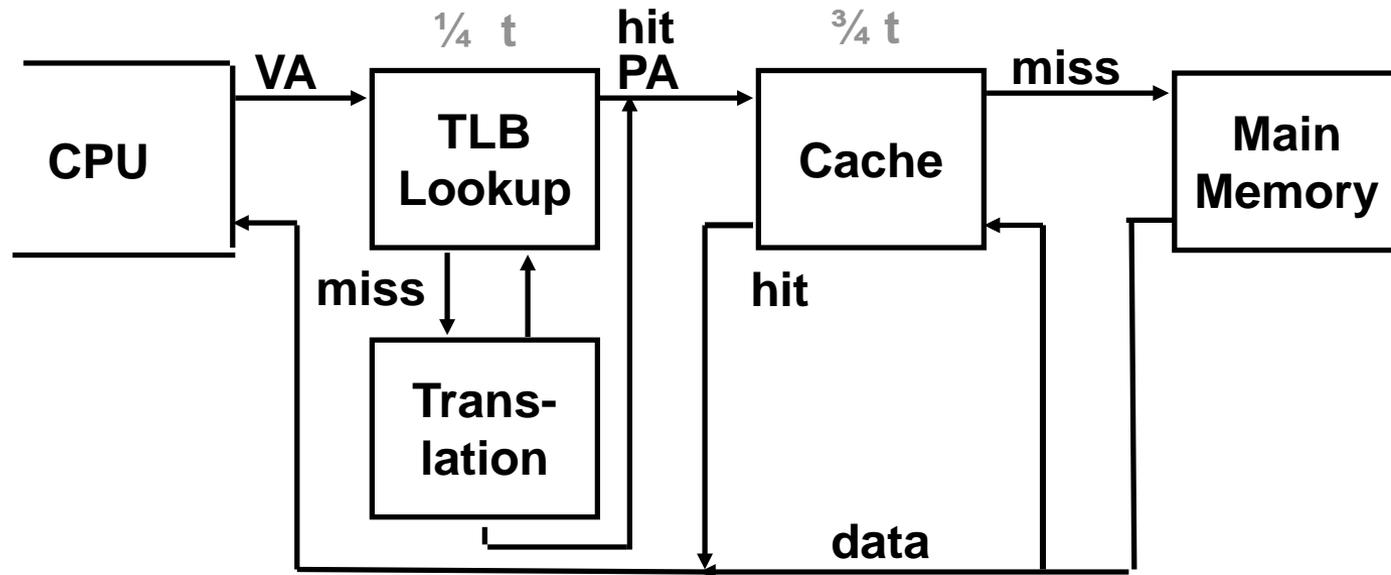**Page Table**
(in physical memory)

**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)

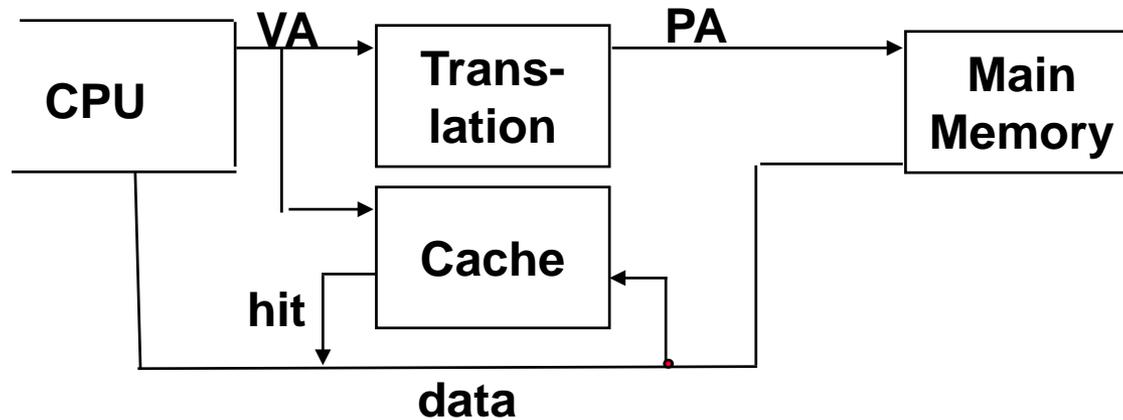  ❑ TLBs are typically not more than 512 entries even on high end machines

# 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

# Why Not a Virtually Addressed Cache?

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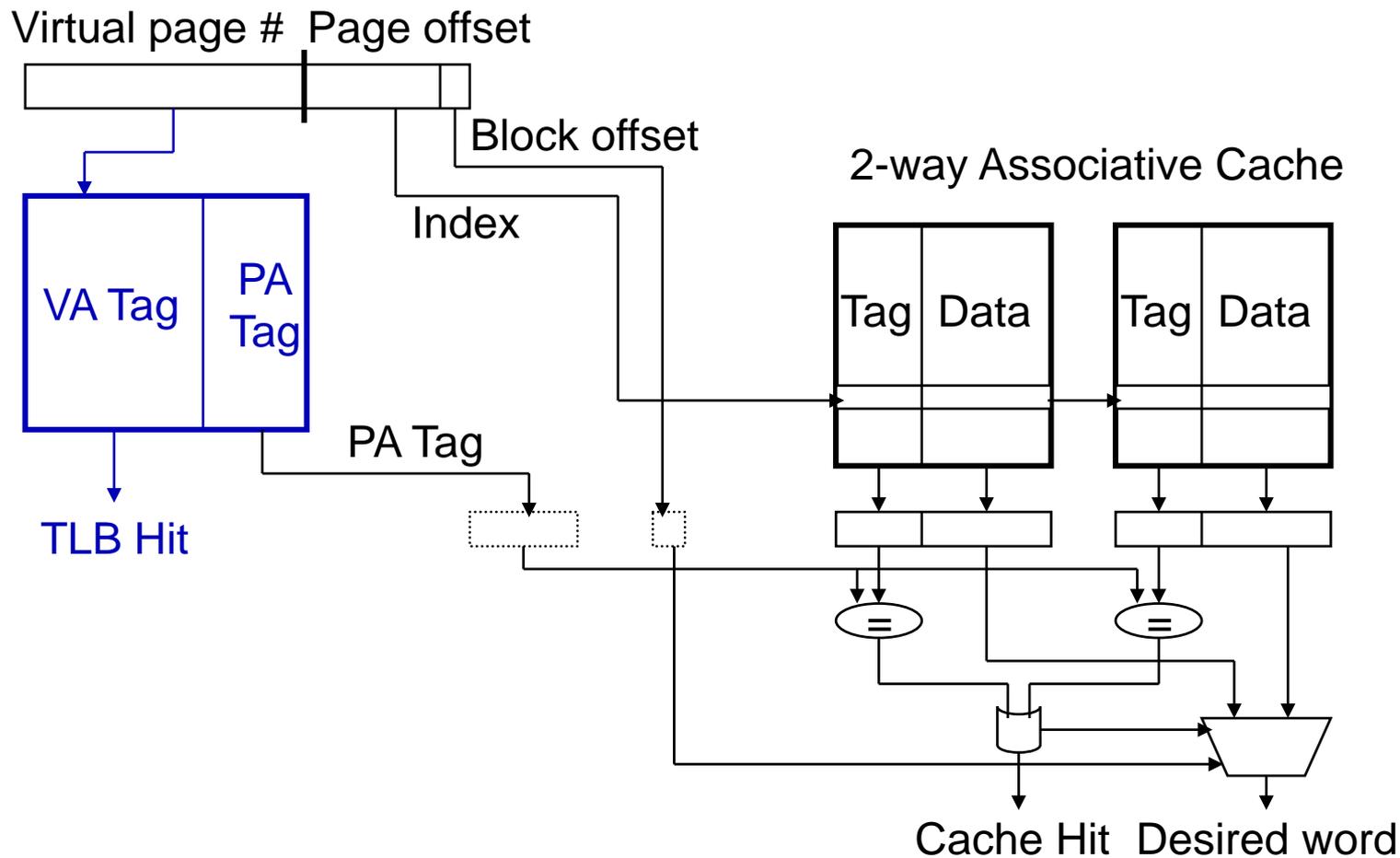


but

◻ 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

2-way Associative Cache

Index

VA Tag | PA Tag

Tag | Data     Tag | Data

PA Tag

TLB Hit

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

# 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

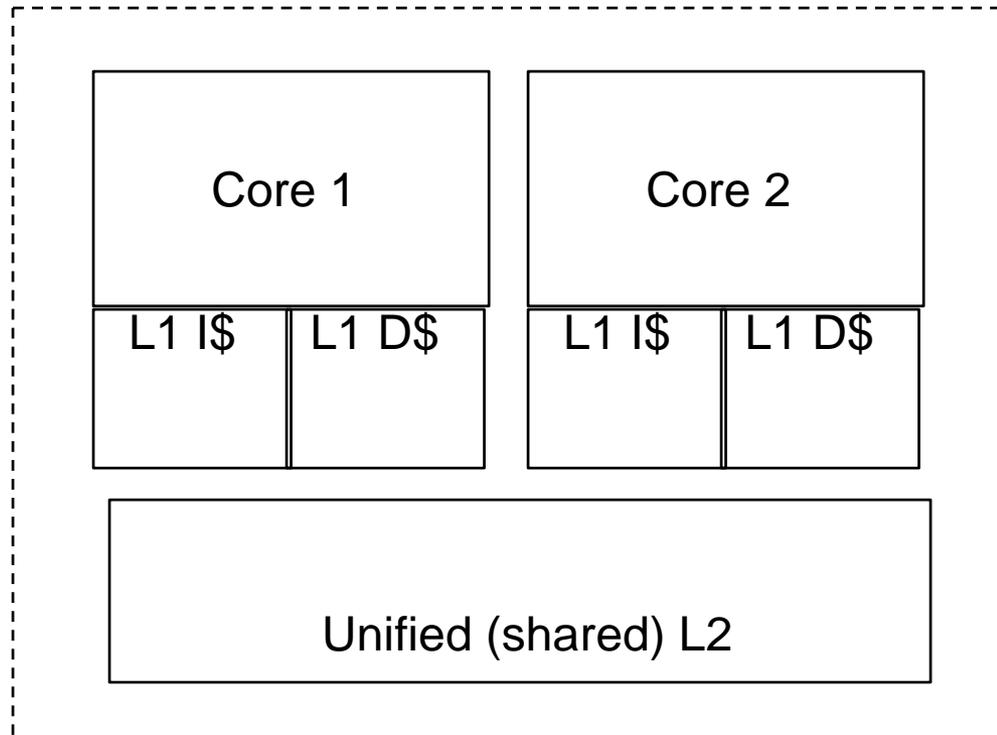| TLB | Page Table | Cache | Possible?  Under what circumstances? |
|---|---|---|---|
| Hit | Hit | Hit | Yes – what we want! |
| Hit | Hit | Miss | Yes – although the page table is not checked if the TLB hits |
| Miss | Hit | Hit | Yes – TLB miss, PA in page table |
| Miss | Hit | Miss | Yes – TLB miss, 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 |

# 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

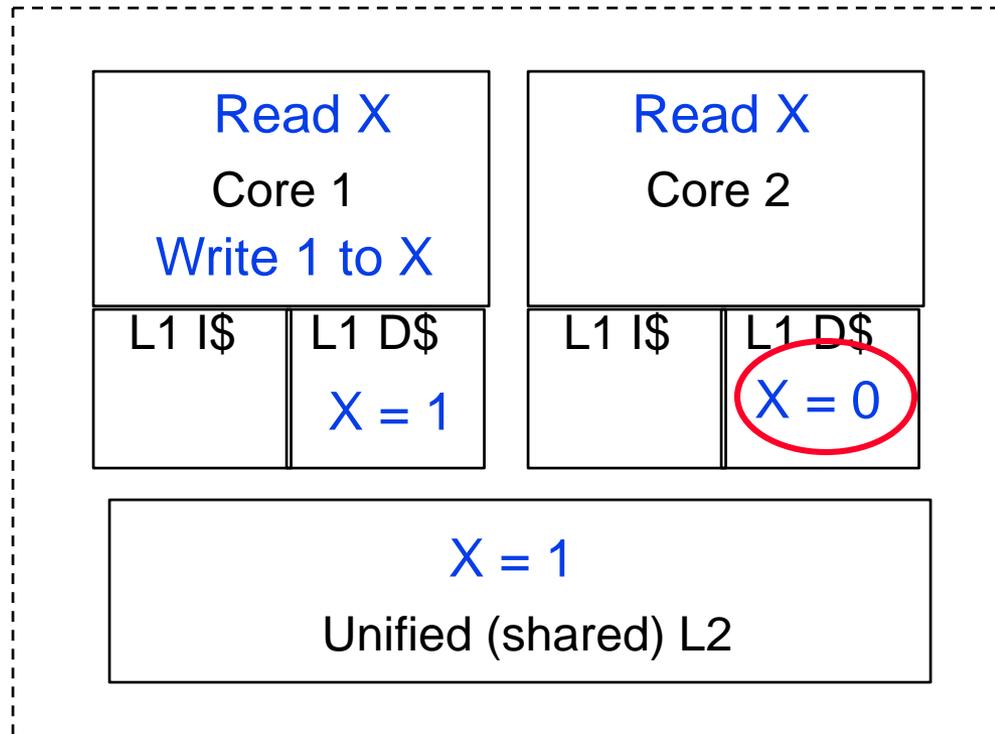|  | **Intel Nehalem** | **AMD Barcelona** |
|---|---|---|
| Address sizes | 48 bits (vir); 44 bits (phy) | 48 bits (vir); 48 bits (phy) |
| Page size | 4KB | 4KB |
| TLB organization | L1 TLB for instructions and L1 TLB for data per core; both are 4-way set assoc.; LRU<br><br>L1 ITLB has 128 entries,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

# Cache Coherence in Multicores

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

# 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)
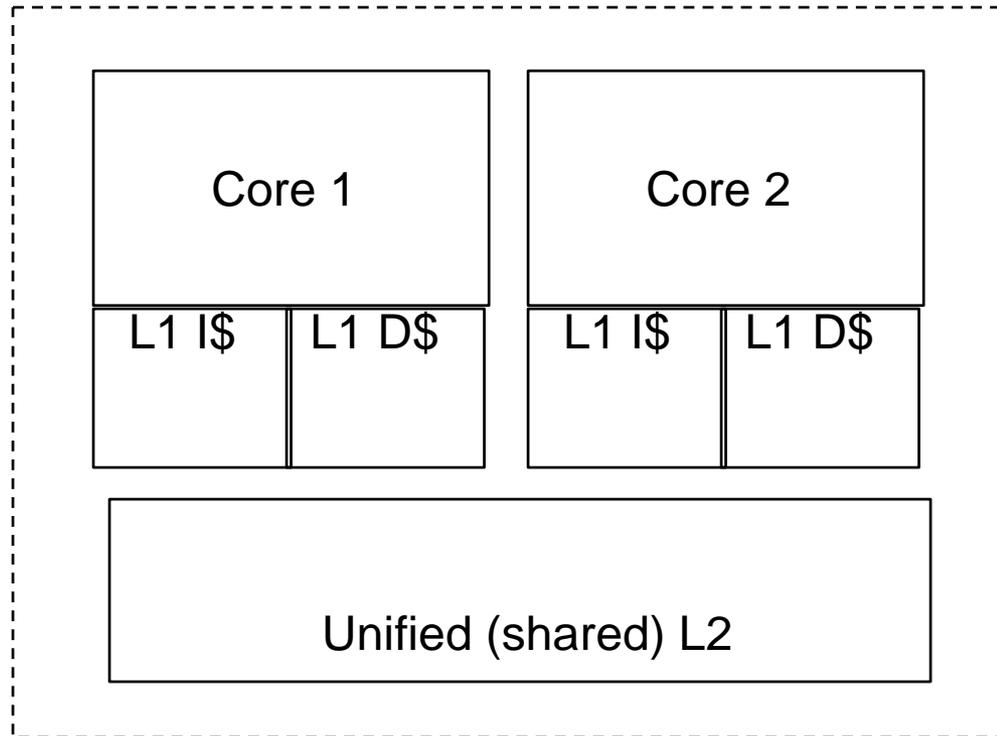
# Cache Coherence Protocols

- ❑ 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
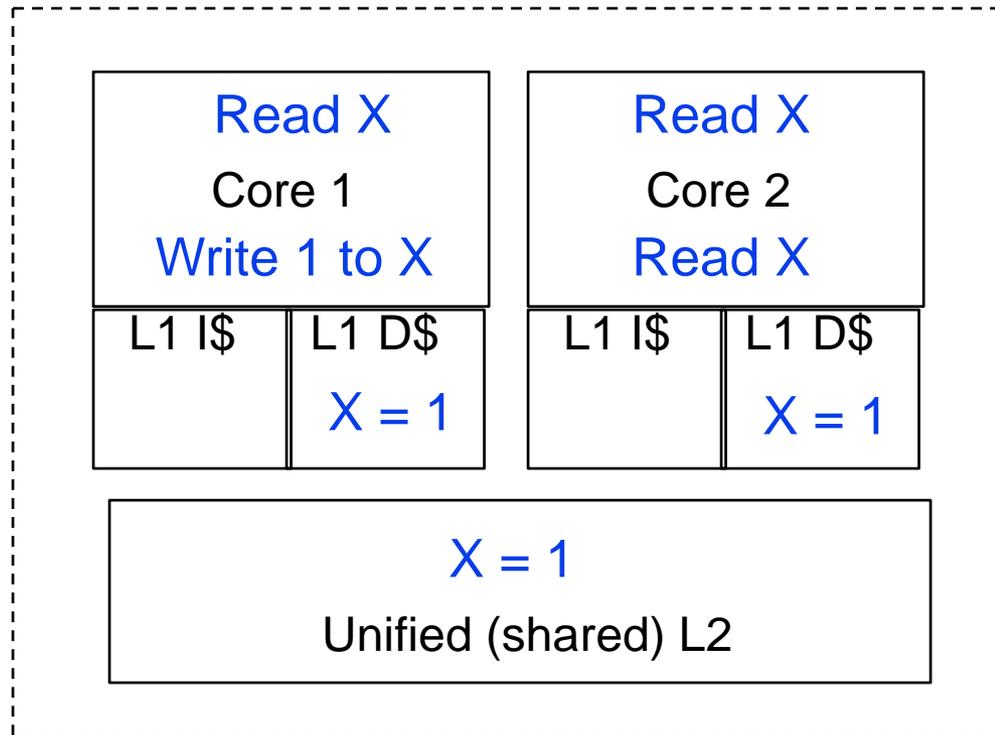
# Handling Writes

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)

# 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
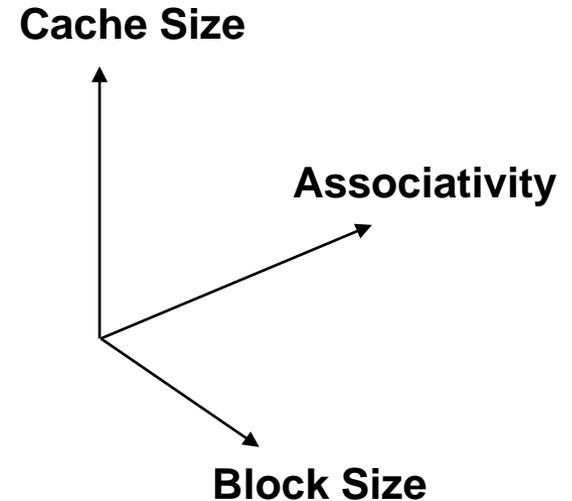
# Summary:  Improving Cache Performance

## 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

❑ Simplicity often wins

**Bad**

**Good**  Factor A          Factor B

**Less**                **More**