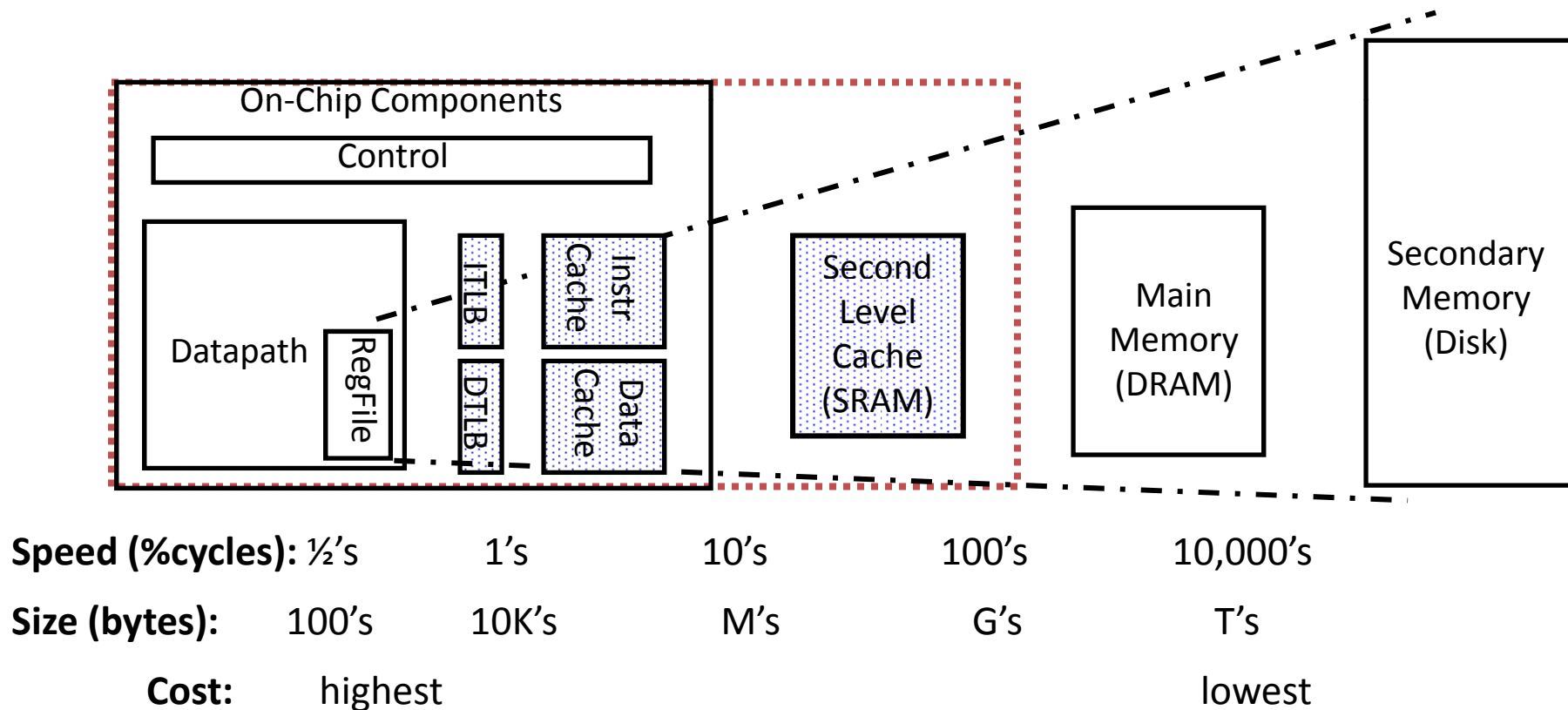# Lecture 14

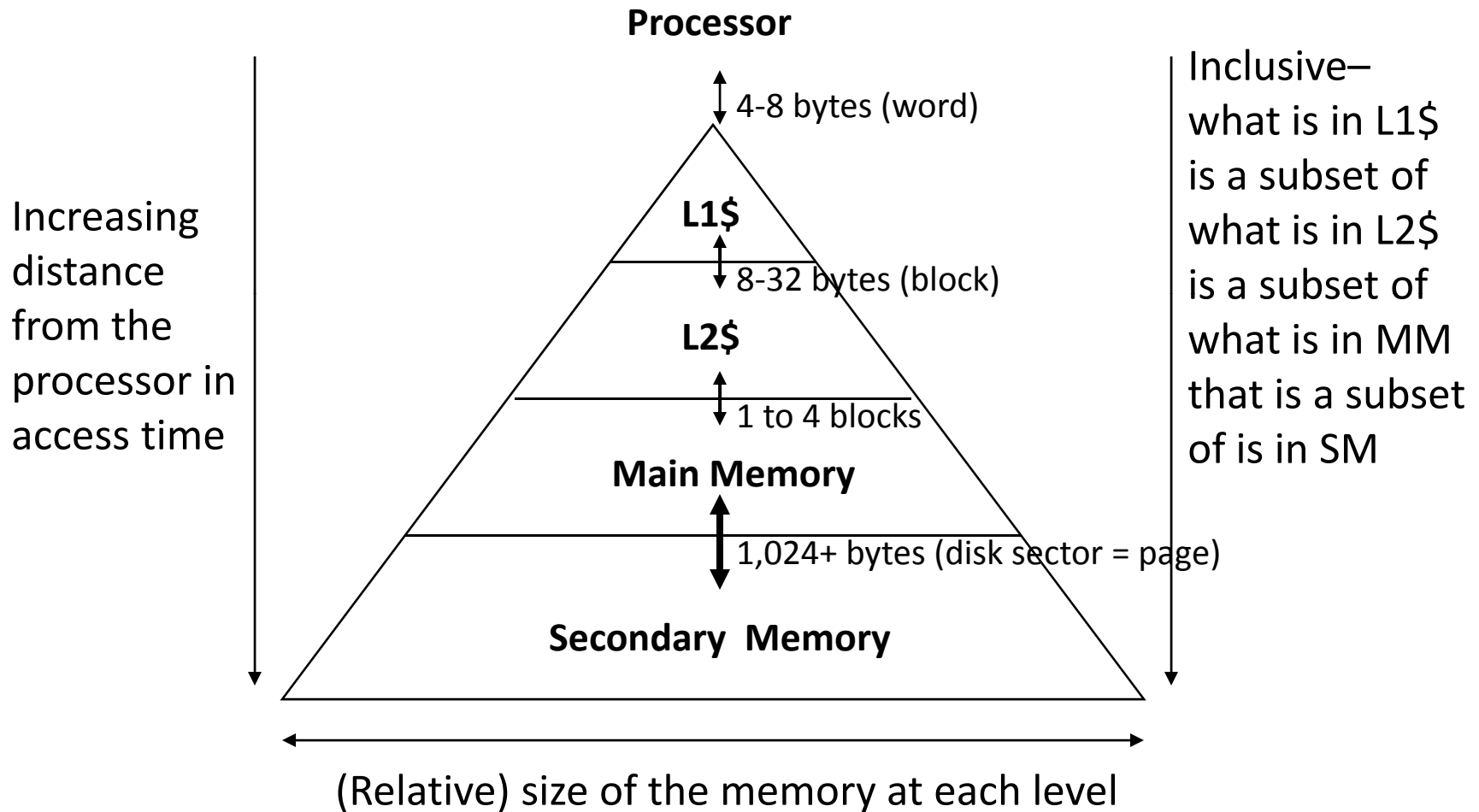# Direct-mapped cache

# 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

| | | | | | |
|---|---|---|---|---|---|
| **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 | |

# Characteristics of the Memory Hierarchy

**Processor**

4-8 bytes (word)

**L1$**

8-32 bytes (block)

**L2$**

1 to 4 blocks

**Main Memory**

1,024+ bytes (disk sector = page)

**Secondary Memory**

Increasing distance from the processor in access time

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

# 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)
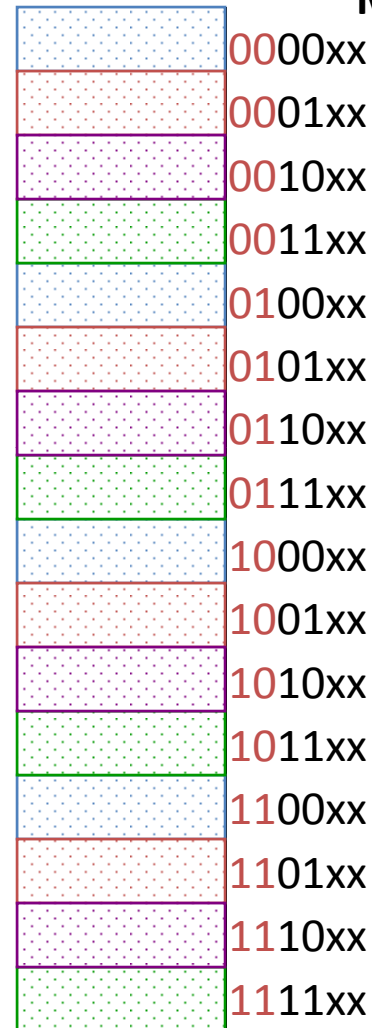
# Caching:  A Simple First Example

**Main Memory**

**Cache**

Index  Valid   Tag    Data

| 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
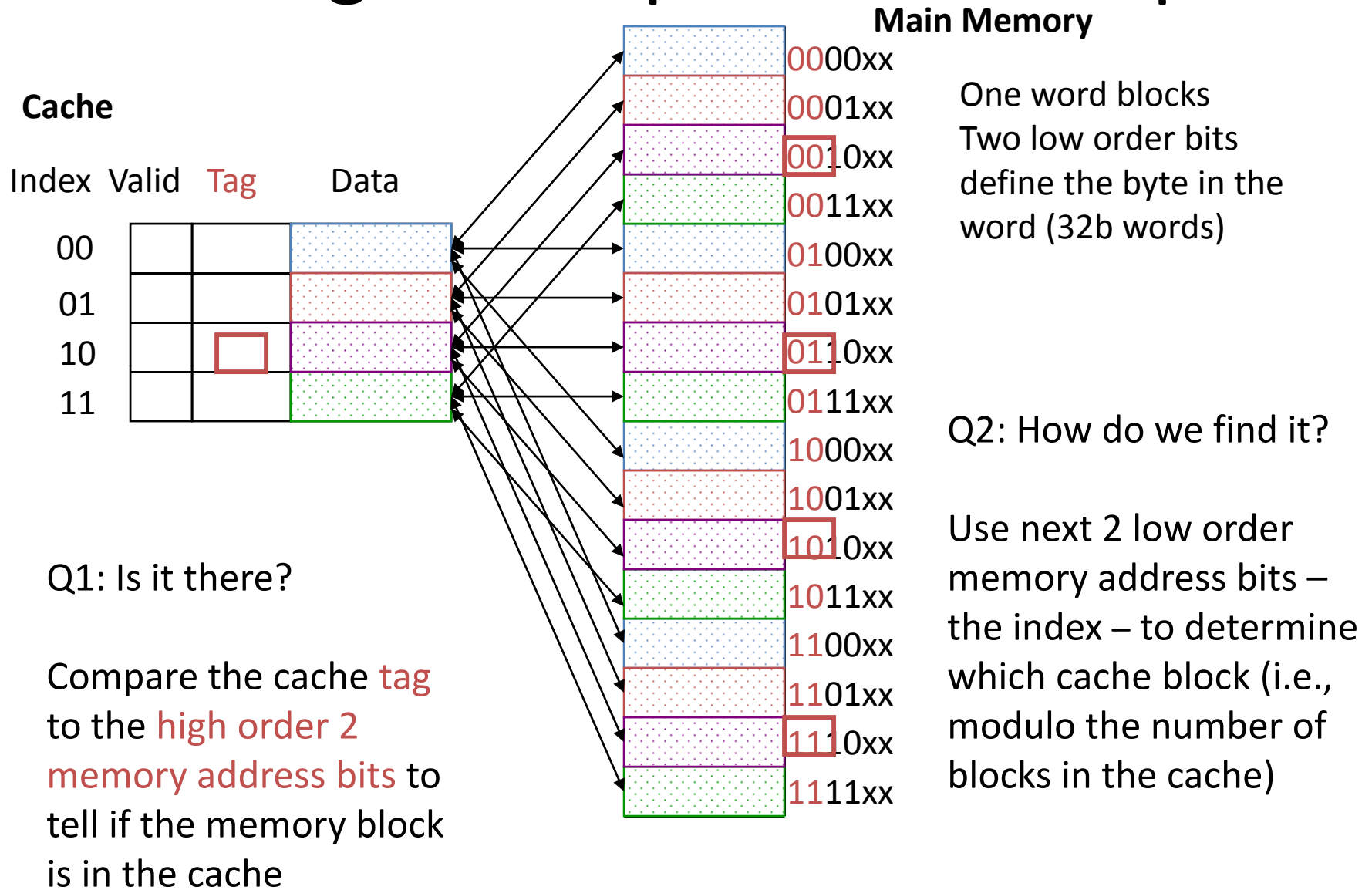1100xx
1101xx
1110xx
1111xx

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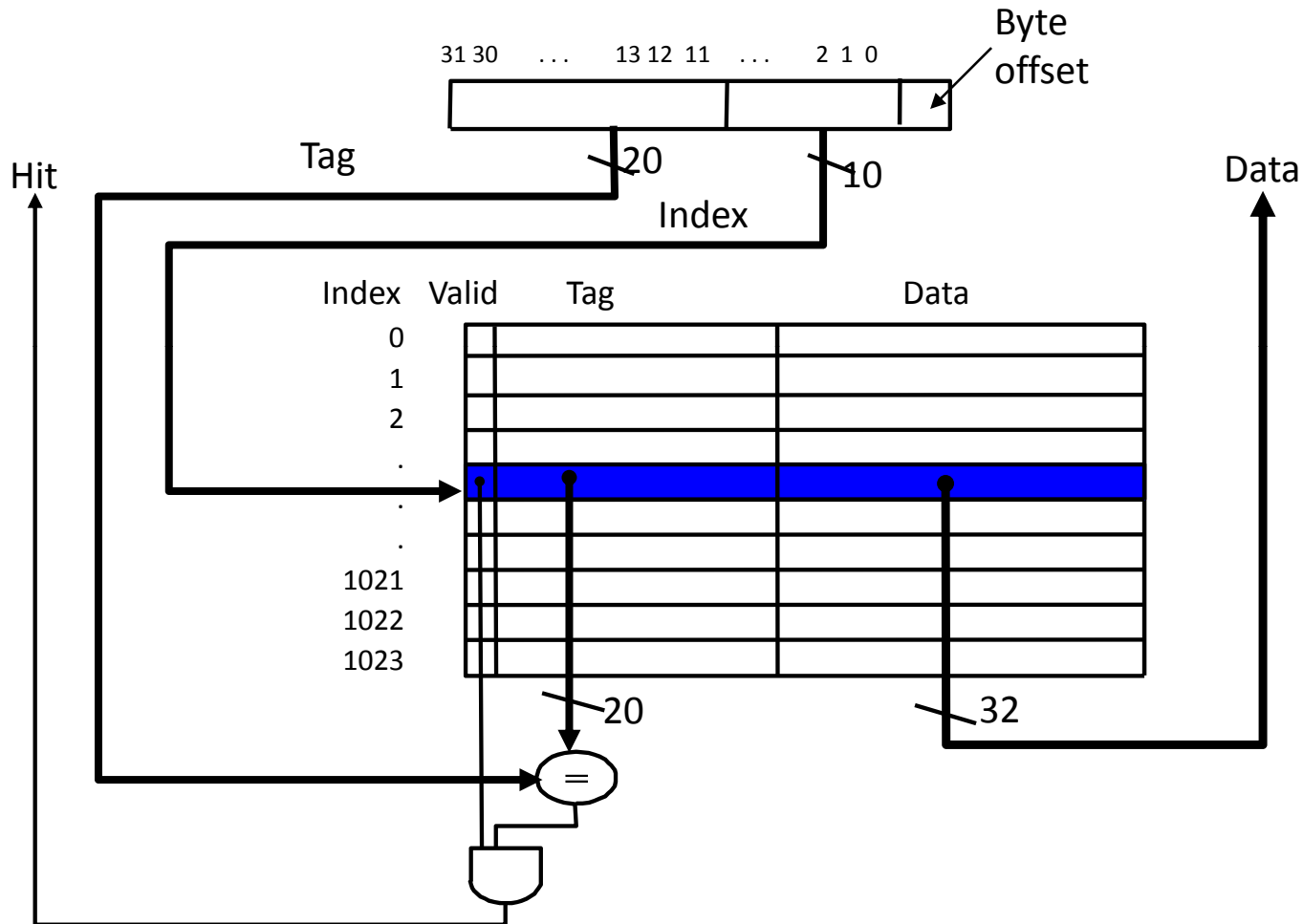
(block address) modulo (# of blocks in the cache)

# MIPS Direct Mapped Cache Example

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

31 30 . . . 13 12 11 . . . 2 1 0

Byte offset

Tag  20

Index  10

Hit

Data

| Index | Valid | Tag | Data |
|-------|-------|-----|------|
| 0 | | | |
| 1 | | | |
| 2 | | | |
| . | | | |
| . | | | |
| . | | | |
| 1021 | | | |
| 1022 | | | |
| 1023 | | | |

20

32

=

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

# 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

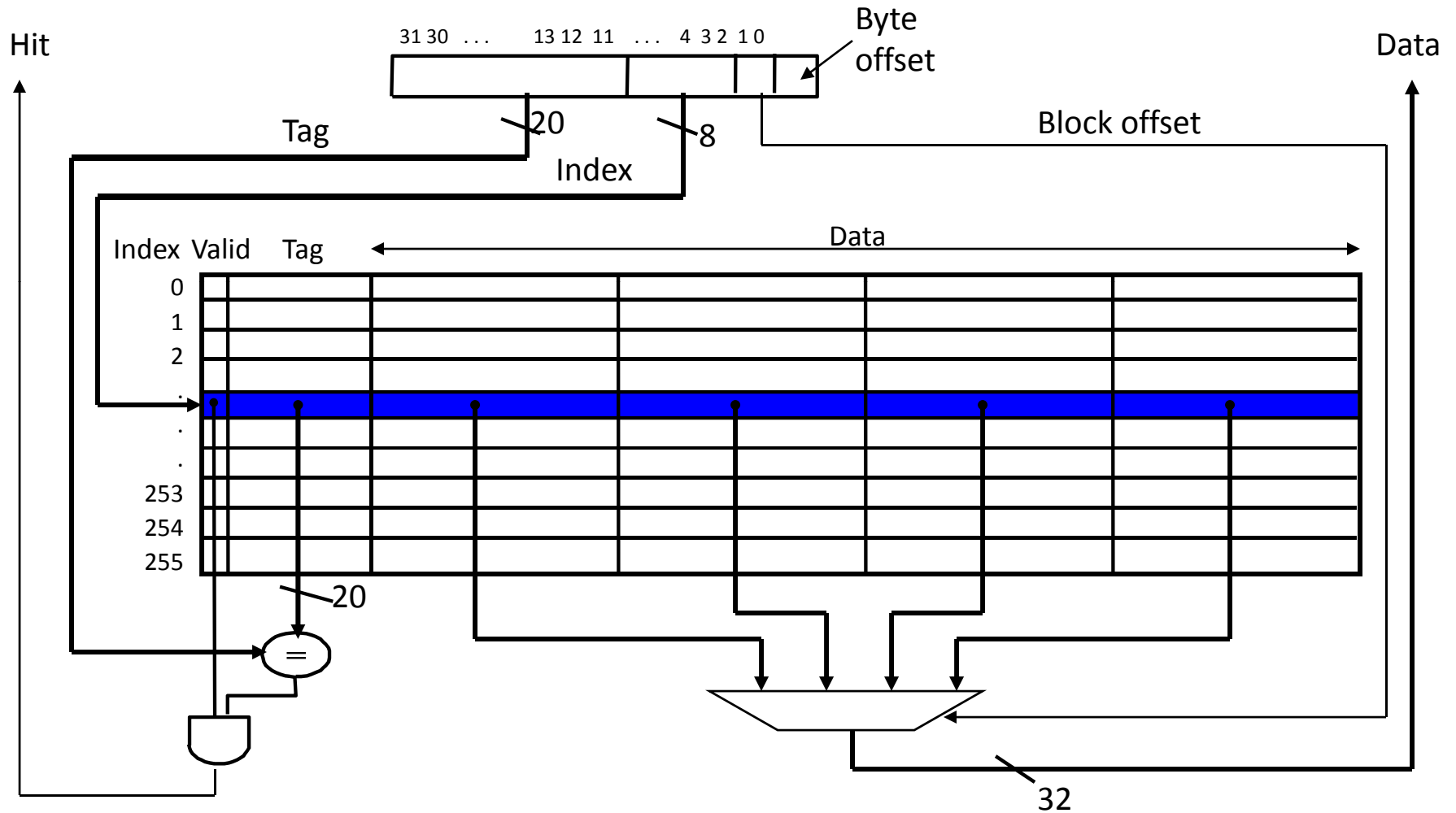| 01 | Mem(4) |
|----|--------|
| 00 | Mem(1) |
| 00 | Mem(2) |
| 00 | Mem(3) |

11                          15

- 8 requests, 6 misses

# Multiword Block Direct Mapped Cache
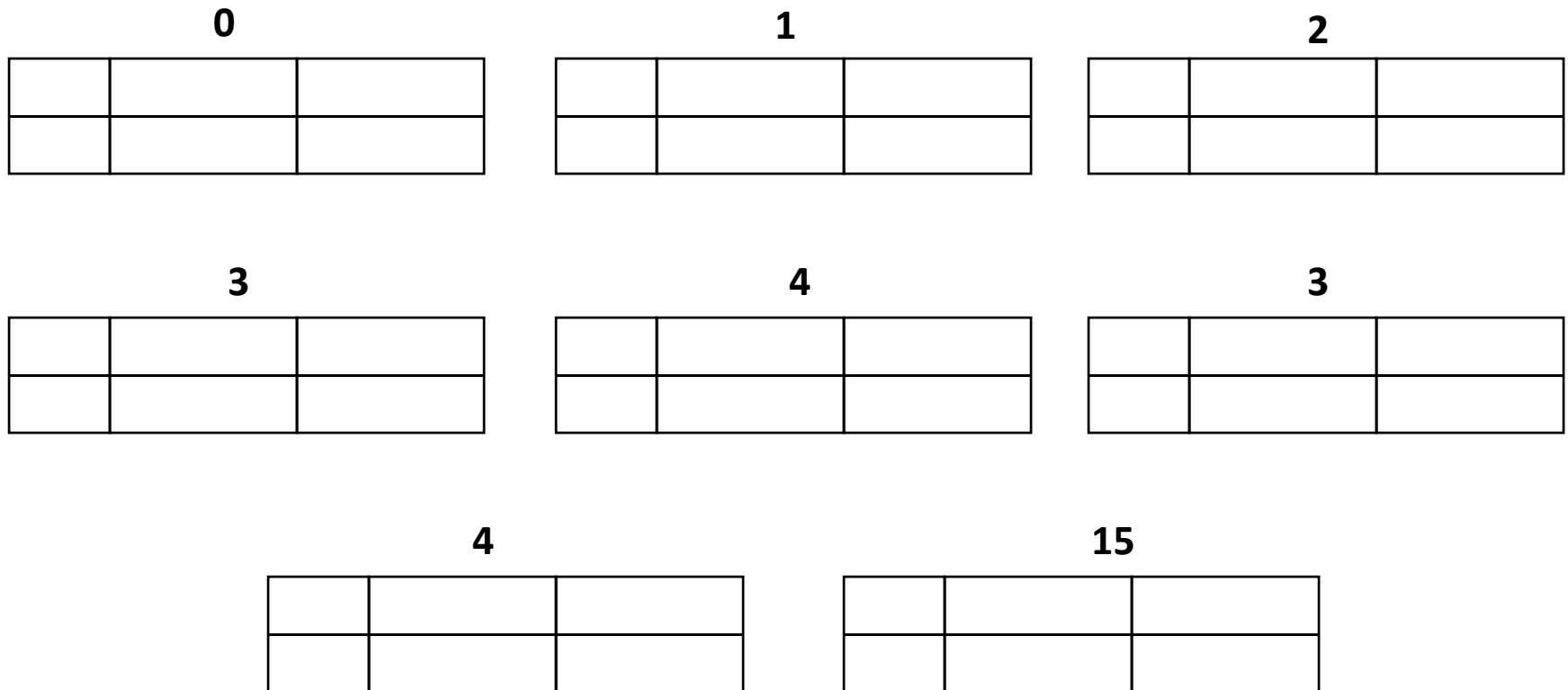
- Four words/block, cache size = 1K words



*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

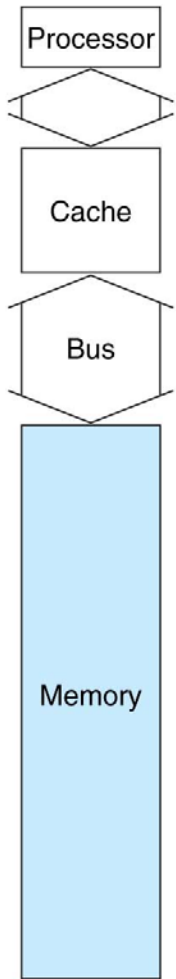11      15      14

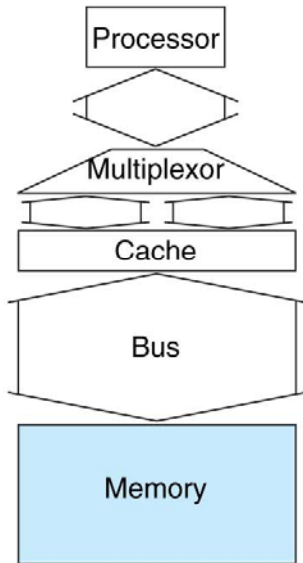| 01 | Mem(5) | Mem(4) |
|----|--------|--------|
| 00 | Mem(3) | Mem(2) |

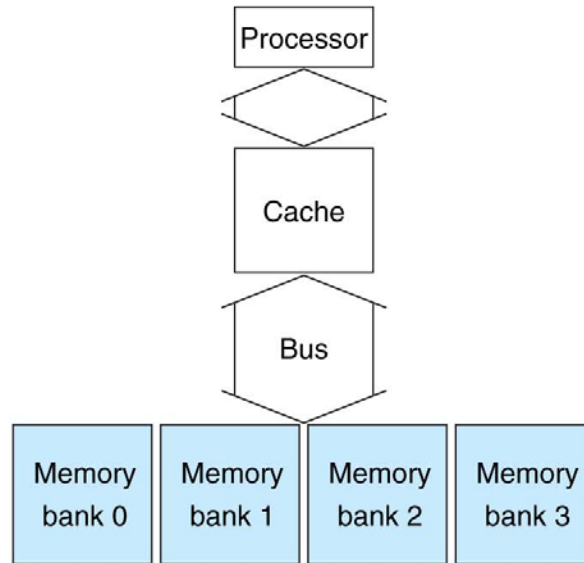- 8 requests, 4 misses

# Miss Rate vs Block Size vs Cache Size



❑ 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)

a. One-word-wide
   memory organization

b. Wider memory organization

c. Interleaved memory organization

# 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)