

THE MEMORY SYSTEM

INTRODUCTION

- ☞ Programs and the data they operate on are held in the memory of the computer.
- ☞ The execution speed of programs is highly dependent on the speed with which instructions and data can be transferred between the processor and the memory.
- ☞ It is also important to have a large memory to facilitate the execution of programs that are large and deal with huge amounts of data.
- ☞ Ideally, the memory would be fast, large and inexpensive. Unfortunately, it is impossible to meet all three of these requirements simultaneously. Increase in speed and size are achieved at increased cost.
- ☞ Much work has gone into developing clever structures that improve the apparent speed and size of the memory, yet keep the cost reasonable.

$$2^n = N$$

SOME BASIC CONCEPTS

$$2^{16} = 65,536$$

☞ The addressing scheme of any computer system determines the maximum size of the main memory that a computer can use.

$$n = \frac{\log \text{MMsize}}{\log_2}$$

Examples:

$$= \frac{\log 65,536}{\log_2} \quad \text{If number of address bits} = 16$$

$$\text{MM Size} = 2^{16} = 65,536$$

= 64K memory locations

2. If number of address bits = 32

$$2^{32} = 4,294,967,296 \quad \text{MM Size} = 2^{32} = 42,94,967,296$$

$$n = \frac{\log 4,294,967,296}{\log_2} \quad = 4G \text{ memory locations}$$

3. If number of address bits = 40

$$= \frac{\log 1,099,511,627,776}{\log_2} \quad \text{MM Size} = 2^{40} = 1,099,511,627,776$$

= 1T memory locations

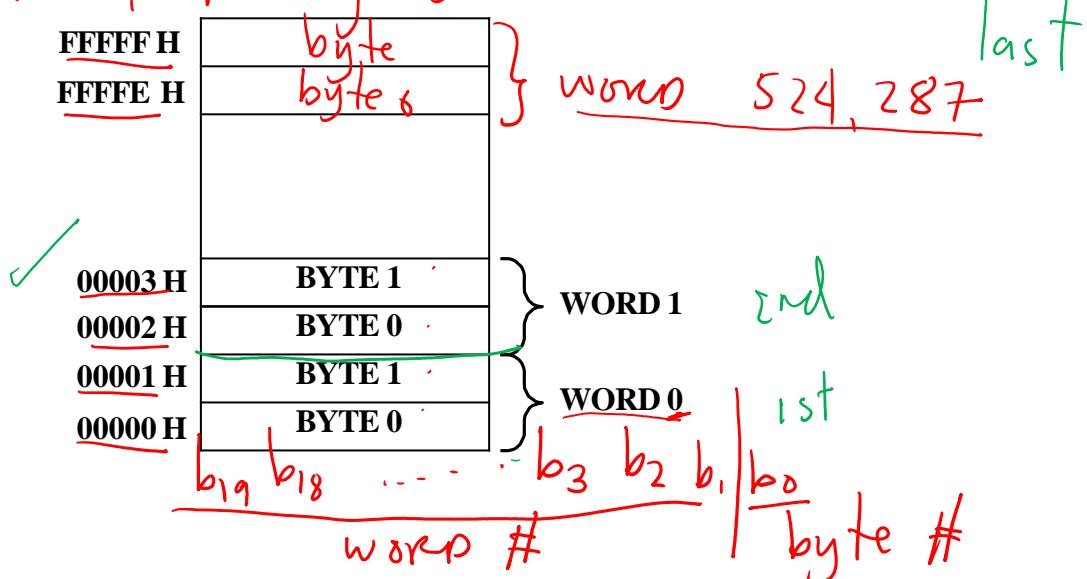
☞ The MM can usually store and retrieve data in word length quantities. However, the processor can access individual memory bytes as long as each byte has a unique address.

Example:

1. For a 16-bit processor

$$\text{Mem size} = 1 \text{ MB} = 1,048,576 \text{ bytes}$$

$$n = \underline{20 \text{ bits}}$$



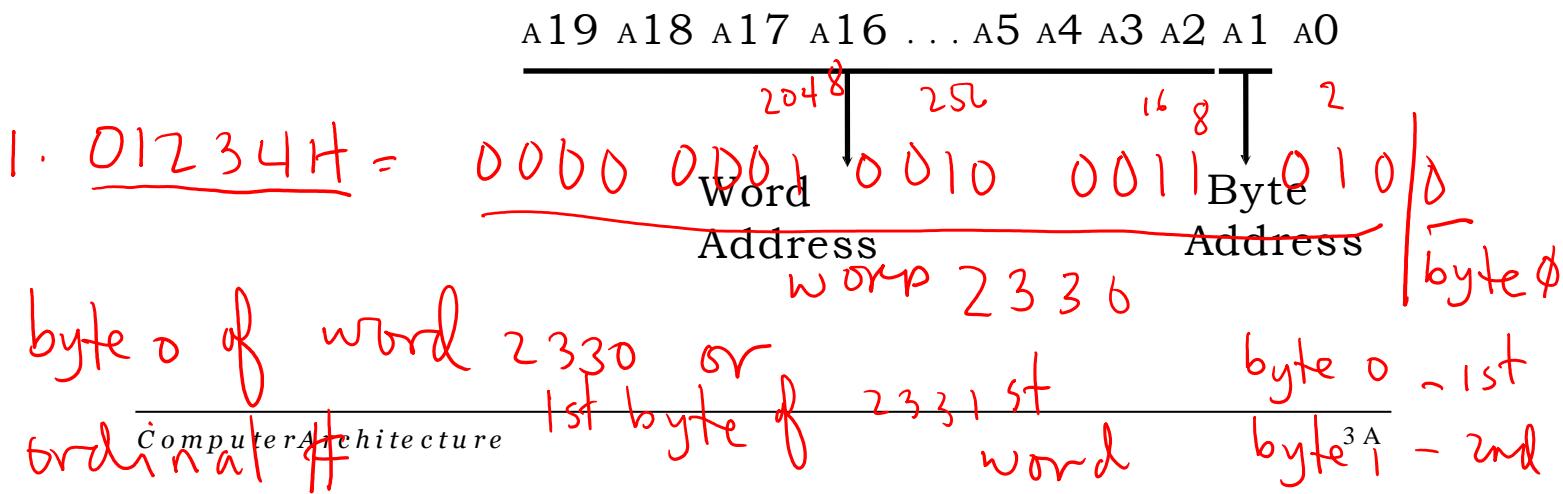
For the first 16-bit word (Word 0):

$$\begin{array}{ll} \text{Low Order Byte Address} & = 0000\dots0000 \\ \text{High Order Byte Address} & = 0000\dots0001 \end{array}$$

For the second 16-bit word (Word 1):

$$\begin{array}{ll} \text{Low Order Byte Address} & = 0000\dots0010 \\ \text{High Order Byte Address} & = 0000\dots0011 \end{array}$$

Therefore:



$$3. \quad \underline{OB0B0H} = \underline{\quad 0000 \quad 1011 \quad 0000 \quad 1011 \quad 000 \quad | 0}$$

word 22,616

byte 0

byte 0 of word 22,616 /

1st byte of 22,617 th word

- 0 - 1st - byte 0
- 1 - 2nd - byte 1

$$2. \quad 00003H = \underline{\quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 001 \quad | 1}$$

word 1

byte 1

byte 1 of word 1 or

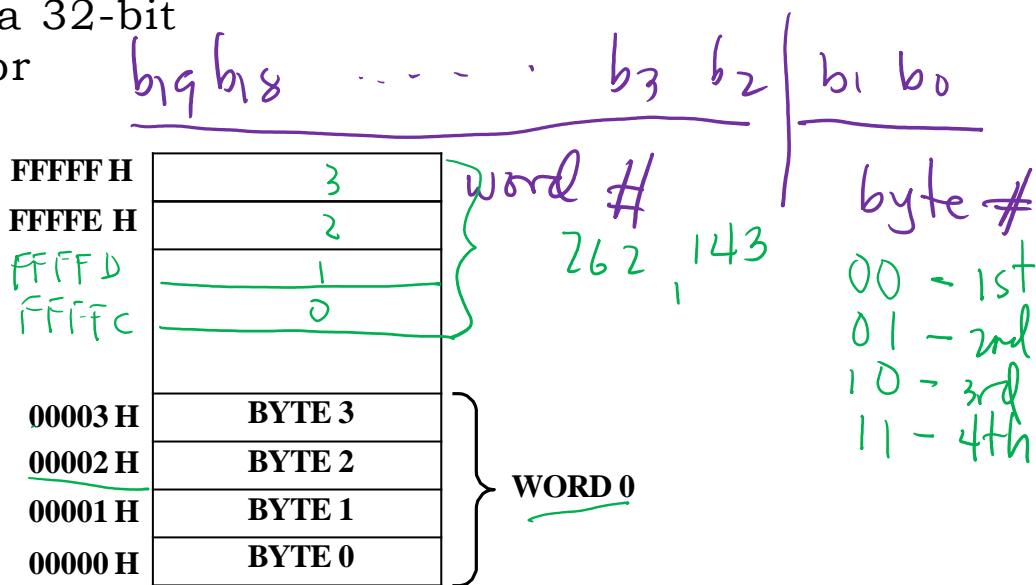
2nd byte of the 2nd word

for 32-bit CPU



2. For a 32-bit processor

$n = \underline{20 \text{ bits}}$



For the first 32-bit word (Word 0):

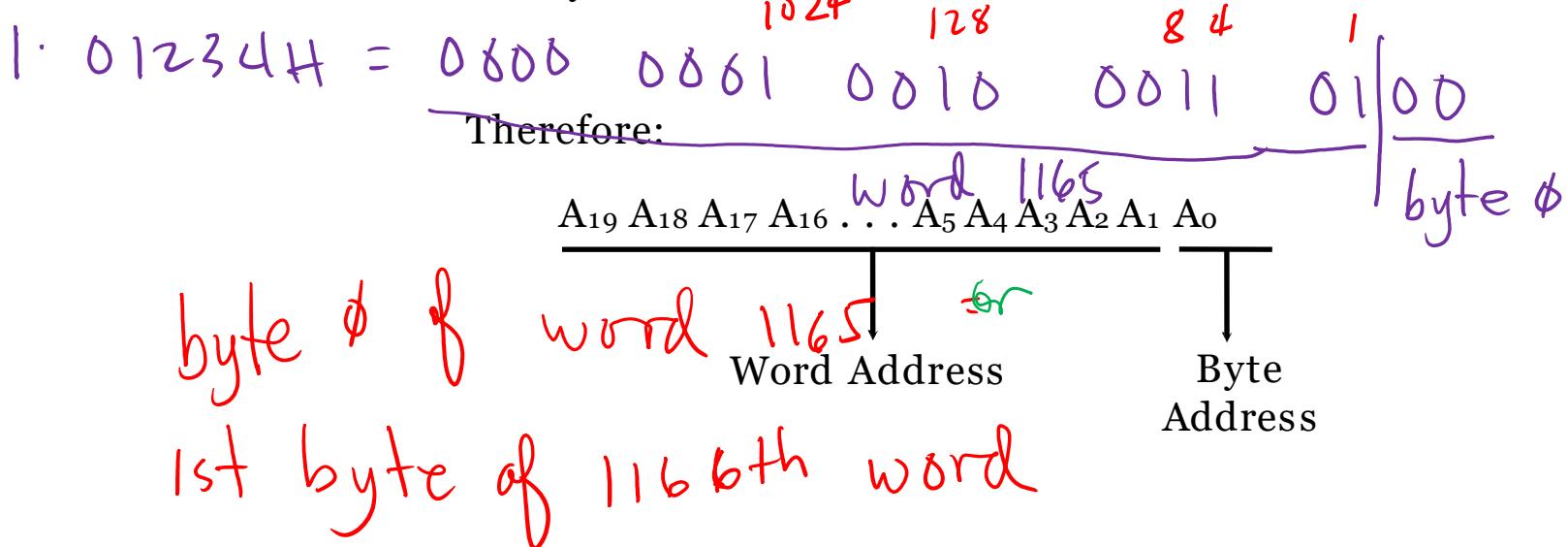
$$1^{\text{s}} \text{ t} \text{ Byte Address} = 0000\ldots0000$$

$$2^{\text{n}} \text{ Byte Address} = 0000\ldots0001$$

$$3^{\text{r}} \text{ d} \text{ Byte Address} = 0000\ldots0010$$

$$4^{\text{t}} \text{ h} \text{ Byte Address} = 0000\ldots0011$$

Example:



$$2. \quad 034A1H = \frac{0000 \ 0011 \ 1024 \ 0100 \ 1010 \ 0001}{\text{Word } 3,368} \quad \text{byte 1}$$

byte 1 of word 3,368 =

2nd byte of the 3,369th word

$$\begin{array}{ccccccccccccc} & 16 & 8 & 4 & 2 & 1 & & & 3 & 2 & 1 & 0 \\ 1024 & 2 & 2 & 2 & 2 & 2 & 1 & 5 & 2 & 2 & 2 & 2 \\ & 512 & 256 & 128 & 64 & 32 & 16 & & 8 & 4 & 2 & 1 \end{array}$$

binary places

byte # (2 bit)

$$\begin{array}{lll} \text{byte 0 - 1st} & = & 00 \\ \text{1 - 2nd} & = & 01 \\ \text{2 - 3rd} & = & 10 \\ \text{3 - 4th / last} & = & 11 \end{array}$$

$$3. \quad 00002H = \frac{0000 \ 0000 \ 0000 \ 0000 \ 00}{\text{Word 0}} \quad \text{byte 2}$$

byte 2 of word 0 =

3rd byte of the 1st word

32-bit CPU

= last byte of word 380

n = 20 bit = 005F3H

0000	0000	0000	0000	0000	0000
0000	0800	0101	1111	0011	
0	0	5	F	3	

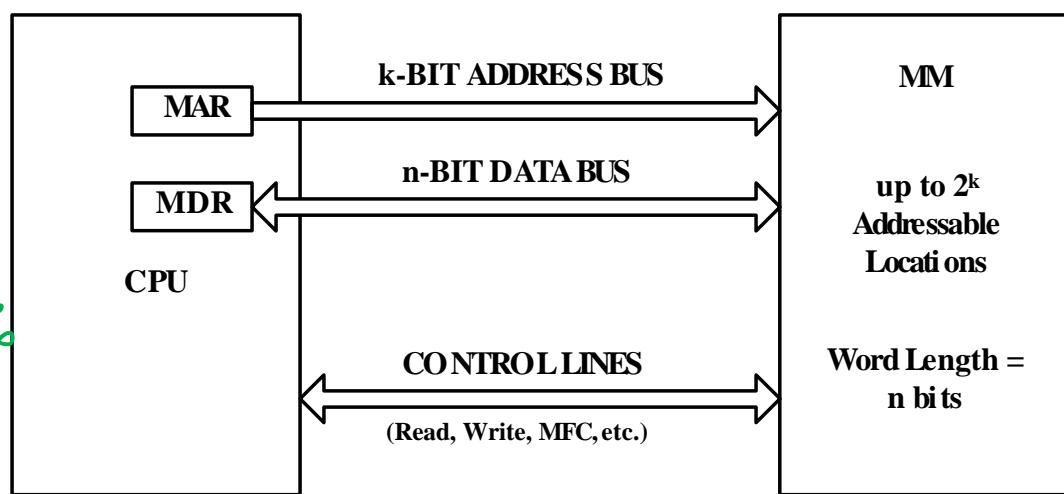
1. 02468H

16-bit

64-bit

b₁₉ b₁₈ ... b₃ | b₂ b₁ b₀

- Connection of the main memory to the processor:



MM size =

1 MB =

1,048,576

address bus - 20 bit

8086

8088

20 bits

- Memory Access Time is the time between the initiation of an operation and the completion of that operation.

data bus
(external)

16-bit

8 bit

Memory Cycle Time is the minimum time delay required between the initiation of two successive memory operations, for example, the time between two successive memory read operations. Cycle time is usually slightly longer than the access time, depending on the implementation details of the memory unit.

Internal

16-bit

80187
80287

16-bit

- ☞ *Random Access Memory* is memory in which the access time of any location is a fixed amount of time that is independent of the location's address. Main memory units are of this type. This distinguishes them from serial, or partly serial, access storage devices such as magnetic tapes and disks. Access time on these devices depends on the address or position of data.
- The processor can usually process instruction and data faster than they can be fetched from memory. Thus, MM becomes a bottleneck.

fast huge inexpensive

SPEED, SIZE, AND COST

- ☞ Fast memory can be achieved if bipolar static RAM chips are used. But these chips are expensive and have relatively small packing densities. Therefore, it is impractical to build a large memory using these chips.
- ☞ The only alternative is to use dynamic RAM chips, which have higher packing densities and are much cheaper. But such memories are significantly slower.
- ☞ The solution therefore is to implement main memory using dynamic RAM technology and use bipolar static RAMs in building smaller memory units where speed is of the essence, such as in *cache memories*.

CACHE MEMORY

- ☞ The speed of the main memory is very low in comparison with the speed of modern processors. For good performance, the processor cannot spend much of its time waiting to access instructions and data in main memory.

- ☞ Therefore, it is important to devise a scheme that reduces the time needed to access the necessary information. Since the speed of the main memory unit is limited by electronic and packaging constraints, the solution must be sought in a different architectural arrangement.

- ☞ An efficient solution is to use a fast cache memory which essentially makes the main memory appear to the processor to be faster than it really is.

- ☞ The effectiveness of the cache mechanism is based on a property of computer programs called the *locality of reference*. Analysis of programs shows that most of their execution time is spent on routines in which many instructions are executed repeatedly. These instructions may constitute a simple loop, nested loops, or a few procedures or functions that repeatedly call each other.

4084

8088

?

T,

i3 -

LS - ↗

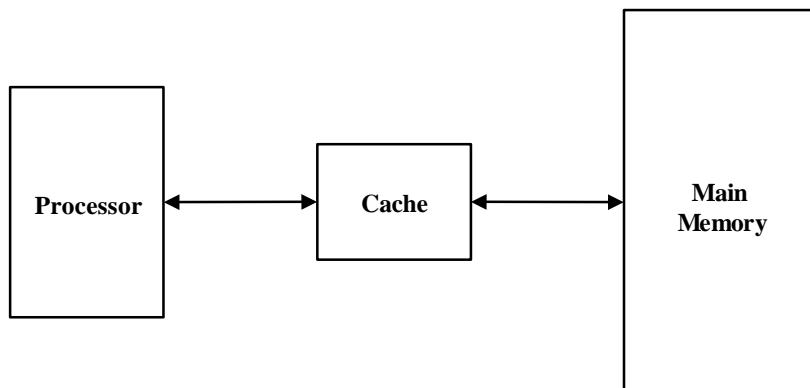
i4 -

i9 -

In other words, locality of reference states that many instructions in a few localized areas of the program are repeatedly executed during some period of time and that the remainder of the program is accessed relatively infrequently.

Locality of reference manifests itself in two ways: *temporal* and *spatial*. Temporal means that a recently executed instruction is likely to be executed again very soon. The spatial aspect means that instructions in close proximity to a recently executed instruction are also likely to be executed soon.

It is therefore logical to put the active segments of a program in a fast memory to reduce total execution time. The fast memory is called the *cache memory*.



☞ Operation of Cache Memory:

The temporal aspect of the locality of reference suggests that whenever an information item is needed (instruction or data) is first needed, this item should be brought into the cache where it will hopefully remain until it is needed again.

The spatial aspect suggests that instead of fetching just one item from the main memory to the cache, it is useful to fetch several items that reside in adjacent addresses as well. The term *block* is used to refer to a set of contiguous address locations of some size. Main memory is therefore divided into blocks of words. Cache memory is smaller than main memory but it can hold or store a number of these blocks at any given time.

When the processor issues a read request, the system transfers the contents of a block of memory words containing the location specified into the cache one word at a time. Subsequently, when the program references any of the locations in this block, the desired contents are read directly from the cache.

When the cache is full and the program references a memory word that is not in the cache, the cache control hardware must decide which block to remove to create space for the new block that contains the referenced word. The collection of rules for making this decision constitutes the *replacement algorithm*.

The processor does not need to know explicitly about the existence of the cache. It simply issues Read or Write requests using addresses that refer to the locations in the memory. The cache control circuitry determines whether the requested word currently exists in the cache.

- ☞ If the requested word is in the cache, the Read or Write operation is performed on the appropriate cache location. In this case, a *read hit* or *write hit* is said to have occurred.
- ☞ In a Read operation, the main memory is not involved. Only the cache is read.
- ☞ In a Write operation, the system can proceed in two ways:
 1. In the first technique, called the *write-through protocol*, the cache location and the main memory location are updated simultaneously.
 2. The second technique is to update only the cache location to mark it as updated (using an associated flag bit called the *dirty* or *modified* bit). The main memory location of the word is updated later, when the block containing this marked word is to be removed from the cache to make room for a new block. This technique is known as the *write-back* or *copy-back protocol*.

- ☞ The write-through protocol is simpler, but it results in unnecessary Write operations in main memory when a given cache word is updated several times during its cache residency.

The write-back protocol may also result in unnecessary Write operations because when a cache block is written back to the memory all words of the block are written back, even if only a single word has been changed while the block was in cache.

- ☞ When the addressed word in a Read operation is not in the cache, a *read miss* occurs. The block of words that contains the requested word is copied from the main memory into the cache. After the entire block is loaded into cache, the particular word requested is forwarded to the processor. Alternatively, this word may be sent to the processor as soon as it is read from main memory. This approach, which is called *load-through*, or *early restart*, reduces the processor's waiting period somewhat, but at the expenses of more complex circuitry.
- ☞ During a Write operation, if the addressed word is not in the cache, a *write miss* occurs. Then, if the write-through protocol is used, the information is written directly into the main memory. In the case of the write-back protocol, the block containing the addressed word is first brought into the cache, and then the desired word in the cache is overwritten with the new information.

Assume: 65,536 bytes

Cache size = 2,048 bytes

block size = 16 bytes / block

$$n = \frac{\log 65,536}{\log 2}$$

$$\# \text{ of } \text{MMBS} = \frac{\text{MM size}}{\text{block size}}$$

$$= \frac{65,536 \text{ bytes}}{16 \text{ bytes}}$$

$$= 4,096 \text{ (or } 4,095 \text{ MMBS)}$$

$$\# \text{ of } \text{COS} = \frac{\text{Cache size}}{\text{block size}}$$

$$= \frac{2,048 \text{ bytes}}{16 \text{ bytes/block}}$$

$$= 128 \text{ (or } 127 \text{ COS)}$$

Assume:

The Memory System

$$\text{MM size} = 65,536 \text{ byte}$$

$$\text{Cache size} = 2048 \text{ byte}$$

$$\underline{\text{block size}} = \underline{16 \text{ bytes/block}}$$

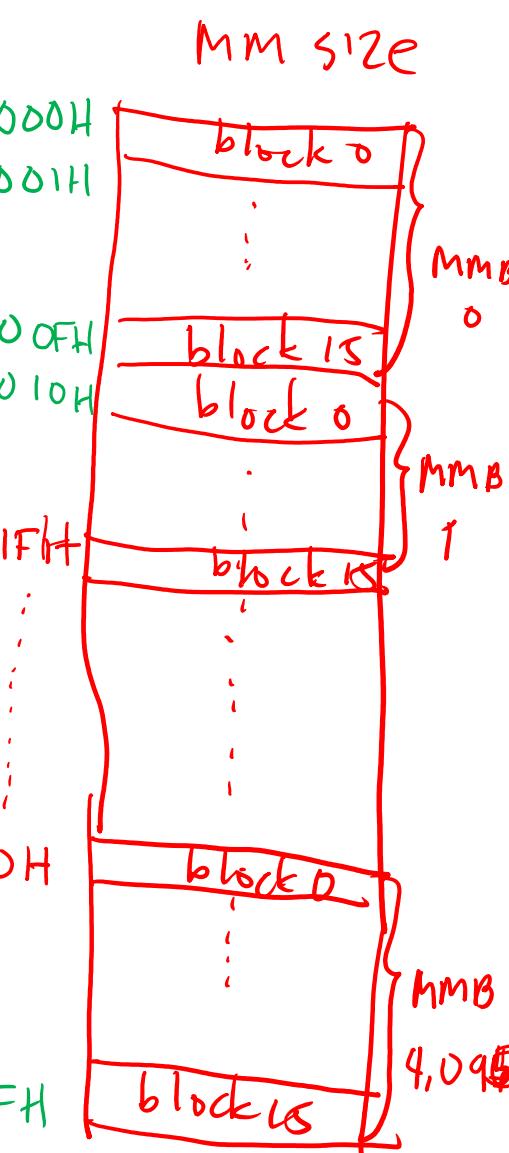
$$n = \frac{\log 65,536}{\log 2} = \boxed{16 \text{ bits}}$$

$$\# \text{ of MMBs} = \frac{\text{MM size}}{\text{block size}}$$

$$= \frac{65,536 \text{ bytes}}{16 \text{ bytes}} \text{ block}$$

$$\# \text{ of MMBS} = \frac{4,096 \text{ blocks}}{(0 - 4,095)}$$

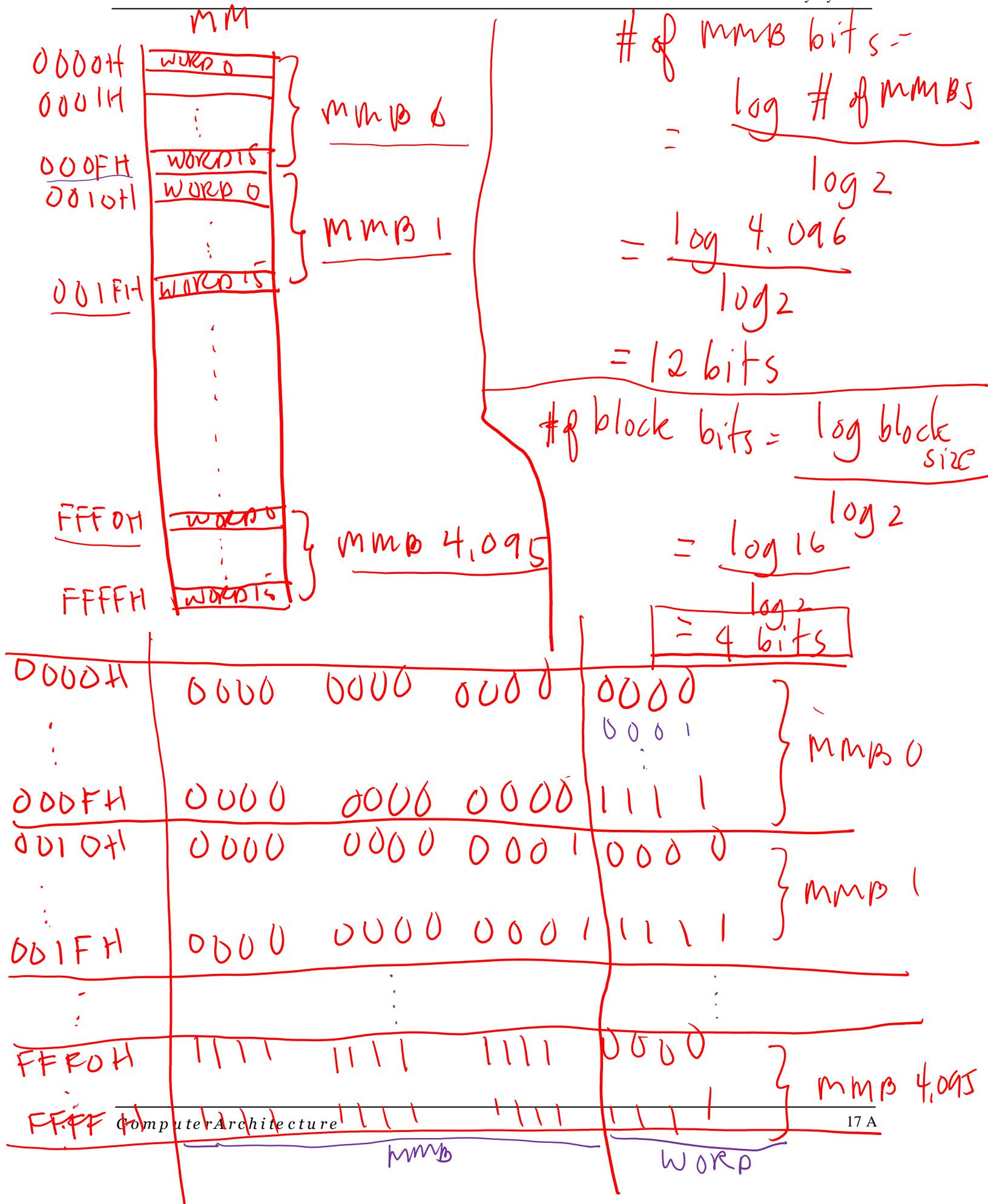
$$\# \text{ of CBs} = \frac{\text{Cache size}}{\text{block size}} =$$



$$\frac{2048 \text{ bytes}}{16 \text{ bytes/block}}$$



$$= \frac{128 \text{ blocks}}{(0 - 127)}$$



$$\# \text{ of MMB bits} = \frac{\log \text{ of Haf MMBs}}{\log_2}$$
$$= \frac{\log 4096}{\log 2}$$

$\boxed{= 12 \text{ bits}}$

$$\# \text{ of word bits} = \frac{\log \text{ block size}}{\log_2}$$
$$= \frac{\log 16}{\log 2}$$

$\boxed{= 4 \text{ bits}}$

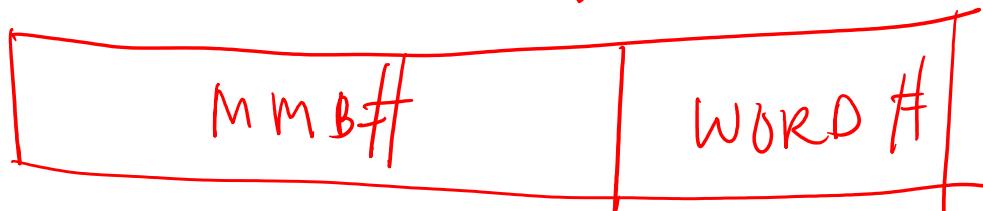
EXAMPLE: mmB # WORD #

1. 001FH = 0000 0000 0001 | 1111
 MMB 1 WORD 15
 word 15 of mmB 1

2. 000FH = 0000 0000 0000 | 1111
 MMB 8 WORD 15
 word 15 of mmB 8

3. 1234H = 0001 0010 0011 | 0100
 MMB 291 WORD 4
 word 4 of mmB 291

n =



32-bit CPU

1. 034164

3. 2BEA2H

2. 1ACE1H

4. 3CA13H
 5. 4D1A4H

FFF0H = 1111 1111 1111 | 0000
MMB 4.095 | WORD 0

EXAMPLE:

$$1 \cdot 001FH = \begin{array}{cccc} 0000 & 0000 & 0001 \\ \hline \text{MMB 1} \end{array} \quad \left| \begin{array}{c} \text{1111} \\ \text{WORD 15} \end{array} \right.$$

word 15 of MMB1

$$2 \cdot 000FH = \begin{array}{cccc} 0000 & 0000 & 0000 \\ \hline \text{MMB 0} \end{array} \quad \left| \begin{array}{c} \text{1111} \\ \text{WORD 15} \end{array} \right.$$

$$3 \cdot 1234H = \begin{array}{cccc} 0001 & 0016 & 0011 \\ \hline \text{MMB 291} \end{array} \quad \left| \begin{array}{c} 0100 \\ \text{WORD 4} \end{array} \right.$$

2^{18} 1024 256 4

$$4 \cdot 03416 = \begin{array}{ccccc} 0000 & 0011 & 0100 & 0001 & 0110 \\ \hline \text{word 3,333} \end{array} \quad \left| \begin{array}{c} \text{byte 2} \end{array} \right.$$

byte 2 of word 3,333 or

3rd byte of 3,334th word

$$5 \cdot 1ACE1H = \begin{array}{cccccc} 0001 & 1010 & 1101 & 1110 & 0001 \\ \hline \text{byte 1 of word 27,448} & \text{word 27,448} & \hline \end{array} \quad \left| \begin{array}{c} \text{byte 1} \end{array} \right.$$

CACHE MAPPING FUNCTIONS

- ☞ A *mapping function* specifies where to place a main memory block in the cache.
- ☞ Case Study:

MM Size = 65,536 words
(16 address bits)

Cache Size = 2,048 words

Block Size = 16 words/block

No. of Cache
Blocks = $2,048 / 16$
 $= 128 \text{ blocks}$
(block 0 to block 127)

No. of MM
Blocks = $65,536 / 16$
 $= 4,096 \text{ blocks}$
(block 0 to 4,095)
block

MM size = 65,536 bytes

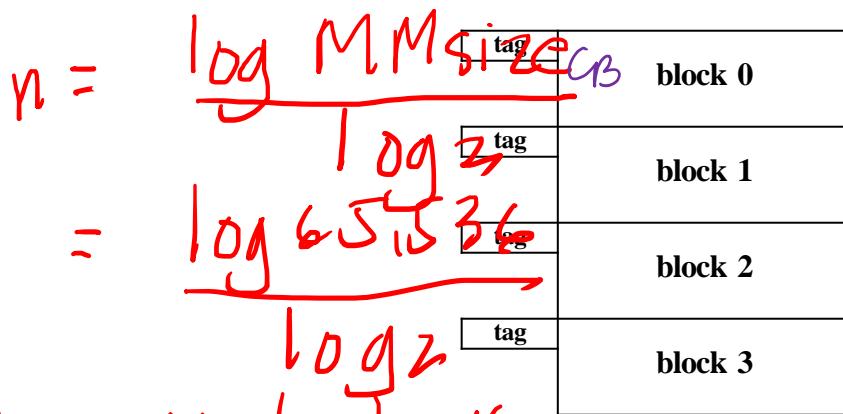
Cache size ~ Mapping Functions
2,048 bytes

Block size 1. Direct Mapping
16 bytes / block

In this technique, block k of the main memory maps onto block k modulo 128 of the cache.

$h = \boxed{\text{TAG} \quad \boxed{\text{BLOCK}} \quad \boxed{\text{WORD}}}$

CACHE



MM blocks 0, 128, 256, ..., 3968

MM blocks 1, 129, 257, ..., 3969

MM blocks 2, 130, 258, ..., 3970

MM blocks 3, 131, 259, ..., 3971

of MMBs = $\frac{\text{MMsize}}{\text{block size}}$

$= \frac{65,536}{16 \text{ bytes}} = 4,096$ (0 - 4,095 MMBs)

of CBs = $\frac{\text{Cache size}}{\text{block size}}$

$= \frac{2,048 \text{ bytes}}{16 \text{ bytes}} = 128$ (0 - 127 CBs)

65,536 bytes
16 bytes / block

MM blocks 125, 253, 381, ..., 4093

MM blocks 126, 254, 382, ..., 4094

MM blocks 127, 255, 383, ..., 4095

Take note that there will be 32 possible MM blocks that can reside in a particular cache block.

2,048 bytes
16 bytes / block

DIRECT MAPPING

The Memory System



ASSUME:

MM SIZE = 65,536 WORDS

CACHE SIZE = 2,048 WORDS

BLOCK SIZE = 16 WORDS / BLOCK

$$n = 16 \text{ bits} //$$

$$\# \text{ of MMBS} = 4,096 (0 - 4,095) -$$

$$\# \text{ of CBS} = 128 (0 - 127) -$$

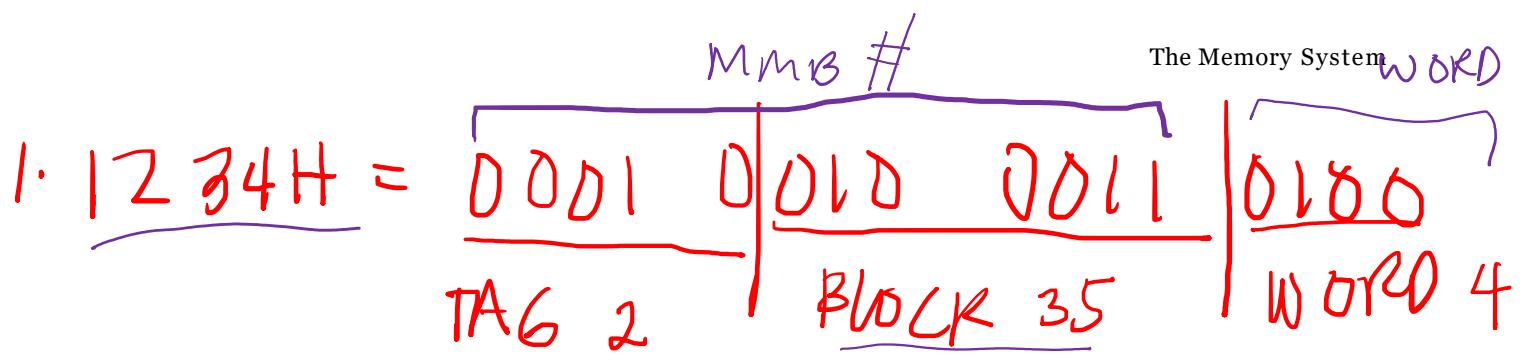
$$\# \text{ of TAG bits} = \log \left(\frac{\# \text{ of MMBS}}{\# \text{ of CBS}} \right)$$

$$= \frac{\log (4,096 / 128)}{\log 2} = 5 \text{ bits} //$$

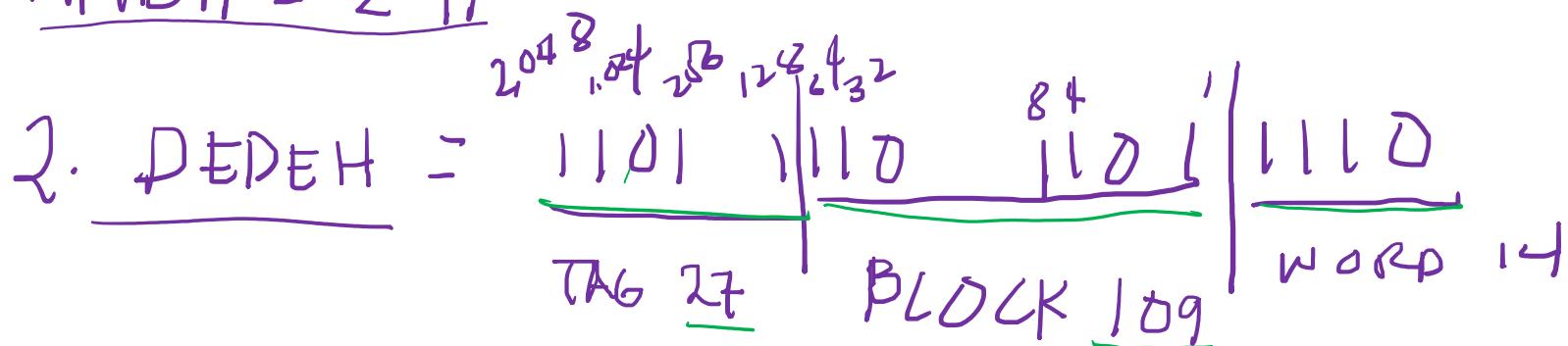
$$\# \text{ of Block bits} = \frac{\log \# \text{ of CBS}}{\log 2} = \frac{\log 128}{\log 2} = 7 \text{ bits}$$

$$\# \text{ of word bits} = \frac{\log \text{ block size}}{\log 2} = \frac{\log 16}{\log 2} = 4 \text{ bits}$$

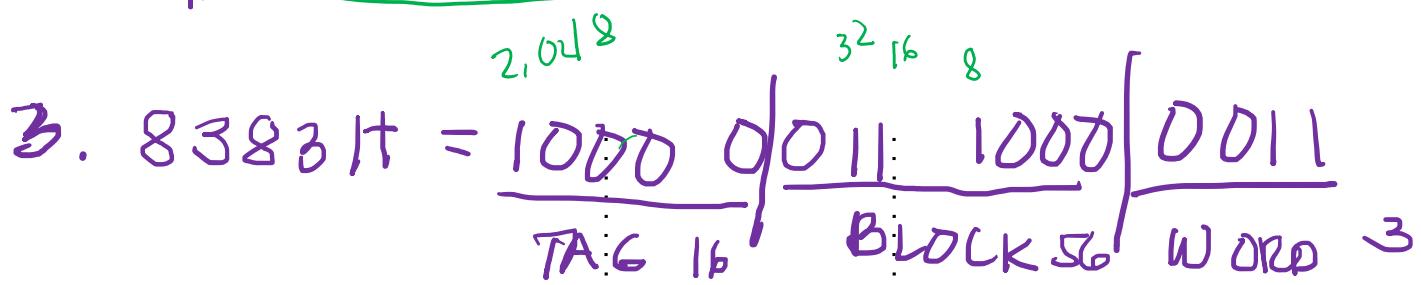
$$\text{mmppf} = \text{TAG} \# \times \# \text{ of CBS} + \text{Block} \#$$



$$\text{MMB\#} = 291$$



$$\text{MMB\#} = 3,565$$



$$\text{MMB\#} = 2104$$

$$\text{MMB\#} = \text{TAG\#} \times \text{\# of CBs} + \text{BLOCK\#}$$

$$1. (2 \times 128) + 35 = 291$$

$$2. (27 \times 128) + 109 = 3565$$

$$3. (16 \times 128) + 56 = 2164$$

1. 1ACEH = 0001 | 010 11 00 | 1110
| TAG 3 | BLOCK 44 | WORD 4

$$MMR\# = \underline{428}$$

2. 2BTAH = 001D | 011 1110 | 101D
| TAG 5 | BLOCK 62 | WORD 10

$$MMR\# = \underline{702}$$

3. 3CA1H =

4. 4D1AH.

Placement of a block in the cache is determined from the memory address. The memory address can be divided into three fields:

	TAG	BLOCK	WORD
MM ADDRESS =	5	7	4

The low-order 4 bits select one of 16 words in a block. When a new block enters the cache, the 7-bit (since there are 2^7 or 128 cache blocks) cache block field determines the cache position in which this block must be stored. The high order 5 bits (since there are 2^5 or 32 possible MM blocks that can reside in a particular cache block) of the memory address of the block are stored in the five *tag bits* associated with its location in the cache.

As execution proceeds, the 7-bit cache block field of each address generated by the processor points to a particular block location in the cache. The high order 5 bits of the address are compared with the tag bits associated with that cache location. If they match, then the desired word is in that block of the cache. If there is no match, then the block containing the required word must first be read from the main memory and loaded into the cache.

The replacement algorithm for the direct mapping technique is trivial. The direct mapping technique is easy to implement but it is not very flexible.

Example:

A 16 MB main memory is to be upgraded by adding a cache memory of 128 KB. Both the main memory and the cache memory are partitioned into blocks of 512 bytes. How many bits make up the tag, block, and word fields of a memory address?

Solution:

$$\begin{aligned}\text{No. of MM Blocks} &= 16 \text{ MB} / 512 \\ &= 32,768 \\ \text{No. of Cache Blocks} &= 128 \text{ KB} / 512 \\ &= 256\end{aligned}$$

$$\begin{aligned}\text{No. of MM Blocks per Cache Block} &= 32,768 / 256 \\ &= 128\end{aligned}$$

$$\begin{aligned}\text{No. of Tag Bits} &= \log 128 / \log 2 \\ &= 7 \text{ bits} \\ \text{No. of Block Bits} &= \log 256 / \log 2 \\ &= 8 \text{ bits} \\ \text{No. of Word Bits} &= \log 512 / \log 2 \\ &= 9 \text{ bits}\end{aligned}$$

Assume:

2. Associative Mapping

MM size - 65,536 bytes

Cache size - 2,048 bytes
A main memory block can potentially reside in any cache block position.

block size - 16 bytes / block

$n = 16 \text{ bits}$

tag	block 0
tag	block 1
tag	block 2
tag	block 3

MM blocks 0, 1, 2, 3, ..., 4095

of MMBS = 4,096 MBs

of CBS = 128 CBS

TAG bits =

$\log \# \text{ MMBS}$

tag	block 125
tag	block 126
tag	block 127

MM blocks 0, 1, 2, 3, ..., 4095

MM blocks 0, 1, 2, 3, ..., 4095

MM blocks 0, 1, 2, 3, ..., 4095

of word bits = $\log \text{ block size}$

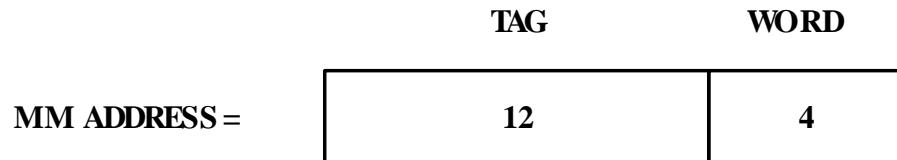
Take note that there will be 4,096 possible main memory blocks that can reside in a particular cache block.

$$\Rightarrow \frac{\log 16}{\log 2} = 4 \text{ bits} //$$

1. $1234H = \underline{0001 \quad 0010 \quad 0011} \Big| \underline{0100}$
 $TAB\# = 291$ WORD 4
MMB# = 291

2. $ABCDH = \underline{1010 \quad 1011 \quad 1100} \Big| \underline{1101}$
 $TAB\# 2748$ WORD 13
MMB# 2748

This scheme requires 12 tag bits (because there are 2^{12} or 4,096 possible main memory blocks that can reside in a particular cache block) to identify a main memory block when it is resident in the cache. The system will compare the tag bits of an address received from the processor to the tag bits of each block of the cache to see if the desired block is present.



Since there is complete freedom in block positioning, a wide range of replacement algorithms is possible. However, it might not be practical to make full use of this freedom, because complex replacement algorithms may be difficult to implement.

The cost of an associative cache is higher than the cost of a direct-mapped cache because of the need to search all 128 tag patterns to determine whether a given block is in the cache. A search of this kind is called an *associative search*. For performance reasons, the tags must be searched in parallel.

Example:

A 16 MB main memory is to be upgraded by adding a cache memory of 128 KB. Both the main memory and the cache memory are partitioned into blocks of 512 bytes. How many bits make up the tag and word fields of a memory address?

Solution:

$$\text{No. of MM Blocks} = 16 \text{ MB} / 512$$

$$= 32,768$$

$$\text{No. of Cache Blocks} = 128 \text{ KB} / 512$$

$$= 256$$

$$\text{No. of Tag Bits} = \log 32,768 / \log 2$$

$$= 15 \text{ bits}$$

$$\text{No. of Word Bits} = \log 512 / \log 2$$

$$= 9 \text{ bits}$$

ASSUMING:

3. Block-Set-Associative Mapping

MMsize - 65,536 bytes



Cache size - 2,048 bytes
 Cache of cache are grouped into sets, and the mapping allows a block of main memory to reside in any cache block of a specific set.
 Block size - 16 bytes/block

Example:
set

$$n = \frac{\log 65,536}{\log 2}$$

$$n = 16 \text{ bits}$$

$$\# \text{ of MMBs} = 4,096 \text{ MMBs}$$

$$\# \text{ of CBS} = 128 \text{ CBS}$$

$$\# \text{ of TAG bits} = \log(\# \text{ of MMBs} / \# \text{ of set})$$

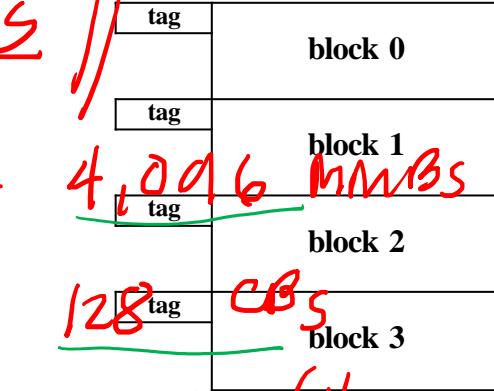
$$\text{Assume } 2 \text{ cache blocks per set}$$

CACHE

= 64 SET

SET 0 = MM blocks 0, 64, 128, ..., 4032

SET 1 = MM blocks 1, 65, 129, ..., 4033



$$= \log(4,096 / 64) = 6 \text{ bits}$$

SET 62 = MM blocks 62, 126, 190, ..., 4094



SET 63 = MM blocks 63, 127, 191, ..., 4095

$$= \log 64 = 6 \text{ bits}$$

$$\# \text{ of Set bits} =$$

$$\# \text{ of word bits} =$$

$$= \frac{\log 16}{\log 2} = 4 \text{ bits}$$

$$\# \text{ of TAG bits} = \frac{\log(\# \text{ of MmBs} / \text{set})}{\log_2}$$

$$= \frac{\log(4,096 / 64)}{\log_2}$$

$$= 6 \text{ bits} //$$

$$\# \text{ of set bit} = \frac{\log(\# \text{ of set})}{\log_2} = \frac{\log 64}{\log_2}$$

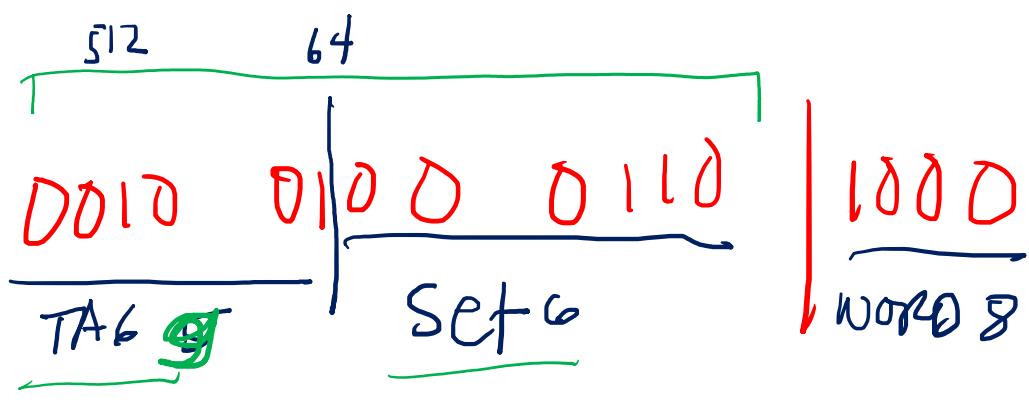
$$\# \text{ of word bits} = \frac{\log \text{ block size}}{\log_2}$$

$$= \frac{\log 16}{\log_2} = 4 \text{ bits} //$$

Example:

1. $2468H =$

$\text{MMB} \# = 582$



BSAM

$$\text{MMB\#} = \text{TAG\#} \times \#\text{Set} + \text{Set\#}$$

$$= 9 \times 64 + 6 = 582$$

1. 1234H = 0001 00 | 10 0011 | D100
 TAG 4 Set 35 WORD 4

$$\text{MMB\#} = \underline{\underline{291}} //$$

..

2. 1ACEH = 0001 10 | 10 1100 | 1110
 TAG 6 SET 44 WORD 14

$$\text{MMB\#} = 428 //$$

$$2 \cdot (4 \times 64) + 35 = 291 //$$

$$3 \cdot (6 \times 64) + 44 = 428 //$$

The contention problem of the direct method is eased by having a few choices for block placement. At the same time, the hardware cost is reduced by decreasing the size of the associative search.

	TAG	SET	WORD
MM ADDRESS =	6	6	4

The 6-bit set field (since there are 2^6 or 64 sets) of the address determines which set of the cache might contain the desired block. The 6-bit tag field (since there are 2^6 or 64 possible main memory blocks that can reside in a particular cache set) of the address must then be associatively compared to the tags of the two blocks of the set if a match occurs signifying block presence.

The extreme condition of having 128 blocks (all cache blocks) per set corresponds to the fully associative technique. The other extreme of having one block per set is the direct-mapping technique.

The block-set-associative mapping technique is the most practical cache mapping technique.

ASSIGNMENT

Example:

DM ✓

AM ✓

BSAM ✓

1. 12 34 5H

Solution:

2. 369 B DH

3. 2 468 A H

4. A B C D E H

5. 4 1 0 7 3 H

A 16 MB main memory is to be upgraded by adding a cache memory of 128 KB. Both the main memory and the cache memory are partitioned into blocks of 512 bytes. Assume that there are 8 blocks to a set. How many bits make up the tag, set, and word fields of a memory address?

$$\begin{aligned} \text{No. of MM Blocks} &= 16 \text{ MB} / 512 \\ &= 32,768 \end{aligned}$$

$$\begin{aligned} \text{No. of Cache Blocks} &= 128 \text{ KB} / 512 \\ &= 256 \end{aligned}$$

$$\begin{aligned} \text{No. of Sets} &= 256 / 8 \\ &= 32 \end{aligned}$$

$$\begin{aligned} \text{No. of MM Blocks per Set} &= 32,768 / 32 \\ &= 1,024 \end{aligned}$$

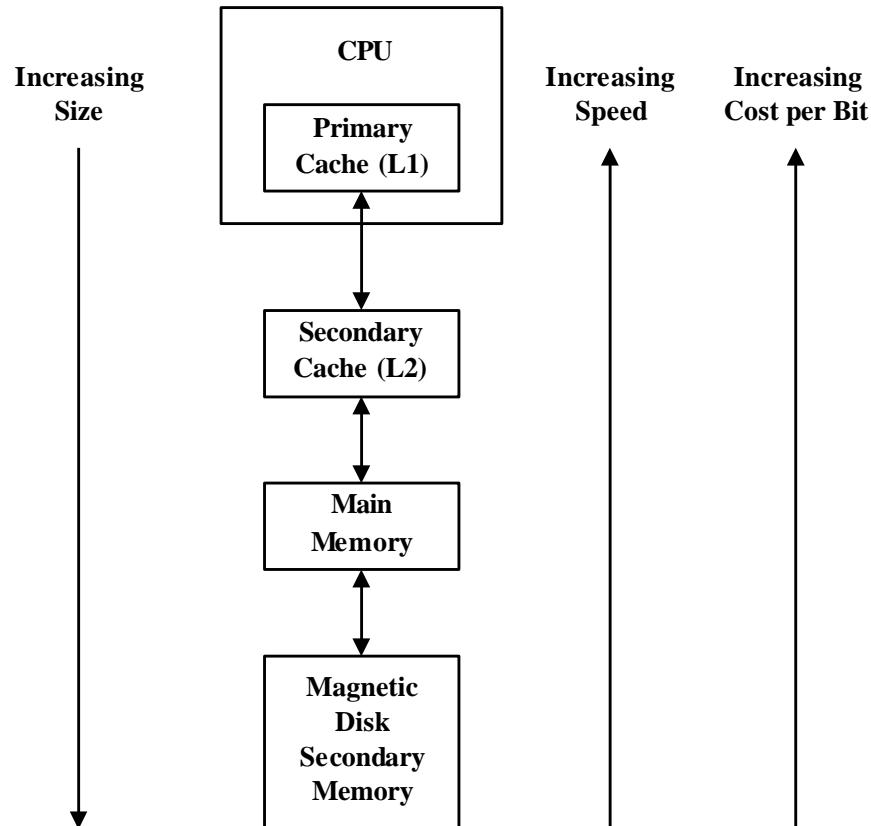
$$\begin{aligned} \text{No. of Tag Bits} &= \log 1,024 / \log 2 \\ &= 10 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{No. of Set Bits} &= \log 32 / \log 2 \\ &= 5 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{No. of Word Bits} &= \log 512 / \log 2 \\ &= 9 \text{ bits} \end{aligned}$$

MEMORY HIERARCHY

- ☞ The entire memory can be viewed as the hierarchy depicted below:



- ☞ The fastest access to data is through the processor registers. If registers are considered to be part of the memory hierarchy, then the processor registers are at the top in terms of speed of access. Of course, the registers provides only a minuscule portion of the required memory.

- ☞ At the next level is the cache memory. There are two types of cache memory. A *primary cache* is located on the processor chip. This cache is small, because it competes for space on the CPU chip, which must implement many other functions. The primary cache is referred to as *Level 1 (L1)* cache.

A *secondary cache* is placed between the primary cache and the main memory. It is referred to as *Level 2 (L2)* cache.

Including a primary cache on the processor chip and using a larger, off-chip, secondary cache is the most common way of designing computers. However, other arrangements can be found in practice. It is possible not to have a cache on the processor chip at all. Also, it is possible to have both L1 and L2 caches on the processor chip.

- ☞ The next level in the hierarchy is the main memory. It is much larger but significantly slower than the cache memory. In a typical memory, the access time for the main memory is about ten times longer than the access time for the L1 cache.
- ☞ Disk devices provide a huge amount of inexpensive storage (secondary storage). They are very slow compared to the semiconductor devices used to implement main memory.



Pentium 4 Caches

The Pentium 4 processor can have up to three levels of caches.

The L1 cache consists of separate data and instruction caches. The data cache has a capacity of 8 KB, organized in a 4-way set-associative manner. Each cache block has 64 bytes. The write-through policy is used to write to the cache.

Data can be accessed from the data cache in two clock cycles. For a 1.4 GHz Pentium, data can be accessed in less than 1.5 ns. The instruction cache does not hold normal machine instructions. Instead, it holds already decoded versions of the instructions.

The L2 cache is a unified cache with a capacity of 256 KB, organized in an 8-way set-associative manner. Each of its blocks comprises 128 bytes. The write-back policy is used to write to the cache. The access latency of this cache is seven clock cycles.

Both L1 and L2 caches are implemented on the processor chip. The architecture also allows for inclusion of an on-chip L3 cache. However, this cache is not implemented in the Pentium 4 chips targeted for desktop computers. It is intended for processor chips used in server systems.

PERFORMANCE CONSIDERATIONS

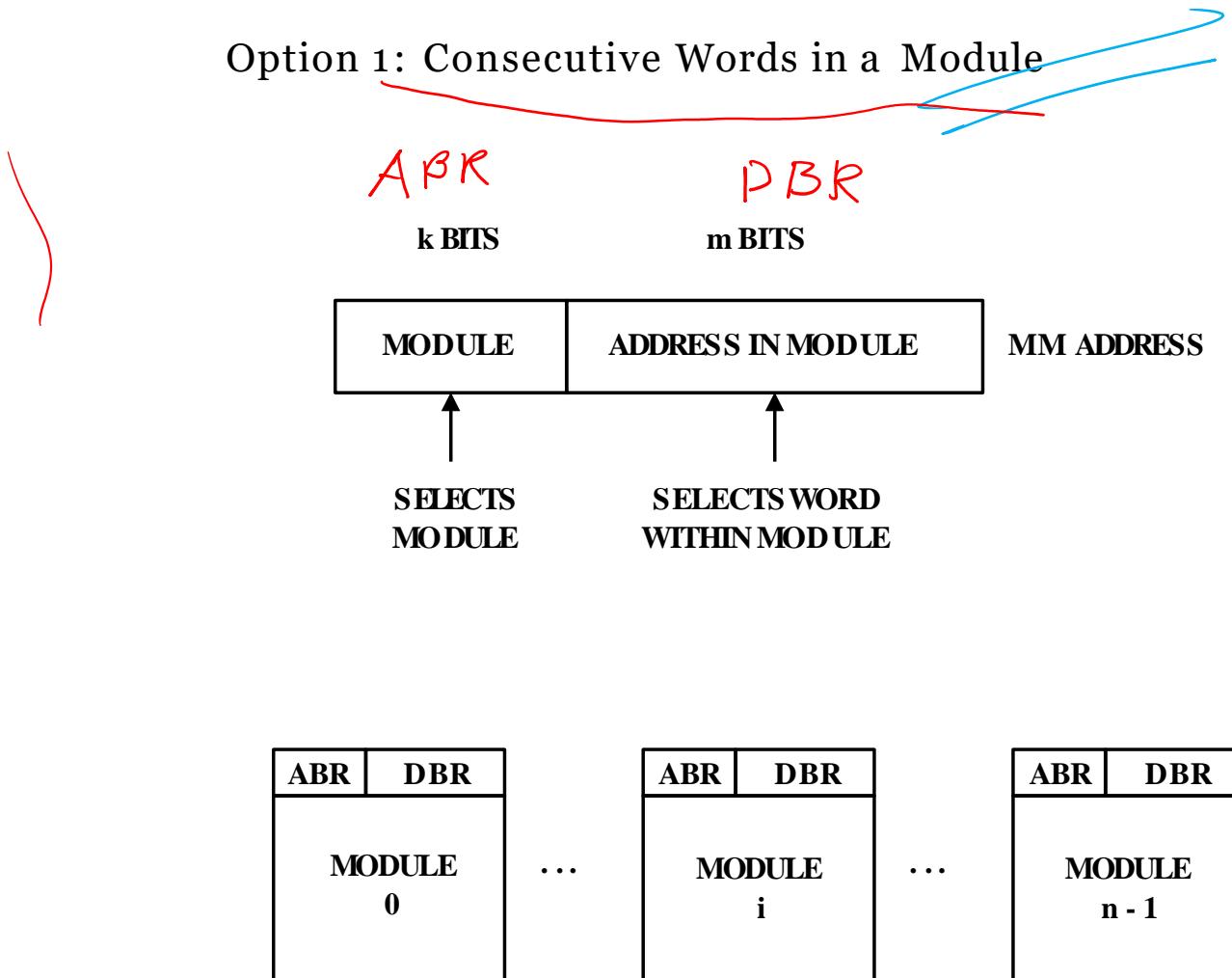
- ☞ Two key factors in the commercial success of a computer are performance and cost; the best possible performance at the lowest possible cost is the objective.
- ☞ A common measure of success is the *price/performance ratio*.
- ☞ Performance depends on how fast machine instructions can be brought into the processor for execution and how fast they can be executed.
- ☞ The memory hierarchy results from the quest for the best price/performance ratio. The main purpose of this hierarchy is to create a memory that the processor sees as having a short access time and a large capacity.
- ☞ The speed and efficiency of data transfer between various levels of the hierarchy are also of great significance. It is beneficial if transfers to and from the faster units can be done at a rate equal to that of the faster unit. This is not possible if both the slow and the fast units are accessed in the same manner, but can be achieved when parallelism is used in the organization of the slower unit.
- ☞ An effective way to introduce parallelism is to use an *interleaved organization*.



Interleaving

If the main memory of a computer is structured as a collection of physically separate modules, each with its own address buffer register (ABR) and data buffer register (DBR), memory access operations may proceed in more than one module at the same time.

Option 1: Consecutive Words in a Module



Example: MM Size = 1,048,576 locations
(20 address lines or bits)

$$\text{MM size} = 1M (1,048,576 \text{ bytes})$$

Option 1 divides the main memory into 32 modules of 32,768 locations each.

$$\text{Module size} = 32,768 \text{ byte/module}$$

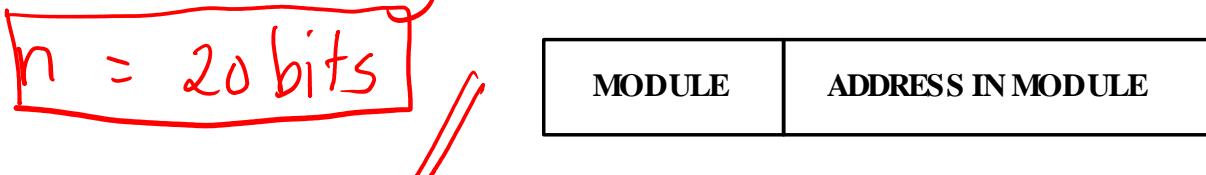
The high-order 5 bits selects one of the 32 modules while low-order 15 bits select one of the 32,768 locations in that module.

$$n = \underline{\log 1,048,576}$$

$$\log_2$$

5 BITS

15 BITS

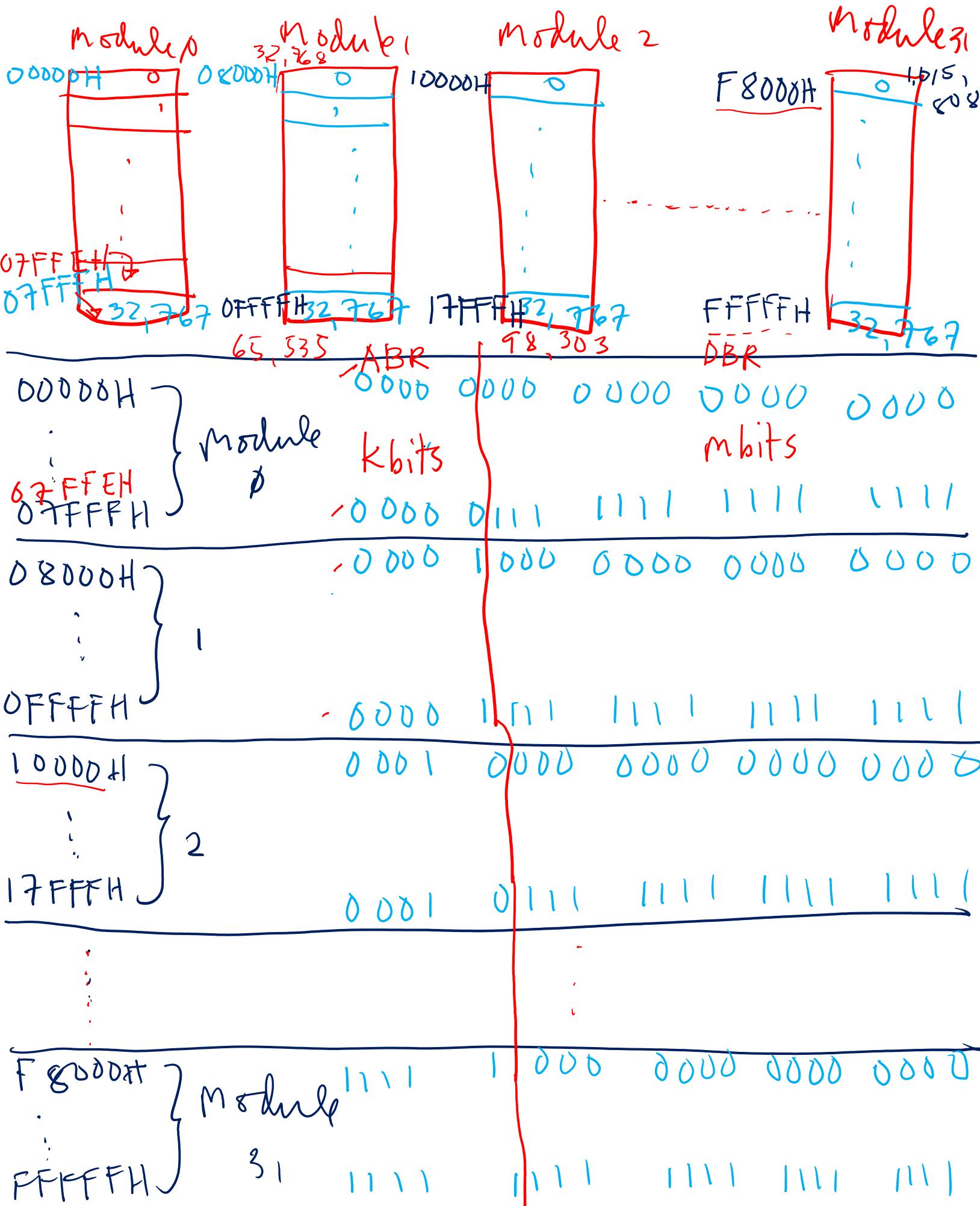


When consecutive locations are accessed, as happens when a block of data is transferred to a cache, only one module is busy. At the same time, however, devices with direct memory access (DMA) capability may be accessing information in other memory modules.

$$\# \text{ of Modules} = \frac{\text{MM size}}{\text{Module size}} = \frac{1,048,576 \text{ bytes}}{32,768 \text{ bytes}} = \boxed{32 \text{ modules}}$$

$$\begin{aligned}\text{\# of address in a} \\ \text{Module bits} &= \frac{\log \text{\# of Module size}}{\log 2} \\ &= \frac{\log 32,768}{\log 2} \\ &= 15 \text{ bits}\end{aligned}$$

$$\begin{aligned}\text{\# of Module bits} &= \frac{\log (\text{\# of Modules})}{\log 2} \\ &= \frac{\log 32}{\log 2} \\ &= 5 \text{ bits}\end{aligned}$$



Example:

1.07 FFFFH

The diagram illustrates the division of a 16-bit address into two parts: a module identifier and a local address. A horizontal line is divided into four segments by three vertical red lines. The first segment contains the binary value '0000'. The second segment contains '111'. The third segment contains '111'. The fourth segment contains '111'. To the left of the first vertical line is the text 'Module φ'. To the right of the second vertical line is the text 'Address in a module'.

$$2. \underline{10000H} =$$

Module 2 | 0001 0000 0000 0000
address in a module 0

$$3 \cdot |12345H| =$$

$\text{E}_{\text{1311}} =$
0001 010 0011 0100 0101
Module 2 9,629

Example: 2

1. word in a module 3 of module 1

$$\begin{array}{r} \underline{0000} & \underline{0000} \\ \underline{0000} & \underline{1000} \\ 0 & 8 \end{array} \quad \begin{array}{r} \underline{0000} \\ \underline{0000} \\ 0 \end{array} \quad \begin{array}{r} \underline{0000} \\ \underline{0000} \\ 0 \end{array} \quad \begin{array}{r} \underline{0000} \\ \underline{0000} \\ 0 \end{array} \quad \begin{array}{r} \underline{0000} \\ \underline{0011} \\ 3 \end{array}$$

08003H

2. word in a module 5 of module 16

3. 10883H = 0061 0800 0080 0068 0011

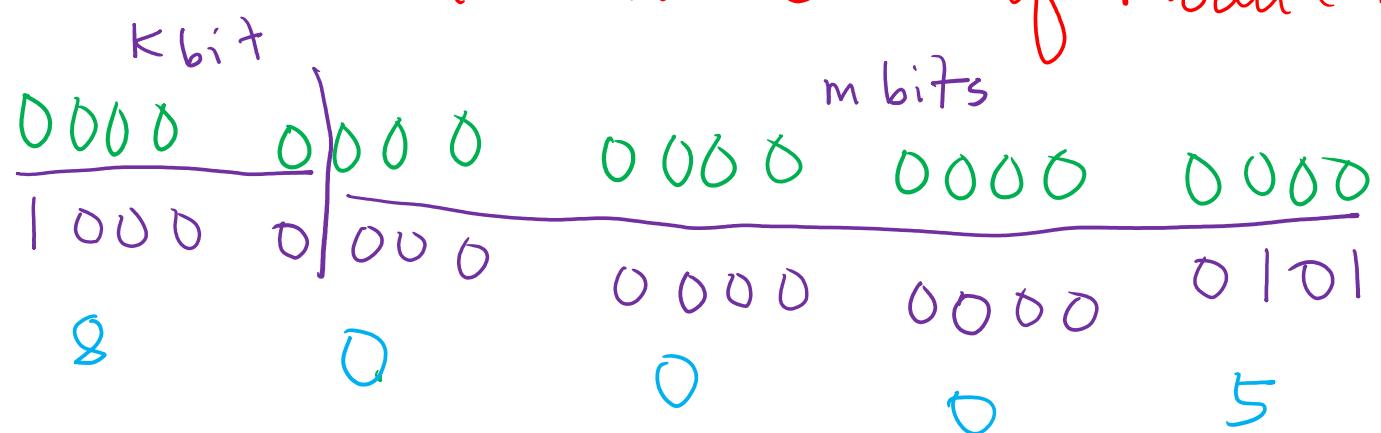
module 2 | words in module 3
module 9

4. Word in a module 80 of the 10th module

$$\begin{array}{r} \underline{0680} & \underline{0000} & \underline{0000} & \underline{0000} & \underline{0000} \\ \underline{0168} & \underline{1000} & \underline{0000} & \underline{0101} & \underline{0000} \\ 4 & 8 & 0 & 5 & 0 \end{array}$$

48030H ✓

2. Word in a module 5 of module 16



80005H

Last words in a module of module 20

$$\begin{array}{r}
 6666 \\
 1010 \\
 \hline
 A
 \end{array}
 \quad
 \begin{array}{r}
 0|000 \\
 0|111 \\
 \hline
 7
 \end{array}
 \quad
 \begin{array}{r}
 0000 \\
 1111 \\
 \hline
 F
 \end{array}
 \quad
 \begin{array}{r}
 0000 \\
 1111 \\
 \hline
 F
 \end{array}
 \quad
 \begin{array}{r}
 0000 \\
 1111 \\
 \hline
 F \quad H
 \end{array}$$

module 20 ≠ 20th module
21st module ≠ module 19

Exercise:

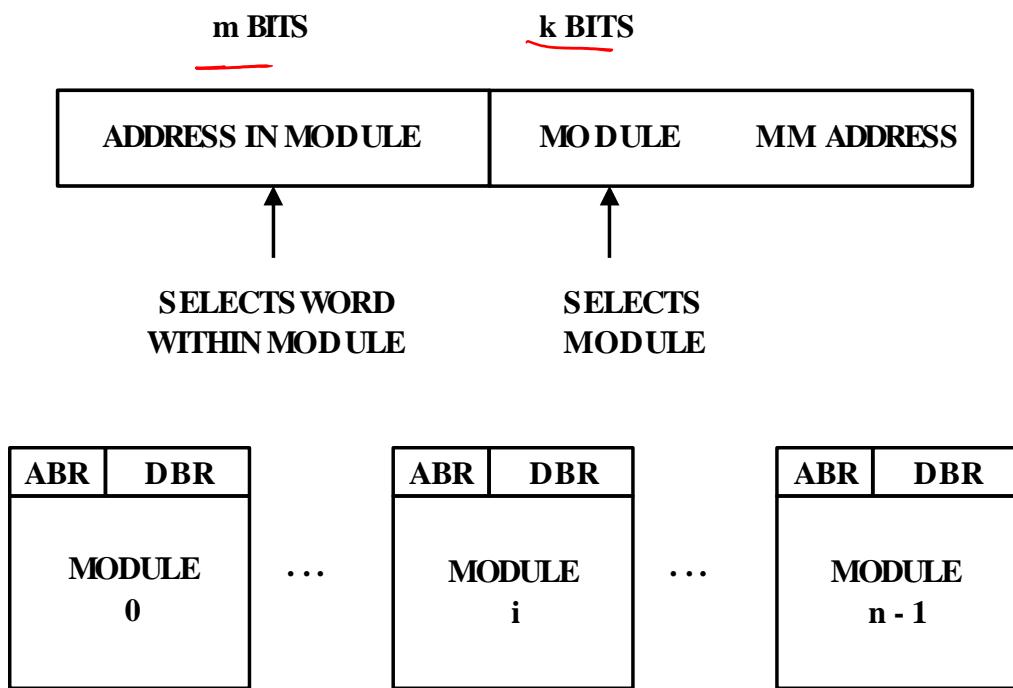
1. 1049A+1

~~0001~~ 0000 | 0100 1001 1010
Module 2 words in a module 1178

2. words in a module 300 of the 20th

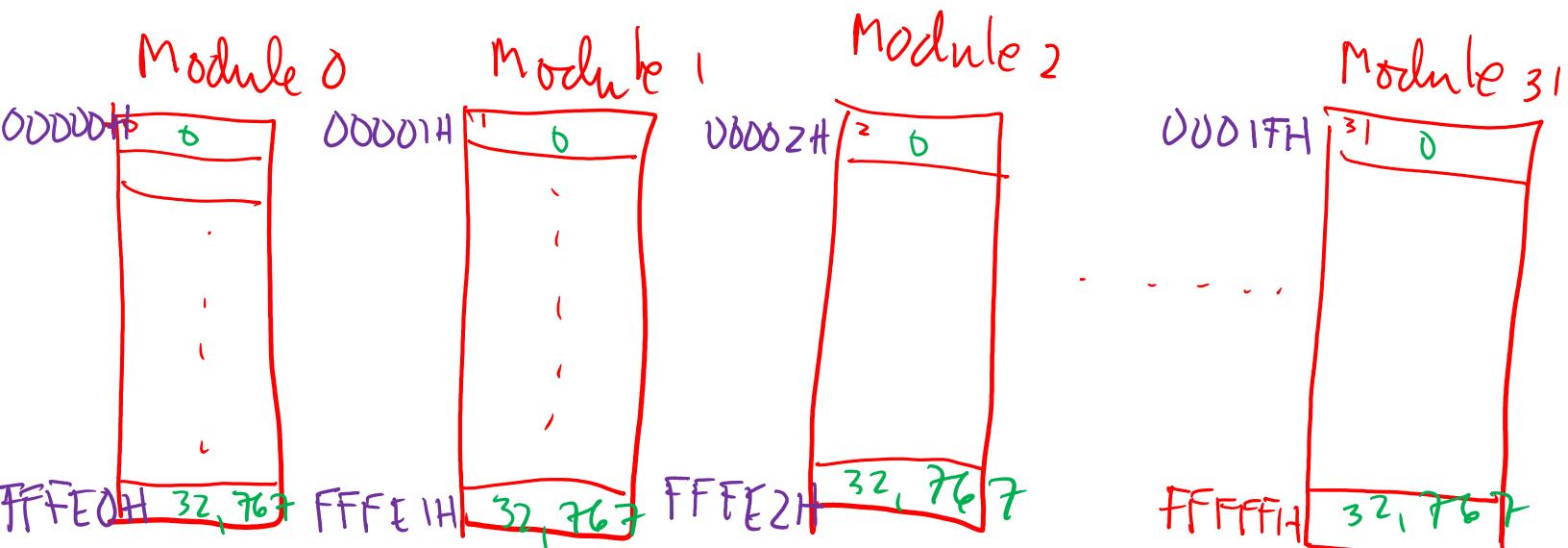
6660 0000 0000 0000 0000 module
 1601 | 000 000 | 6010 | 100 = 9812CH

Option 2: Consecutive Words in Consecutive Modules or Memory Interleaving



In the second option, the low order bits of the memory address select the a module, and the high order bits select a location within that module. In this way, consecutive addresses are located in successive modules. Thus, any component of the system that generates requests for access to consecutive memory locations can keep several modules busy at one time. This results in both faster access to a block of data and higher average utilization of the memory system as a whole.

To take advantage of memory interleaving, the system must be capable of initiating a memory access operation while waiting for the completion of a previous memory access.



Examples:

$00020H = \underbrace{0000 \ 0000 \ 0000}_{\text{words in a module}} \underbrace{0010 \ 0000}_{\text{address in a module}} \ 0$

1

$12345H = \underbrace{\begin{matrix} 0001 & 0010 & 0011 & 0100 \end{matrix}}_{\text{address in a module}} \underbrace{0101}_{\text{Module 5}}$

$\begin{matrix} 1^2 \\ 512 \\ 2^{57} \\ 128 \\ 4 \end{matrix}$

$\begin{matrix} 900 \\ \text{words/address in a module 900 of} \\ \text{the last module} \end{matrix}$

$\begin{matrix} 0000 & 0000 & 0000 & 0000 \\ 0000 & 0111 & 0000 & 1001 \end{matrix}$

$\begin{matrix} 0 \\ 512 \\ 7 \end{matrix}$

$\begin{matrix} 0 \\ 2^8 \\ 0 \end{matrix}$

$\begin{matrix} 0 \\ 4 \\ 9 \end{matrix}$

$\begin{matrix} 0000 \\ 1111 \end{matrix} =$

$F0709FH$

EXERCISE : 231st address in
a module =
address in a module 230

1. $340F9H$ = address in a module 6663

0011 0100 0000 1111 1001
address in a module Module 25

2. 260th address in a module of
the 30th module

0000 0000 0000 0000 0000
0000 0010 0000 0111 1101

0207DH

ordinal number

1st module =

2nd module =

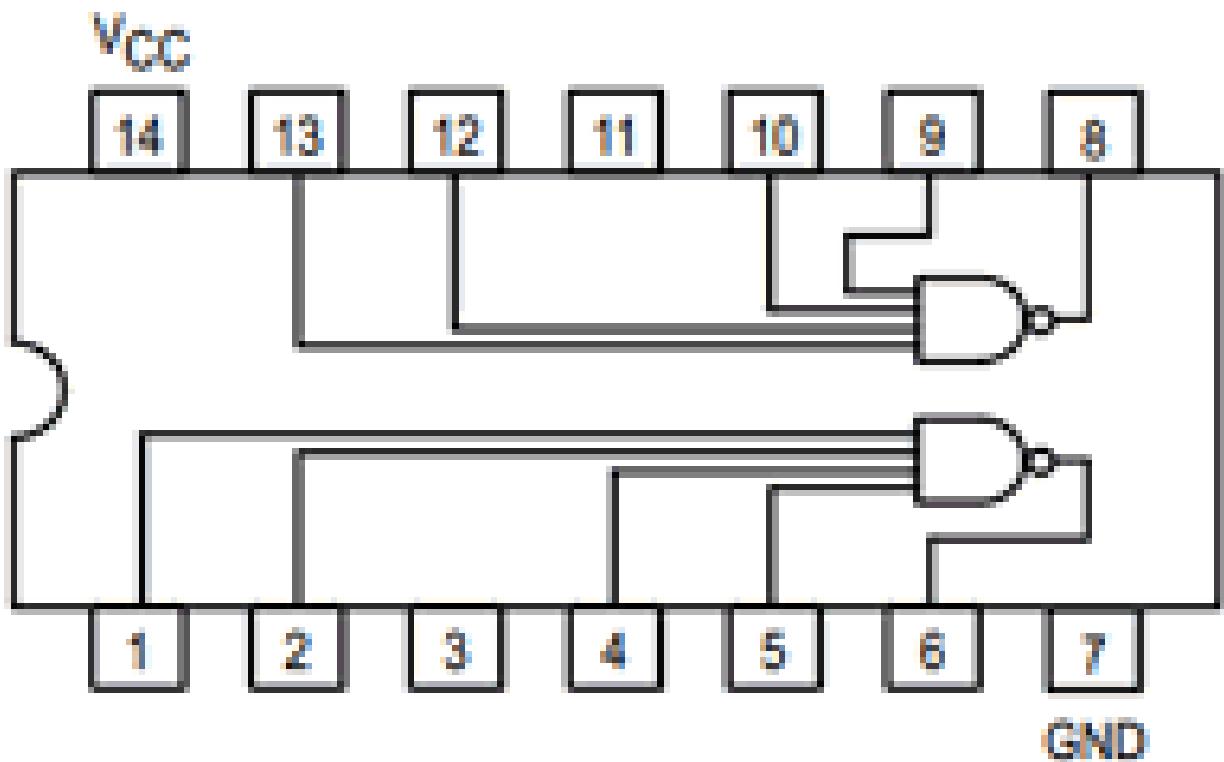
29th module =

module 0 ✓

module 1 ✓

module 28 ✓

This is a 4-input NAND gate-74LS20
Pin Configuration



$4 \rightarrow 1$

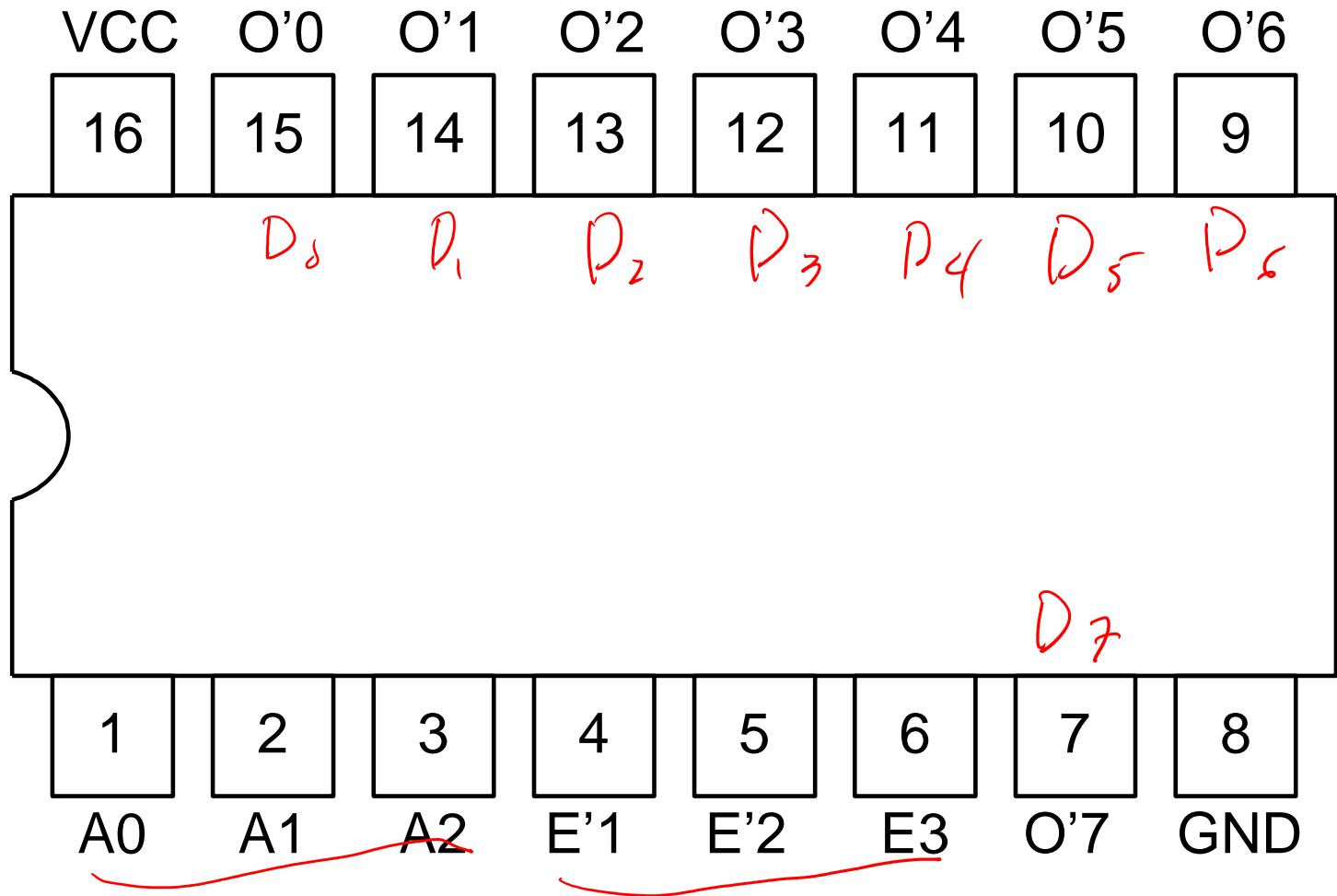
$2^2 \rightarrow 1^2$

$a \quad b \quad c \quad | \quad P_0 \quad P_1 \quad P_2 \quad P_3 \quad P_4 \quad P_5 - P_6 \quad P_7$

0	0	0		1	0	0	0	0	0	0	0	0
0	0	1		0	1	0	0	0	0	0	0	0
0	1	0		0	0	1	0	0	0	0	0	0
0	1	1		0	0	0	1	0	0	0	0	0
1	0	0		0	0	0	0	1	0	0	0	0
1	0	1		0	0	0	0	0	1	0	0	0
1	1	0		0	0	0	0	0	0	1	0	0
1	1	1		0	0	0	0	0	0	0	1	0

$n \rightarrow 2^n$ decoder

74LS138 Pin Configuration



Ungsod

Lim

Discaya

Mahaguay

Tanedo

Jawer

OPTION 1

OPTION 2

$$1. \ 3CA13H = \frac{0011}{\text{Module}} \quad | \quad \begin{array}{cccccc} 00 & 1010 & 0001 & 0011 \\ & & & & & \\ & & & & & 18963 \end{array}$$

$$2. \ OA1B2H \quad 7$$

3 - last word in a module of

Option 2:

$$\begin{array}{c} 2^{18} \text{ module } 9 \\ 3CA13H = \frac{0011}{4096} \quad \frac{1106}{1024} \quad \frac{1010}{512} \quad \frac{0001}{64} \quad \frac{0011}{16} \\ \text{Word in a module} \quad | \quad \text{Module } 19 \end{array}$$

$$\begin{array}{c} 8192 \quad 7,760 \\ 0000 \quad | \quad \frac{1010}{256} \quad \frac{0001}{128} \quad \frac{1011}{16} \quad \frac{0010}{2} \\ \text{Module } 1 \quad \text{Word in a module } 8626 \end{array}$$

$$\begin{array}{c} 1024 \quad 256 \\ 0000 \quad 1010 \quad 0001 \quad | \quad \frac{1011}{8} \quad \frac{1}{4} \quad \frac{1}{1} \\ \text{Word in a module } 1293 \quad \text{Module } 18 \end{array}$$

Last word in a module of memory

OPTION 1:

$$\begin{array}{r} 0000 \\ 0100 \\ \hline 4 \end{array} \quad \begin{array}{r} 0000 \\ | \\ 111 \\ \hline F \end{array} \quad \begin{array}{r} 0000 \\ ||| \\ 111 \\ \hline F \end{array} \quad \begin{array}{r} 0000 \\ | \\ 111 \\ \hline F \end{array} \quad \begin{array}{r} 0000 \\ | \\ 111 \\ \hline F \end{array}$$

4FFFFH

OPTION 2:

$$\begin{array}{r} 0000 \\ | \\ 111 \\ \hline F \end{array} \quad \begin{array}{r} 0000 \\ | \\ 111 \\ \hline F \end{array} \quad \begin{array}{r} 0000 \\ | \\ 111 \\ \hline F \end{array} \quad \begin{array}{r} 0000 \\ | \\ 110 \\ \hline E \end{array} \quad \begin{array}{r} 0000 \\ | \\ 1001 \\ \hline 9 \end{array}$$

FFFE9H