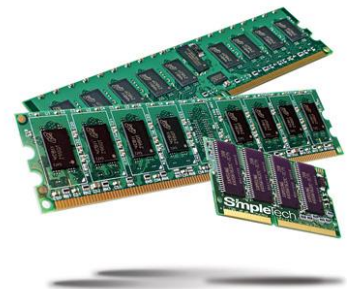
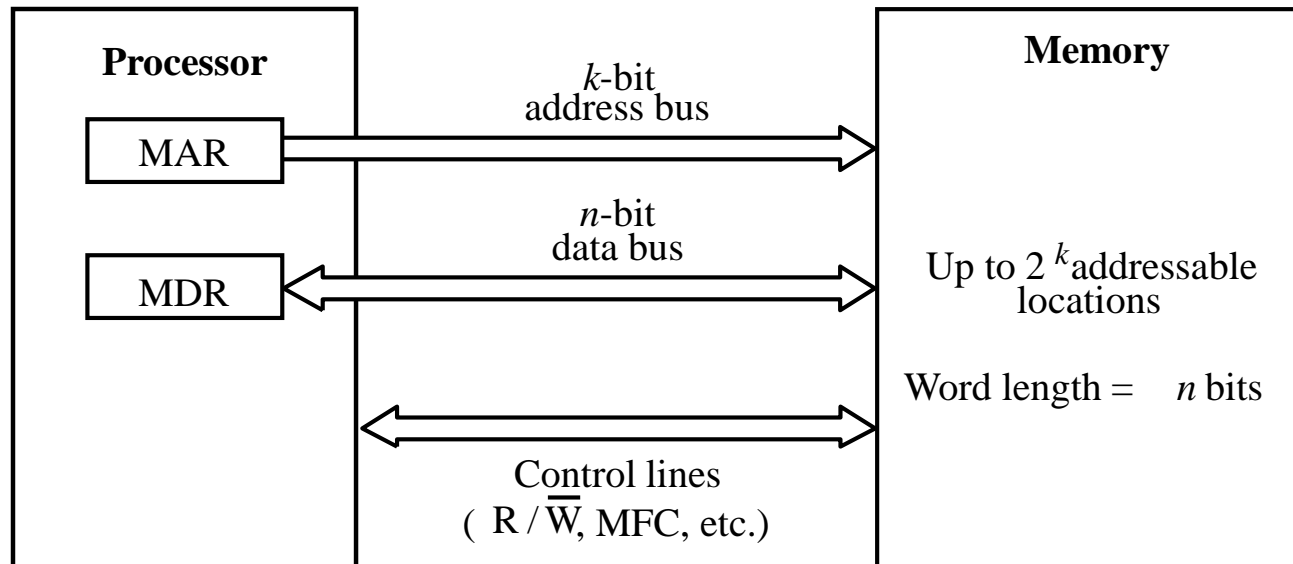


Fundamental Concepts

The Memory System

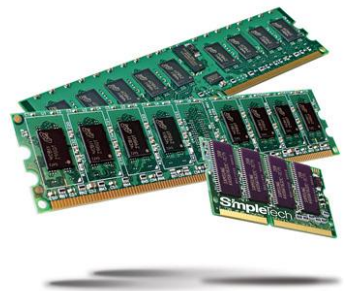
Some basic concepts

- Maximum size of the Main Memory
- byte-addressable
- CPU-Main Memory Connection



Some basic concepts(Contd.,)

- Measures for the speed of a memory:
 - memory access time.
 - memory cycle time.
- An important design issue is to provide a computer system with as large and fast a memory as possible, within a given cost target.
- Several techniques to increase the effective size and speed of the memory:
 - Cache memory (to increase the effective speed).
 - Virtual memory (to increase the effective size).

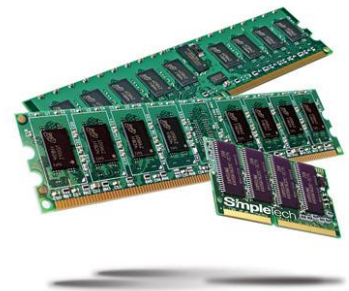


Semiconductor RAM memories

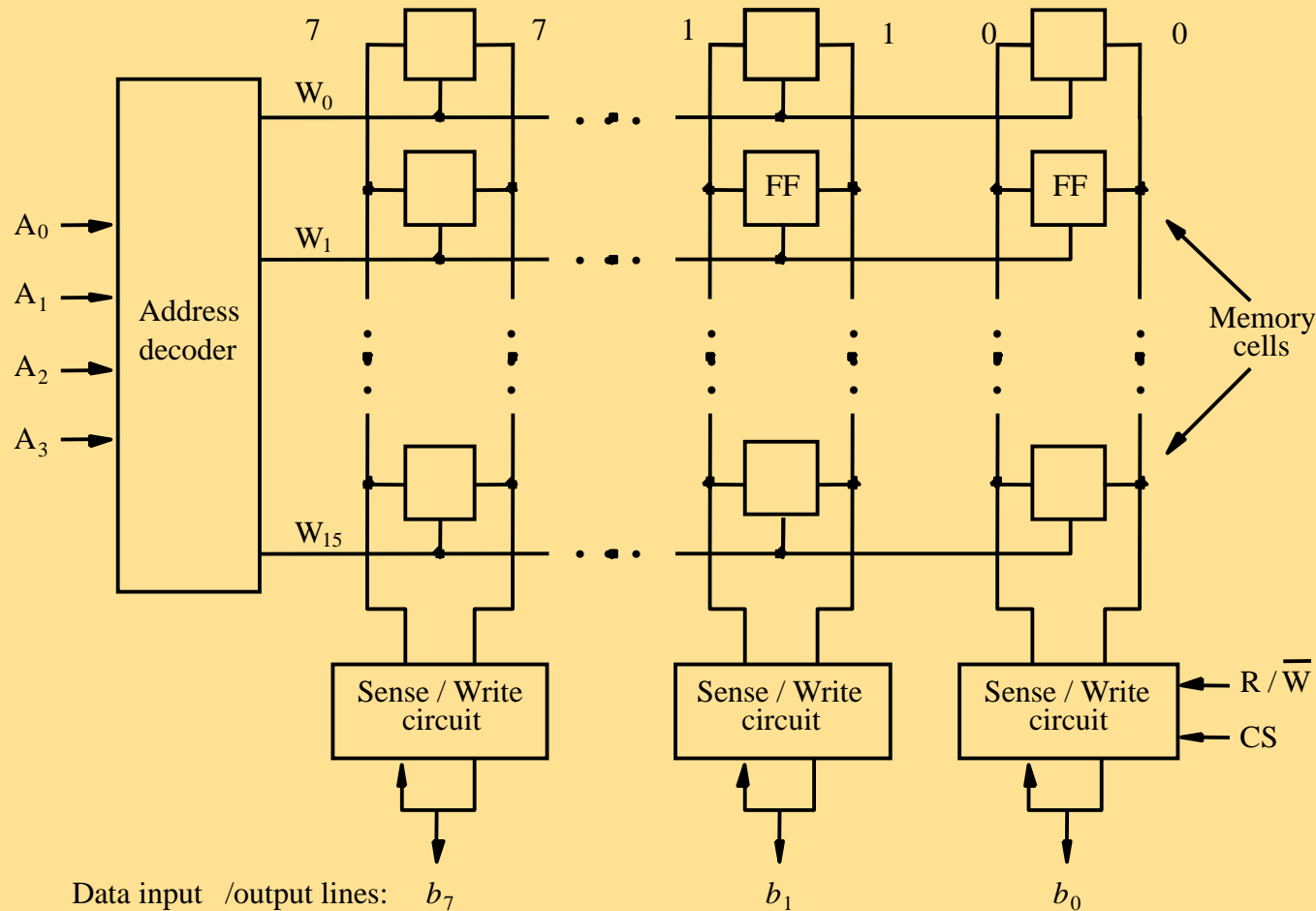
The Memory System

Internal organization of memory chips

- Each memory cell can hold one bit of information.
- Memory cells are organized in the form of an array.
- One row is one memory word.
- All cells of a row are connected to a common line, known as the “word line”.
- Word line is connected to the address decoder.
- Sense/write circuits are connected to the data input/output lines of the memory chip.



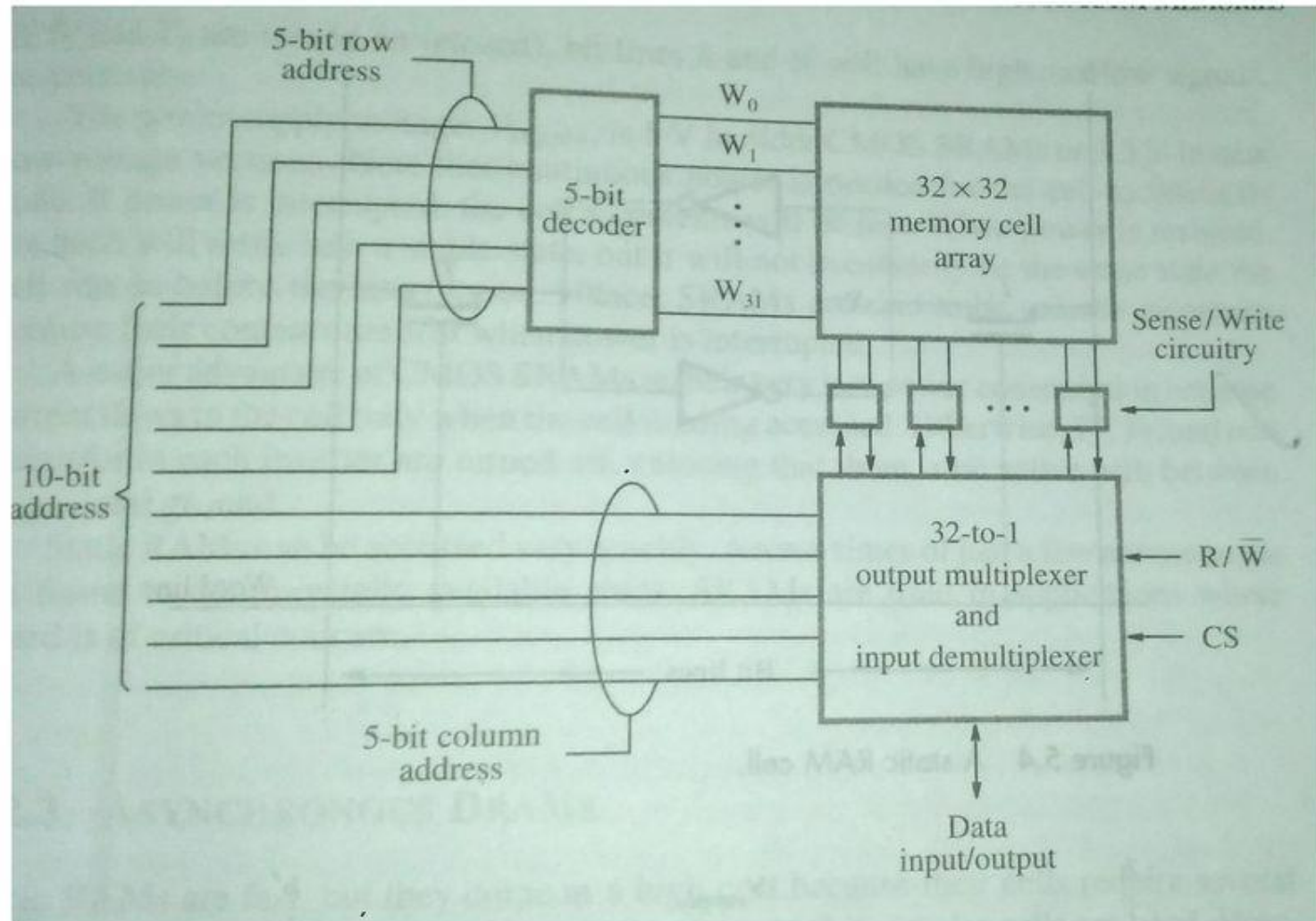
Internal organization of memory chips (Contd.,)





-
- ❑ Memory chip consisting of 16 words of 8 bit each.
 - ❑ This is referred to as a 16x8 organization.
 - ❑ It can store 128 bits and requires 14 external connections for address, data and control lines.
 - ❑ It also needed 2 lines for power supply and ground.
 - ❑ Consider a memory circuit with 1K cells.
 - ◆ this can be organized as a 128x8 memory.
 - ◆ Requiring a total of 19 external connections

Organization of 1K X 1 Memory Chip



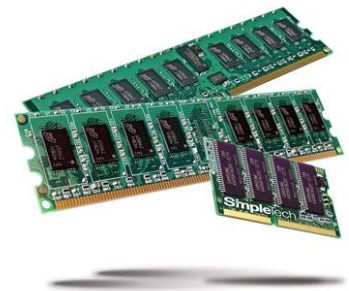
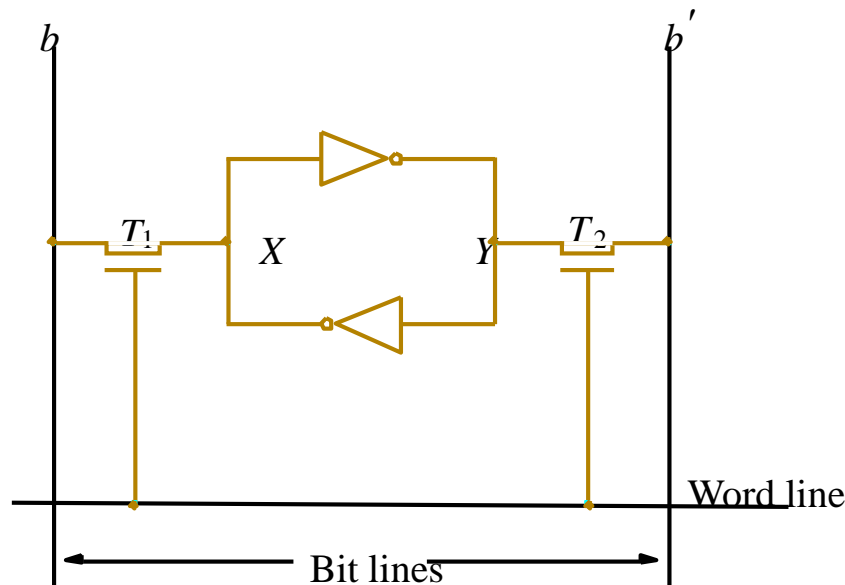


◆ Same can be organized into a **1Kx1** format

- 10 bit address is needed
 - But there is only one data line.
 - Requiring a total of **15 external** connections
- ◆ The required 10-bit address is divided into 2 groups of 5-bits each to form a row and column addresses for the cell array.
- ◆ Row address selects a row of 32 cells, all of which are selected in parallel
- ◆ According to the column address only one of these cells is connected to the external data line by the output multiplexer and the input demultiplexer

SRAM Cell

- Two transistor inverters are cross connected to implement a basic flip-flop.
- The cell is connected to one word line and two bits lines by transistors T_1 and T_2
- When word line is at ground level, the transistors are turned off and the latch retains its state
- Read operation: In order to read state of SRAM cell, the word line is activated to close switches T_1 and T_2 . Sense/Write circuits at the bottom monitor the state of b and b'





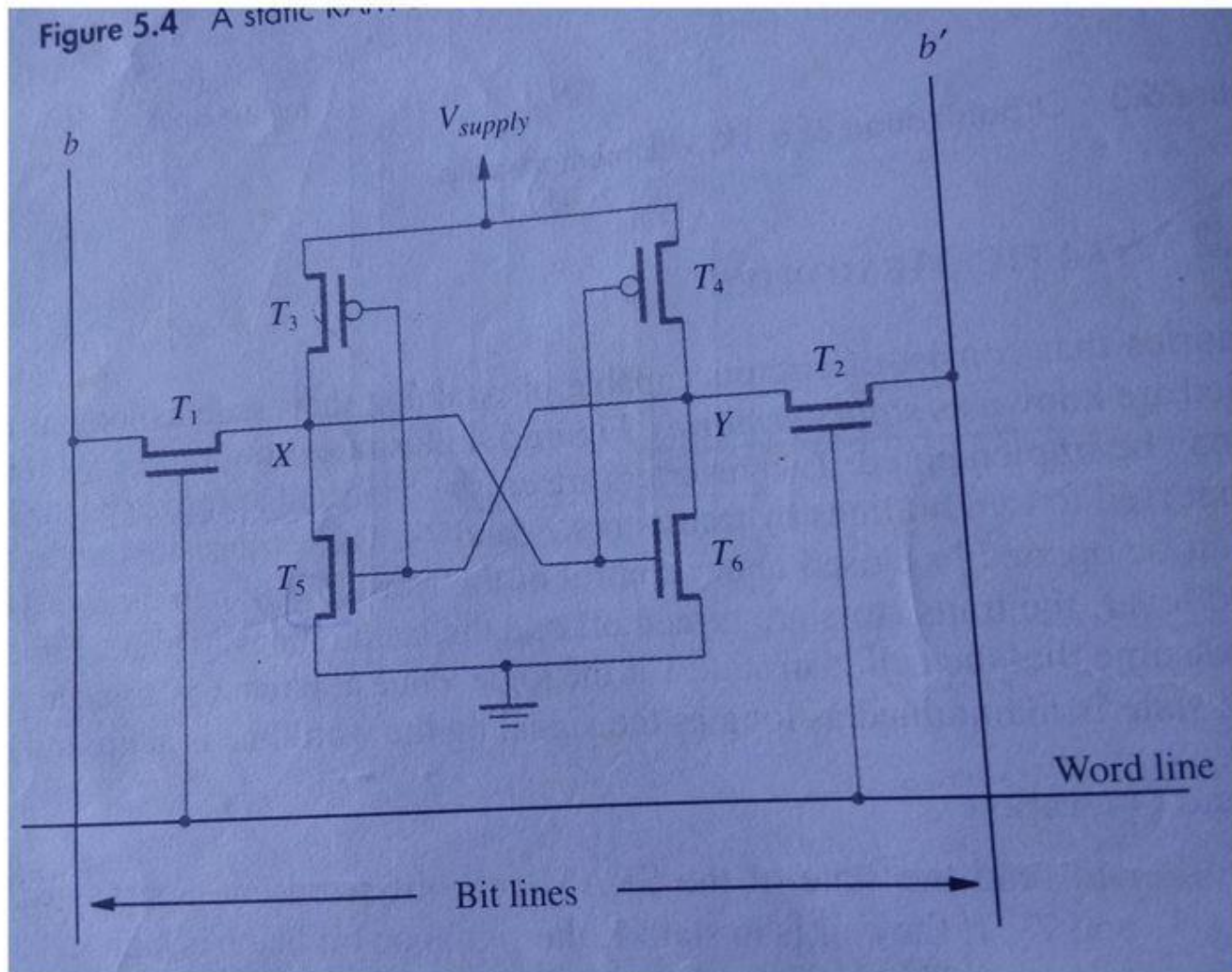
❑ Write Operation:

- ◆ The state of the cell is set by placing the appropriate value on bit line b and its complement on b' and then activating the word line.
- ◆ This forces the cell into the corresponding state. The required signal on the bit lines are generated by Sense / Write circuit.



- ❑ Transistor pairs (T3, T5) and (T4, T6) form the inverters in the latch.
- ❑ In state 1, the voltage at point X is high by having T3, T6 on and T4, T5 are OFF.
- ❑ Thus T1 and T2 returned ON (Closed), bit line b and b' will have high and low signals respectively.
- ❑ The CMOS requires 5V (in older version) or 3.3.V (in new version) of power supply voltage.
- ❑ The continuous power is needed for the cell to retain its state.

CMOS Memory Cell





☐ **Merit:**

- ◆ It has low power consumption because the current flows in the cell only when the cell is being accessed.
- ◆ Static RAMs can be accessed quickly. Its access time is few nanoseconds.

☐ **Demerit:**

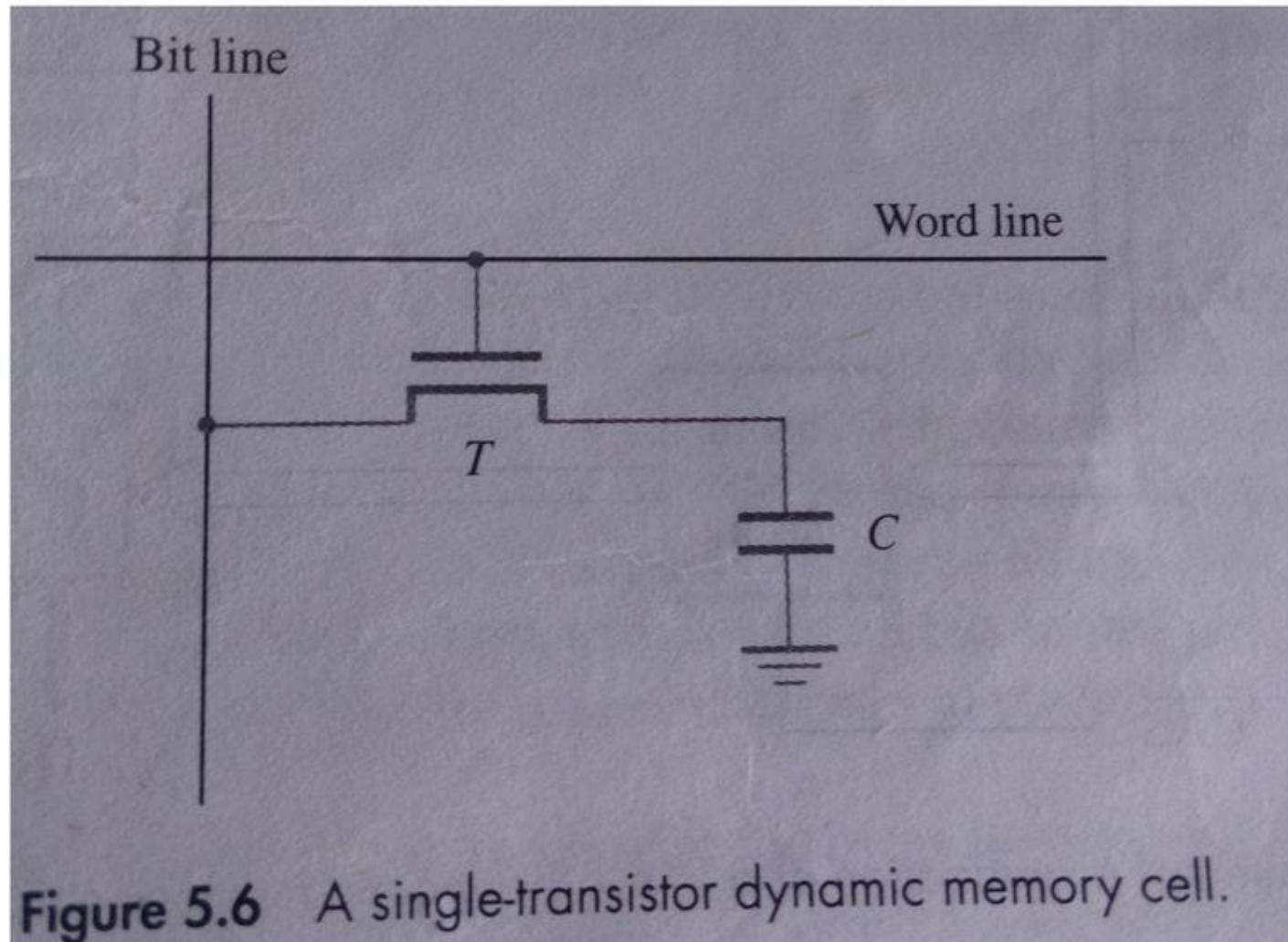
- ◆ SRAMs are said to be volatile memories because their contents are lost when the power is interrupted.





ASYNCHRONOUS DRAM

- ❑ *Less expensive RAM's* can be implemented if simpler cells are used.
- ❑ Such cells cannot retain their state indefinitely.
- ❑ Hence they are called *Dynamic RAM's* (DRAM).
- ❑ The information stored in a dynamic memory cell in the form of a charge on a capacitor and this charge can be maintained only for a few milliseconds.
- ❑ The contents must be periodically refreshed by restoring this capacitor charge to its full value.





❑ Write Operation:

- ◆ In order to store information in the cell, the transistor T is turned on and the appropriate voltage is applied to the bit line, which charges the capacitor.
- ◆ After the transistor is turned off, the capacitor begins to discharge which is caused by the capacitor's own leakage resistance.
- ◆ Hence the information stored in the cell can be retrieved correctly before the threshold value of the capacitor drops down,



❑ Read Operation

- ◆ During a read operation, the transistor is turned on and a sense amplifier connected to the bit line detects whether the charge on the capacitor is above the threshold value.

❑ A 16- megabit DRAM chip configured as 2M x 8, is shown in Figure

- ◆ If charge on capacitor $>$ threshold value \rightarrow Bit line will have logic value **1**.
- ◆ If charge on capacitor $<$ threshold value \rightarrow Bit line will set to logic value **0**.

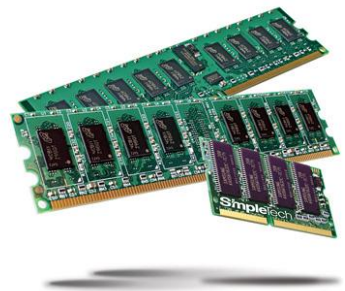
Asynchronous DRAMs

■ Static RAMs (SRAMs):

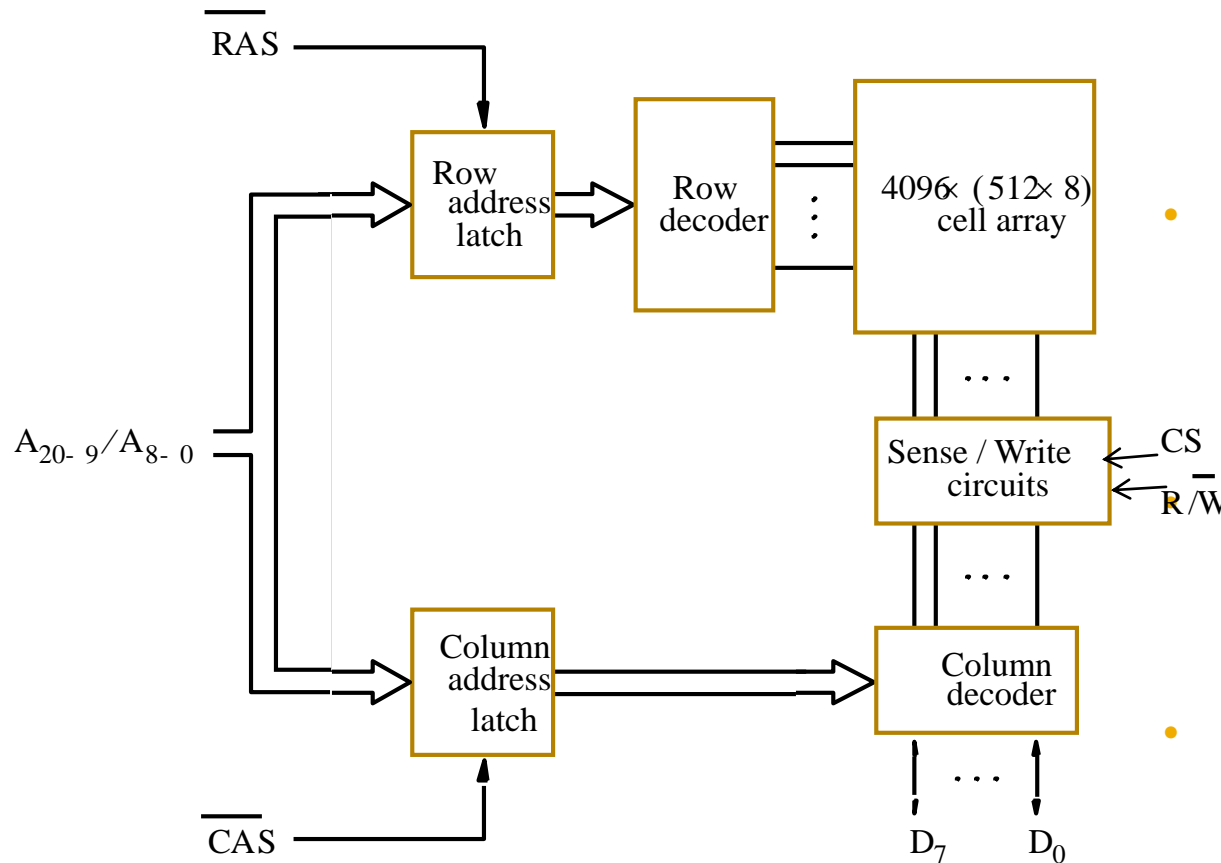
- Consist of circuits that are capable of retaining their state as long as the power is applied.
- Volatile memories, because their contents are lost when power is interrupted.
- Access times of static RAMs are in the range of few nanoseconds.
- However, the cost is usually high.

■ Dynamic RAMs (DRAMs):

- Do not retain their state indefinitely.
- Contents must be periodically refreshed.
- Contents may be refreshed while accessing them for reading.



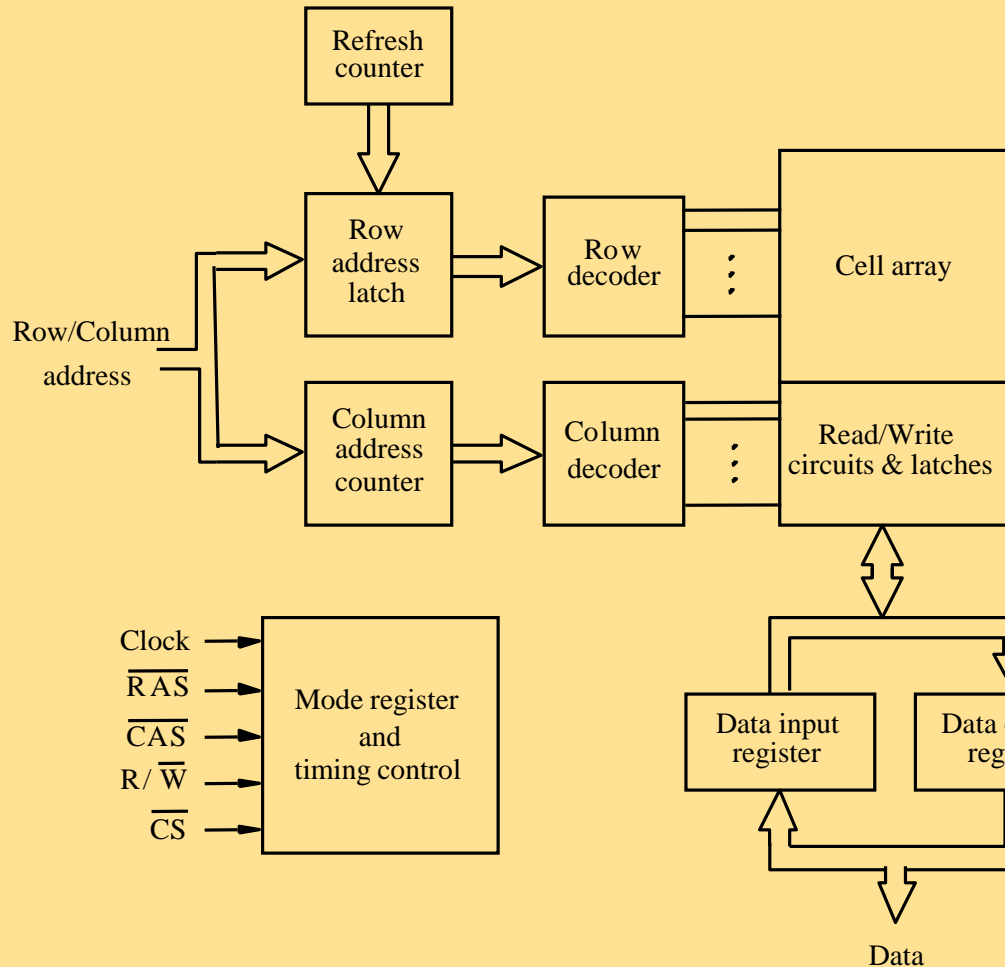
Asynchronous DRAMs



- *Each row can store 512 bytes. 12 bits to select a row, and 9 bits to select a group in a row. Total of 21 bits.*
- *First apply the row address, RAS signal latches the row address. Then apply the column address, CAS signal latches the address.*
- *Timing of the memory unit is controlled by a specialized unit which generates RAS and CAS.*
- *This is asynchronous DRAM*



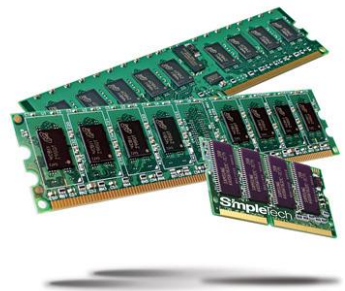
Synchronous DRAMs



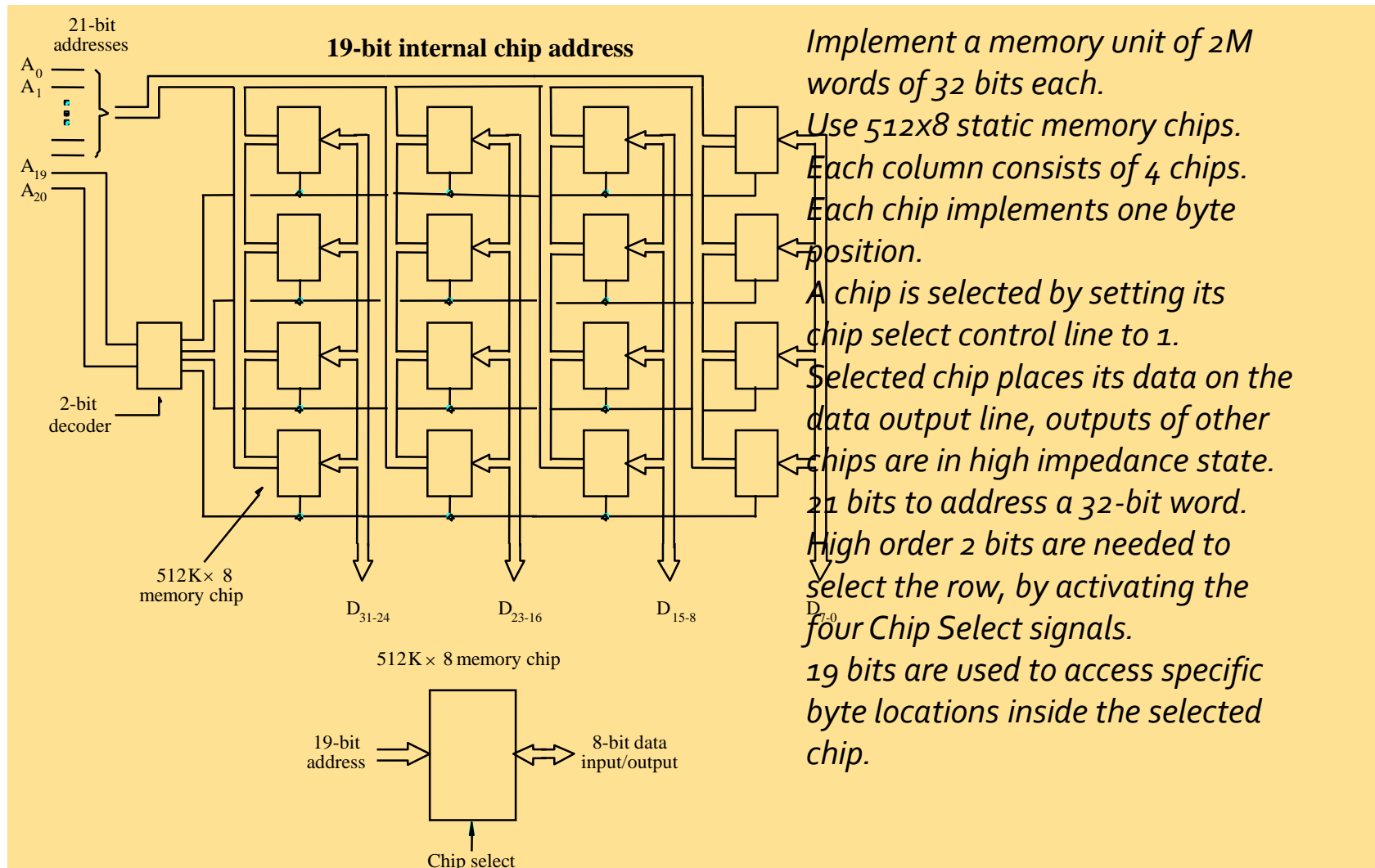
- Operation is directly synchronized with processor clock signal.
- The outputs of the sense circuits are connected to a latch.
- During a Read operation, the contents of the cells in a row are loaded onto the latches.
- During a refresh operation, the contents of the cells are refreshed without changing the contents of the latches.
- Data held in the latches correspond to the selected columns are transferred to the output.
- For a burst mode of operation, successive columns are selected using column address counter and clock. CAS signal need not be generated externally. A new data is placed during raising edge of the clock

Latency, Bandwidth, and DDRSDRAMs

- Memory latency is the time it takes to transfer a word of data to or from memory
- Memory bandwidth is the number of bits or bytes that can be transferred in one second.
- DDRSDRAMs
 - Cell array is organized in two banks

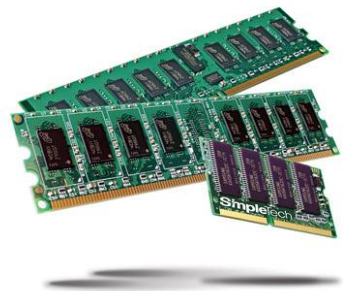


Static memories



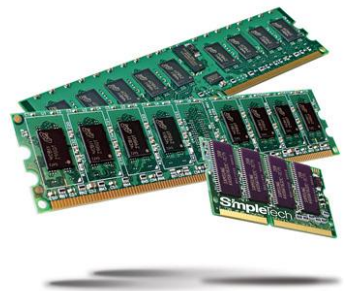
Dynamic memories

- Large dynamic mem sys can be implemented using DRAM chips like static mem sys.
- Placing large memory systems directly on the motherboard will occupy a large amount of space.
 - Also, this arrangement is inflexible since the memory system cannot be expanded easily.
- Packaging considerations have led to development of larger mem units known as SIMMs (Single In-line Memory Modules) and DIMMs (Dual In-line Memory Modules).
- Memory modules are an assembly of memory chips on a small board that plugs vertically onto a single socket on the motherboard.
 - Occupy less space on the motherboard.
 - Allows for easy expansion by replacement.

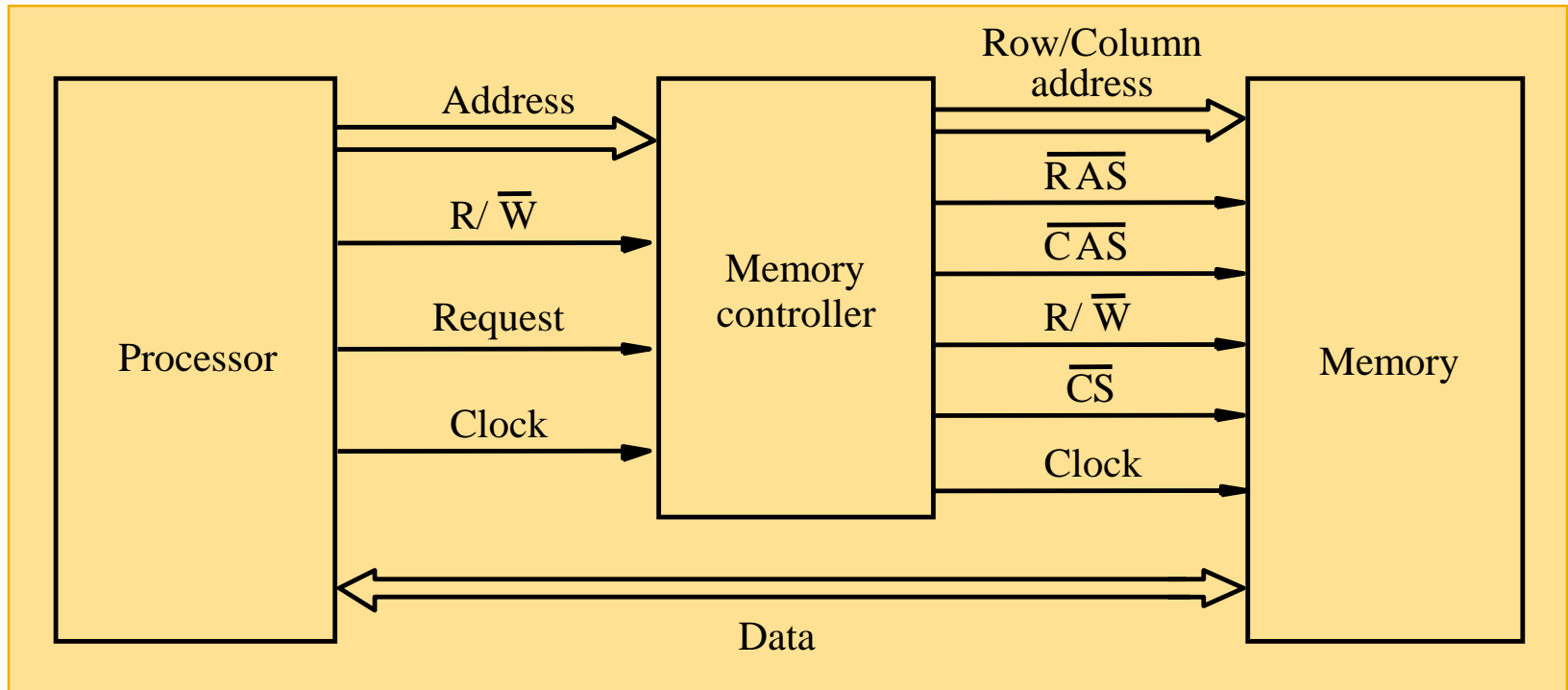


Memory controller

- Recall that in a dynamic memory chip, to reduce the number of pins, multiplexed addresses are used.
- Address is divided into two parts:
 - High-order address bits select a row in the array.
 - They are provided first, and latched using RAS signal.
 - Low-order address bits select a column in the row.
 - They are provided later, and latched using CAS signal.
- However, a **processor issues all address bits at the same time.**
- In order to **achieve the multiplexing, memory controller circuit is inserted between the processor and memory.**



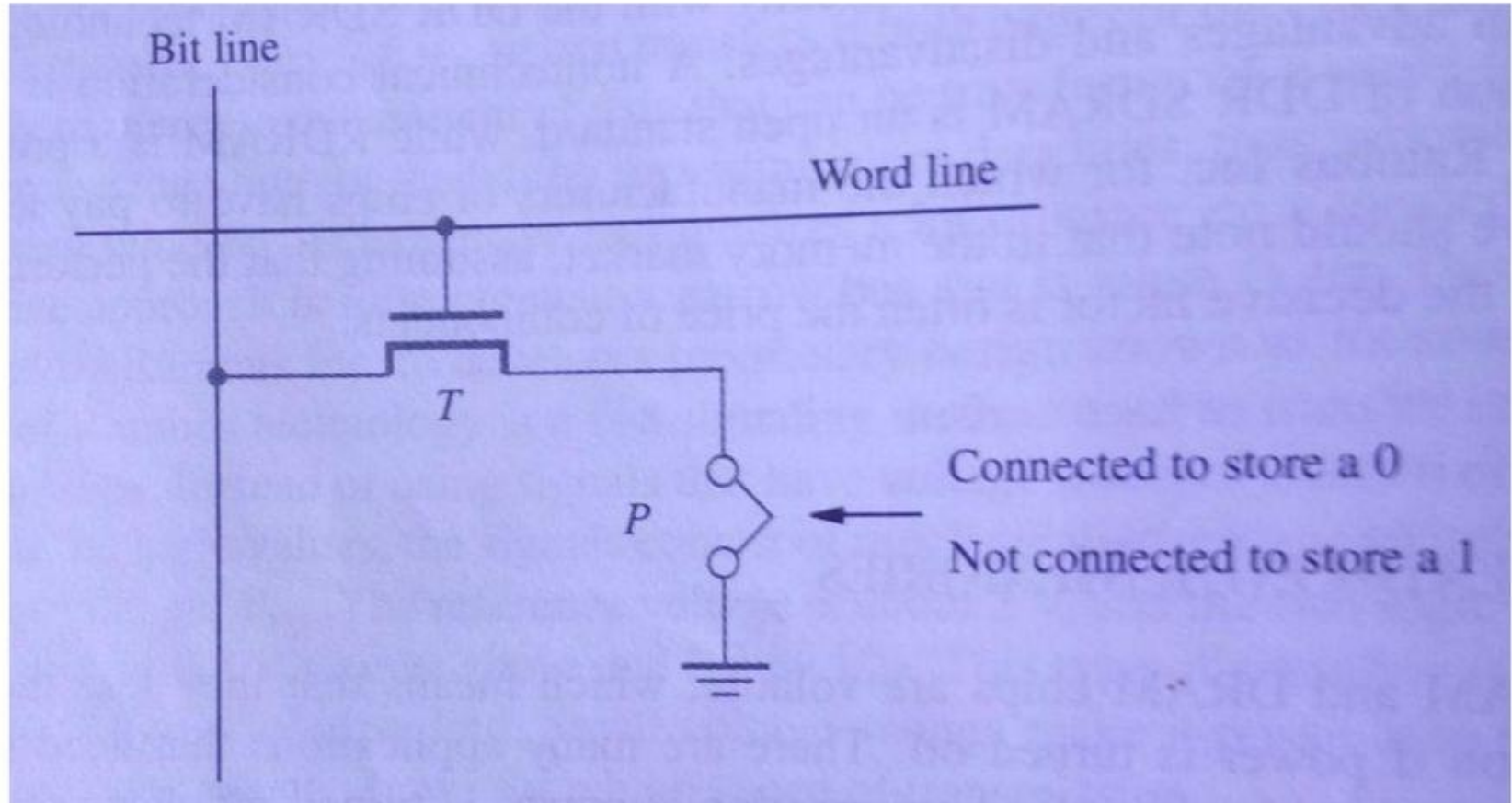
Memory controller (contd..)



Read-Only Memories (ROMs)

The Memory System

ROM CELL





- ❑ **At Logic value 0** → Transistor (T) is connected to the ground point (P).
- ❑ Transistor switch is closed and voltage on bit line nearly drops to zero.
- ❑ **At Logic value 1** → Transistor switch is open. The bit line remains at high voltage.
- ❑ To read the state of the cell, the word line is activated.
- ❑ A Sense circuit at the end of the bit line generates the proper output value.



PROM: Programmable ROM:

- ☐ PROM allows the data to be loaded by the user.
- ☐ Programmability is achieved by inserting a fuse at point P in a ROM cell.
- ☐ Before it is programmed, the memory contains all 0's.
- ☐ The user can insert 1's at the required location by burning out the fuse at these locations using high current pulse.
- ☐ This process is irreversible.
- ☐ **Merit:**
 - ◆ It provides flexibility.
 - ◆ It is faster.
 - ◆ It is less expensive because they can be programmed directly by the user.





EPROM - Erasable reprogrammable ROM:

- ❑ EPROM allows the stored data to be erased and new data to be loaded.
- ❑ In an EPROM cell, a connection to ground is always made at $_P'$ and a special transistor is used, which has the ability to function either as a normal transistor or as a disabled transistor that is always turned off.
- ❑ This transistor can be programmed to behave as a permanently open switch, by injecting charge into it that becomes trapped inside.
- ❑ Erasure requires dissipating the charges trapped in the transistor of memory cells.
- ❑ This can be done by exposing the chip to ultraviolet light, so that EPROM chips are mounted in packages that have transparent windows.



□ Merits:

- ◆ It provides flexibility during the development phase of digital system.
- ◆ It is capable of retaining the stored information for a long time.

□ Demerits:

- ◆ The chip must be physically removed from the circuit for reprogramming and its entire contents are erased by UV light.

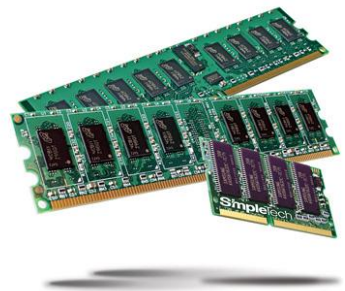


EEPROM: -Electrically Erasable ROM:

- ❑ **EEPROM** (also written **E2PROM** and pronounced "e-e-prom," "double-e prom," "e-squared," or simply "e-prom") stands for **E**lectrically **E**rasable **P**rogrammable **R**ead- **O**nly **M**emory and is a type of non-volatile memory used in computers and other electronic devices to store small amounts of data that must be saved when power is removed, e.g., calibration tables or device configuration

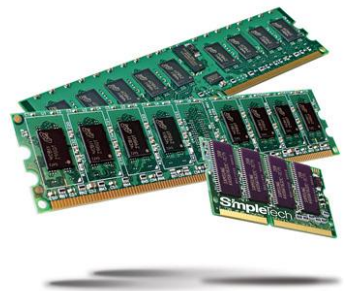
Read-Only Memories (ROMs)

- SRAM and SDRAM chips are volatile:
 - Lose the contents when the power is turned off.
- Many applications need memory devices to retain contents after the power is turned off.
 - For example, computer is turned on, the operating system must be loaded from the disk into the memory.
 - Store instructions which would load the OS from the disk.
 - Need to store these instructions so that they will not be lost after the power is turned off.
 - We need to store the instructions into a non-volatile memory.
- Non-volatile memory is read in the same manner as volatile memory.
 - Separate writing process is needed to place information in this memory.
 - Normal operation involves only reading of data, this type of memory is called Read-Only memory (ROM).



Read-Only Memories (Contd.,)

- Read-Only Memory:
 - Data are written into a ROM when it is manufactured.
- Programmable Read-Only Memory (PROM):
 - Allow the data to be loaded by a user.
 - Process of inserting the data is irreversible.
 - Storing information specific to a user in a ROM is expensive.
 - Providing programming capability to a user may be better.
- Erasable Programmable Read-Only Memory (EPROM):
 - Stored data to be erased and new data to be loaded.
 - Flexibility, useful during the development phase of digital systems.
 - Erasable, reprogrammable ROM.
 - Erasure requires exposing the ROM to UV light.



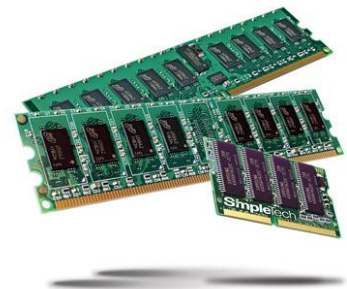
Read-Only Memories (Contd.,)

■ Electrically Erasable Programmable Read-Only Memory (EEPROM):

- To erase the contents of EPROMs, they have to be exposed to ultraviolet light.
- Physically removed from the circuit.
- EEPROMs the contents can be stored and erased electrically.

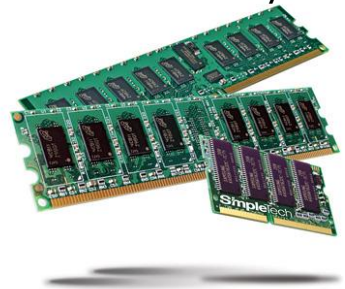
■ Flash memory:

- Has similar approach to EEPROM.
- Read the contents of a single cell, but write the contents of an entire block of cells.
- Flash devices have greater density.
 - Higher capacity and low storage cost per bit.
- Power consumption of flash memory is very low, making it attractive for use in equipment that is battery-driven.
- Single flash chips are not sufficiently large, so larger memory modules are implemented using flash cards and flash drives.

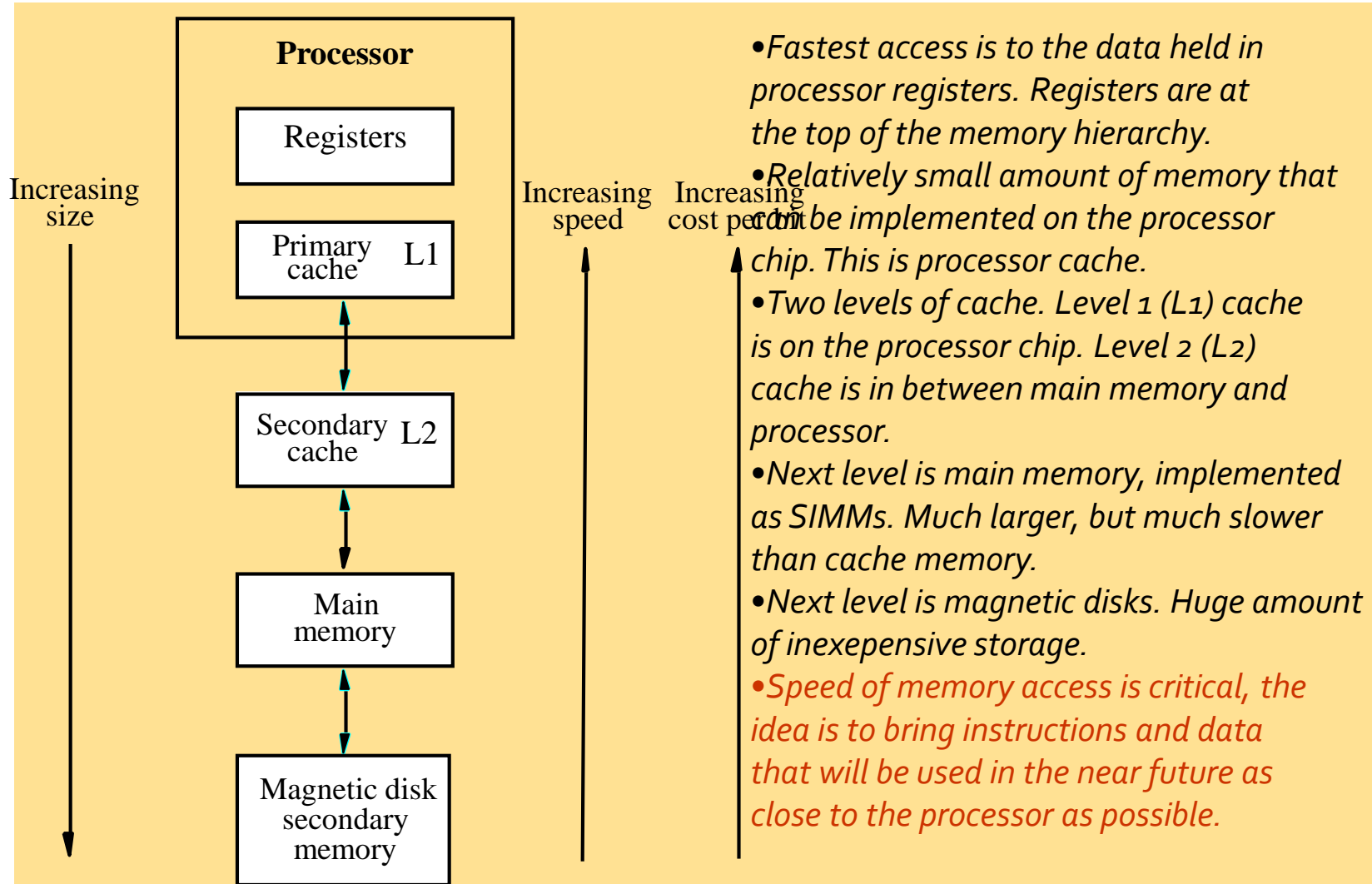


Speed, Size, and Cost

- A big challenge in the design of a computer system is to provide a sufficiently large memory, with a reasonable speed at an affordable cost.
- Static RAM:
 - Very fast, but expensive, because a basic SRAM cell has a complex circuit making it impossible to pack a large number of cells onto a single chip.
- Dynamic RAM:
 - Simpler basic cell circuit, hence are much less expensive, but significantly slower than SRAMs.
- Magnetic disks:
 - Storage provided by DRAMs is higher than SRAMs, but is still less than what is necessary.
 - Secondary storage such as magnetic disks provide a large amount of storage, but is much slower than DRAMs.



Memory Hierarchy

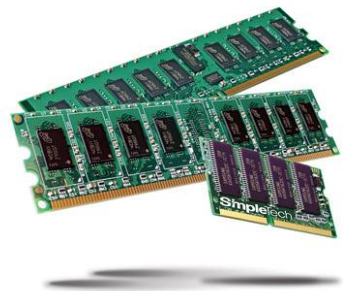


Cache Memories

The Memory System

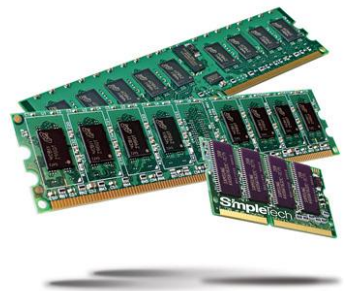
Cache Memories

- Processor is much faster than the main memory.
 - As a result, the processor has to spend much of its time waiting while instructions and data are being fetched from the main memory.
 - Major obstacle towards achieving good performance.
- Speed of the main memory cannot be increased beyond a certain point.
- Cache memory is an architectural arrangement which makes the main memory appear faster to the processor than it really is.
- Cache memory is based on the property of computer programs known as "locality of reference".

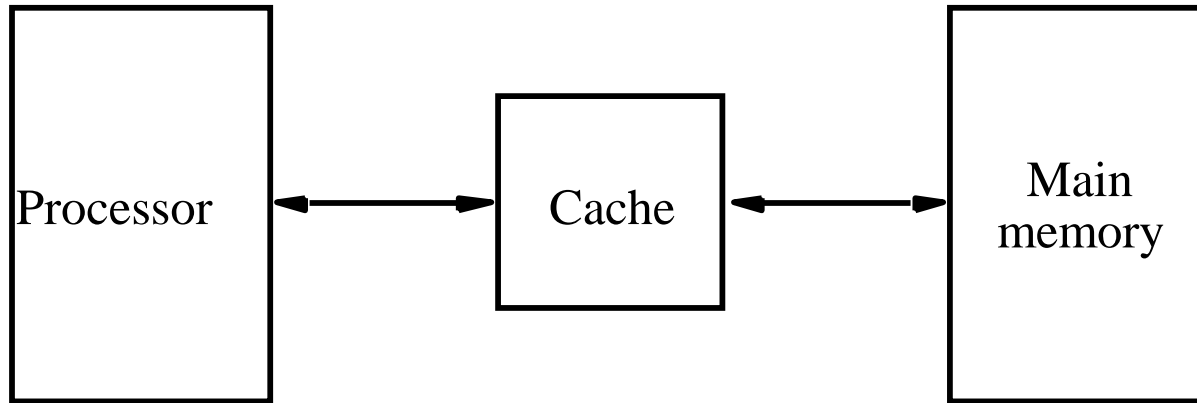


Locality of Reference

- Analysis of programs indicates that many instructions in localized areas of a program are executed repeatedly during some period of time, while the others are accessed relatively less frequently.
 - These instructions may be the ones in a loop, nested loop or few procedures calling each other repeatedly.
 - This is called "locality of reference".
- Temporal locality of reference:
 - Recently executed instruction is likely to be executed again very soon.
- Spatial locality of reference:
 - Instructions with addresses close to a recently instruction are likely to be executed soon.



Cache memories

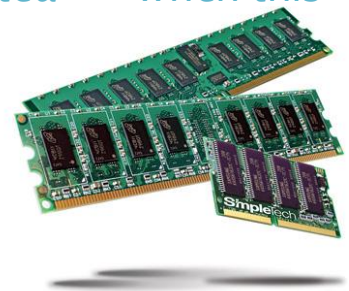


- Processor issues a Read request, a block of words is transferred from the main memory to the cache, one word at a time.
- Subsequent references to the data in this block of words are found in the cache.
- At any given time, only some blocks in the main memory are held in the cache. Which blocks in the main memory are in the cache is determined by a "mapping function".
- When the cache is full, and a block of words needs to be transferred from the main memory, some block of words in the cache must be replaced. This is determined by a "replacement algorithm".



Cache hit

- Existence of a cache is transparent to the processor. The processor issues Read and Write requests in the same manner.
- If the data is in the cache it is called a Read or Write hit.
- Read hit:
 - The data is obtained from the cache.
- Write hit:
 - Cache has a replica of the contents of the main memory.
 - Contents of the cache and the main memory may be updated simultaneously. This is the write-through protocol.
 - Update the contents of the cache, and mark it as updated by setting a bit known as the dirty bit or modified bit. The contents of the main memory are updated when this block is replaced. This is write-back or copy-back protocol.



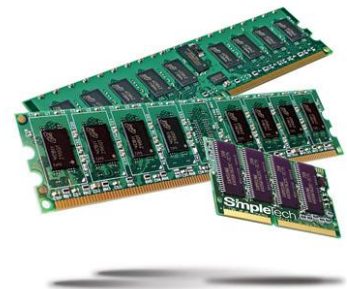
Cache miss

- *If the data is not present in the cache, then a Read miss or Write miss occurs.*
- *Read miss:*
 - *Block of words containing this requested word is transferred from the memory.*
 - *After the block is transferred, the desired word is forwarded to the processor.*
 - *The desired word may also be forwarded to the processor as soon as it is transferred without waiting for the entire block to be transferred. This is called load-through or early-restart.*
- *Write-miss:*
 - *Write-through protocol is used, then the contents of the main memory are updated directly.*
 - *If write-back protocol is used, the block containing the addressed word is first brought into the cache. The desired word is overwritten with new information.*



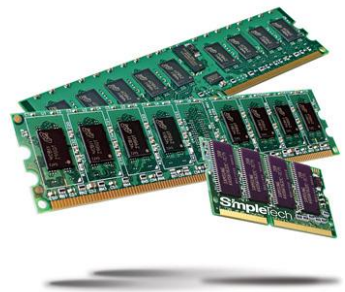
Cache Coherence Problem

- A bit called as "valid bit" is provided for each block.
- If the block contains valid data, then the bit is set to 1, else it is 0.
- Valid bits are set to 0, when the power is just turned on.
- When a block is loaded into the cache for the first time, the valid bit is set to 1.
- Data transfers between main memory and disk occur directly bypassing the cache.
- When the data on a disk changes, the main memory block is also updated.
- However, if the data is also resident in the cache, then the valid bit is set to 0.
- What happens if the data in the disk and main memory changes and the write-back protocol is being used?
- In this case, the data in the cache may also have changed and is indicated by the dirty bit.
- The copies of the data in the cache, and the main memory are different. This is called the cache coherence problem.
- One option is to force a write-back before the main memory is updated from the disk.

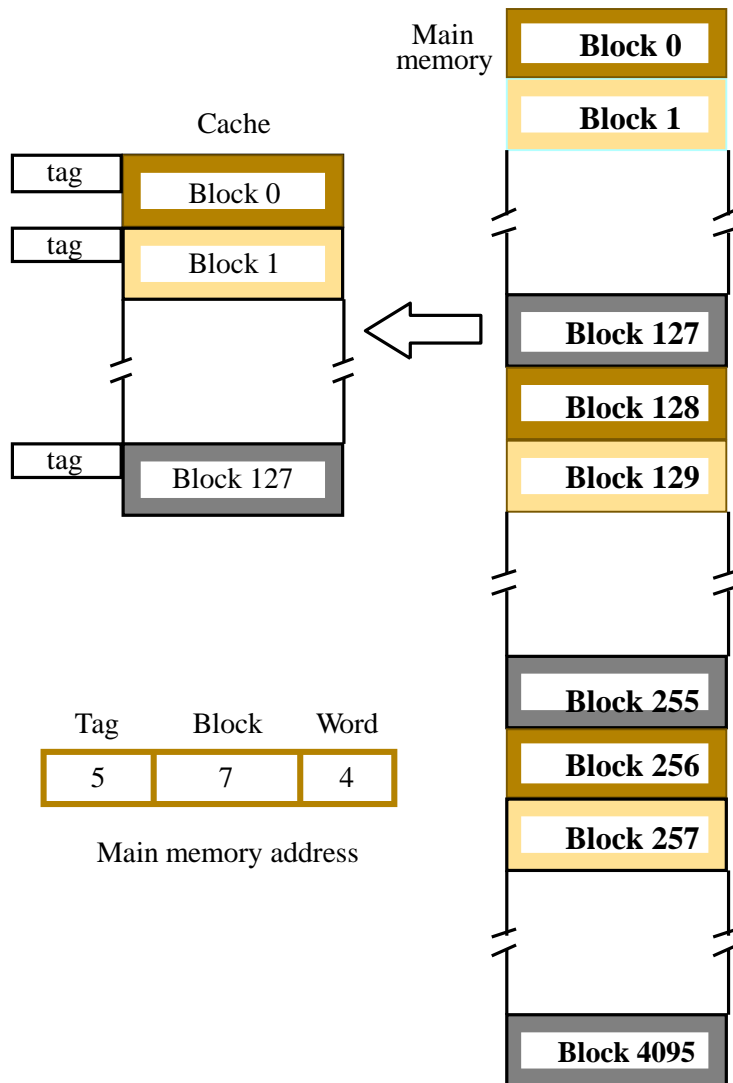


Mapping functions

- Mapping functions determine how memory blocks are placed in the cache.
- A simple processor example:
 - Cache consisting of 128 blocks of 16 words each.
 - Total size of cache is 2048 (2K) words.
 - Main memory is addressable by a 16-bit address.
 - Main memory has 64K words.
 - Main memory has 4K blocks of 16 words each.
- Three mapping functions:
 - Direct mapping
 - Associative mapping
 - Set-associative mapping.

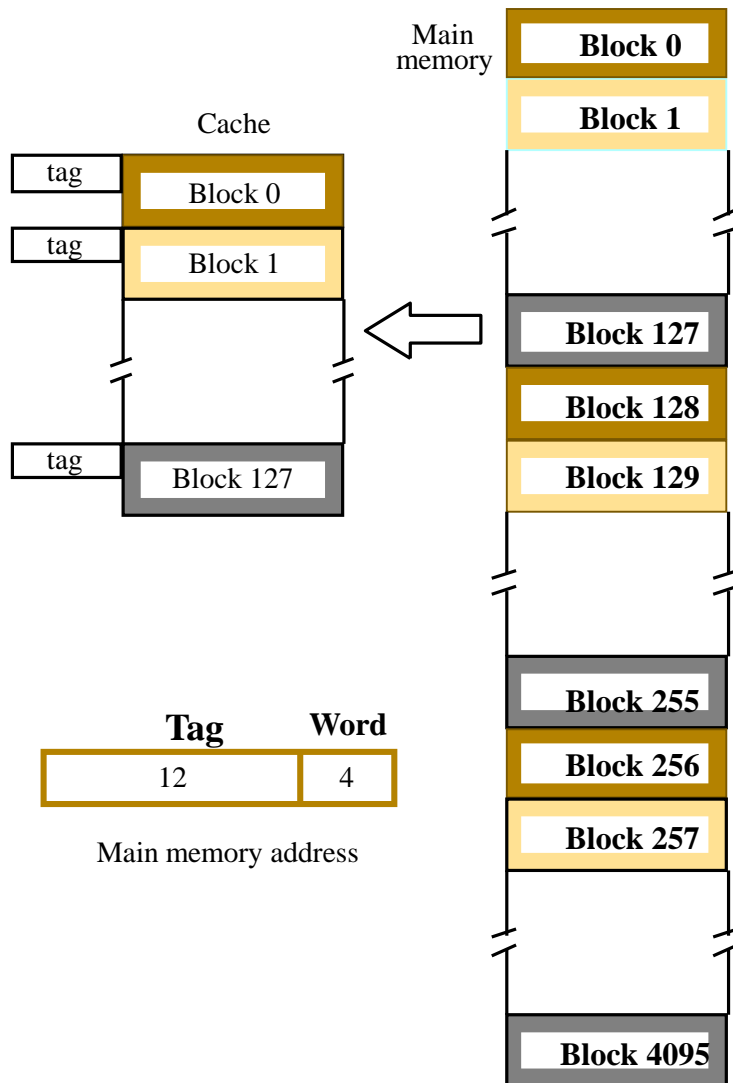


Direct mapping



- Block j of the main memory maps to j modulo 128 of the cache. 0 maps to 0, 129 maps to 1.
- More than one memory block is mapped onto the same position in the cache.
- May lead to contention for cache blocks even if the cache is not full.
- Resolve the contention by allowing new block to replace the old block, leading to a trivial replacement algorithm.
- Memory address is divided into three fields:
 - Low order 4 bits determine one of the 16 words in a block.
 - When a new block is brought into the cache, the the next 7 bits determine which cache block this new block is placed in.
 - High order 5 bits determine which of the possible 32 blocks is currently present in the cache. These are tag bits.
- Simple to implement but not very flexible.

Associative mapping



- Main memory block can be placed into any cache position.

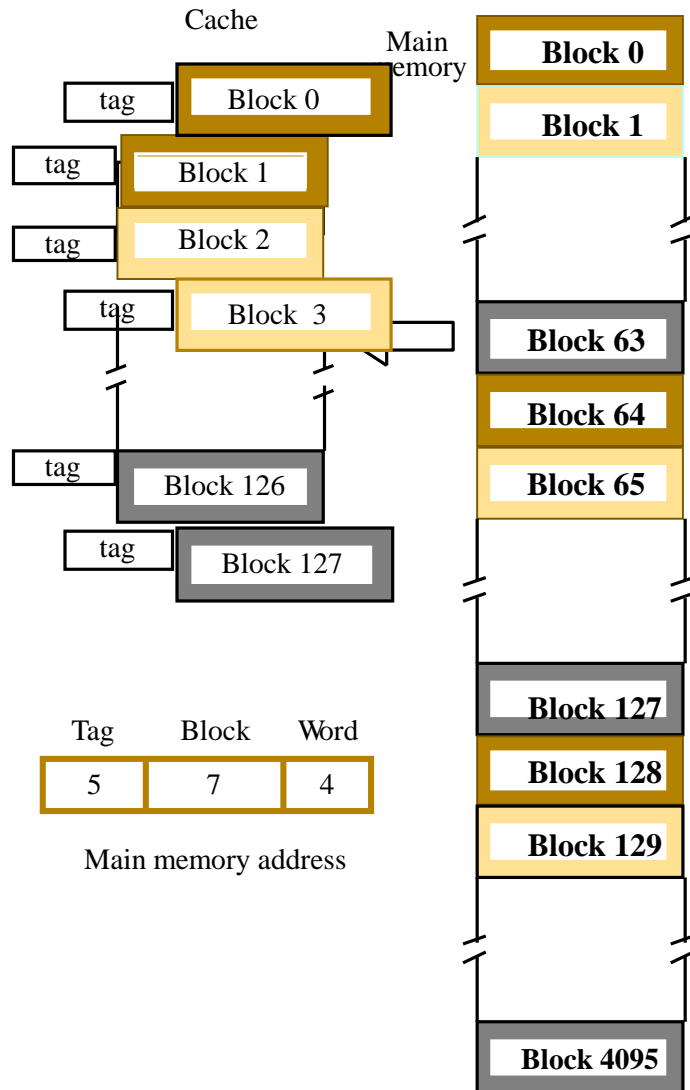
- Memory address is divided into two fields:
 - Low order 4 bits identify the word within a block.
 - High order 12 bits or tag bits identify a memory block when it is resident in the cache.

- Flexible, and uses cache space efficiently.

- Replacement algorithms can be used to replace an existing block in the cache when the cache is full.

- Cost is higher than direct-mapped cache because of the need to search all 128 patterns to determine whether a given block is in the cache.

Set-Associative mapping



Blocks of cache are grouped into sets.

Mapping function allows a block of the main memory to reside in any block of a specific set.

Divide the cache into 64 sets, with two blocks per set. Memory block 0, 64, 128 etc. map to block 0, and they can occupy either of the two positions.

Memory address is divided into three fields:

- 6 bit field determines the set number.*
- High order 6 bit fields are compared to the tag fields of the two blocks in a set.*

Set-associative mapping combination of direct and associative mapping.

Number of blocks per set is a design parameter.

- One extreme is to have all the blocks in one set, requiring no set bits (fully associative mapping).*
- Other extreme is to have one block per set, is the same as direct mapping.*