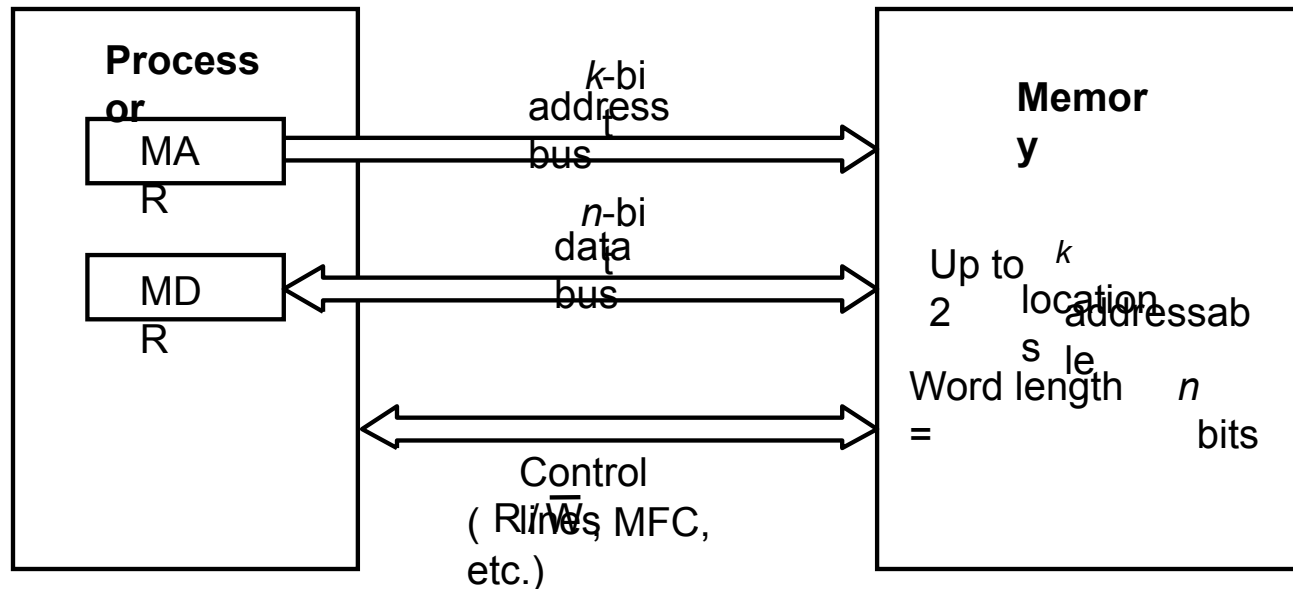


Fundamental Concepts

The Memory System

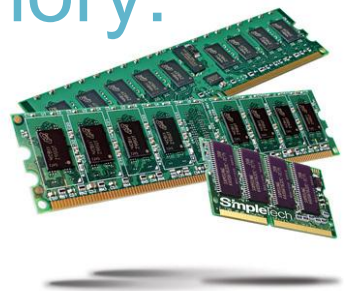
Some basic concepts

- 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).

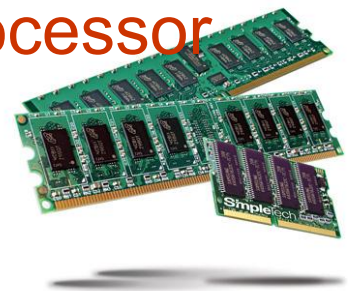


Semiconductor RAM memories

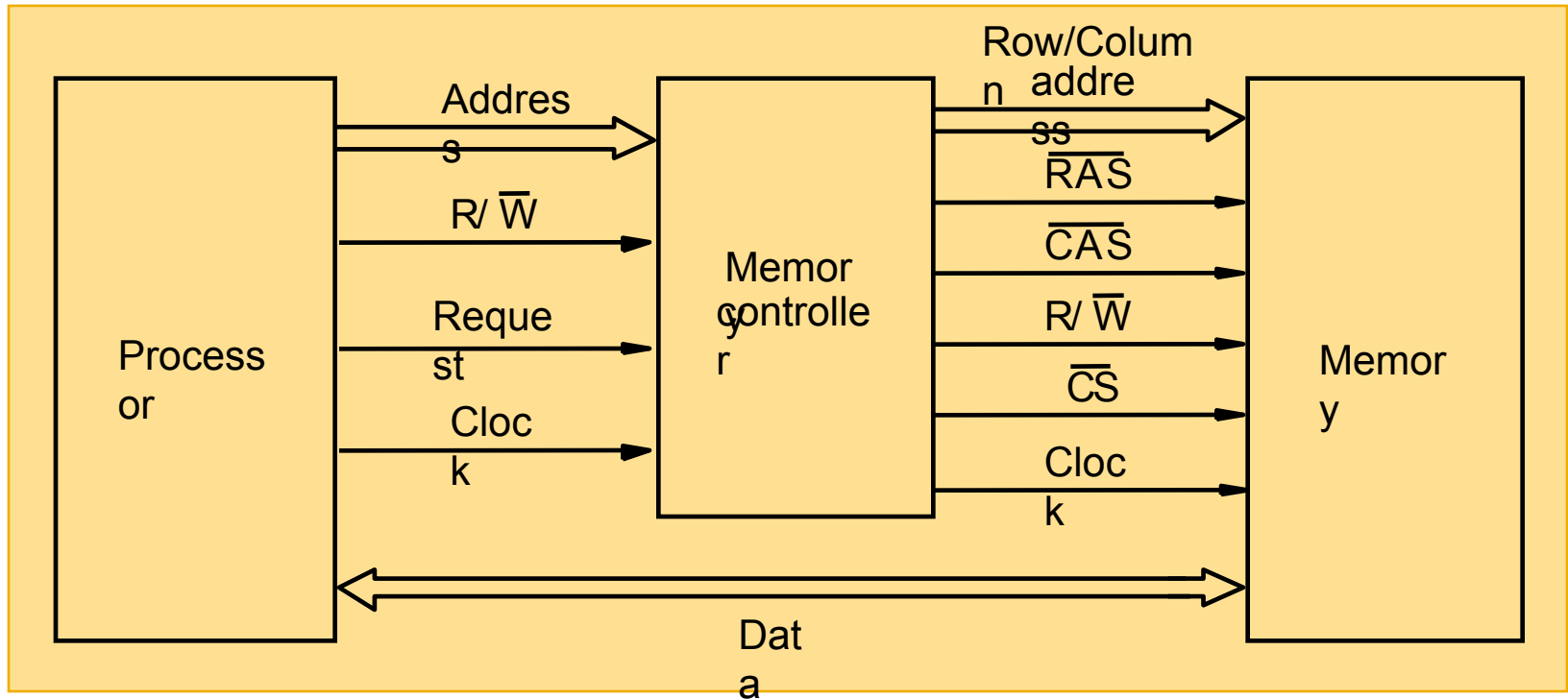
The Memory System

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

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:

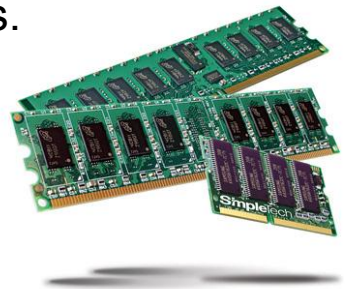
- Data are written into a ROM when it is manufactured.

■ 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.

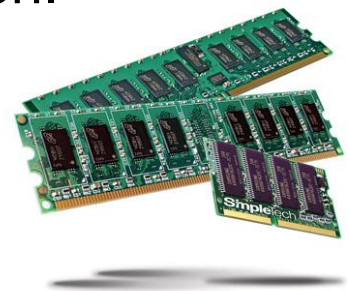
■ Erasable Programmable Read-Only Memory (EPROM):

- 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.



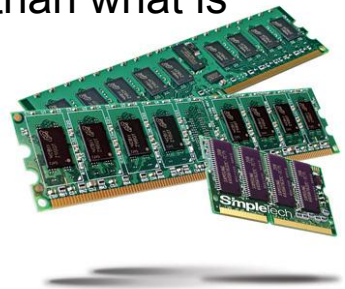
Read-Only Memories (Contd.,)

- 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.

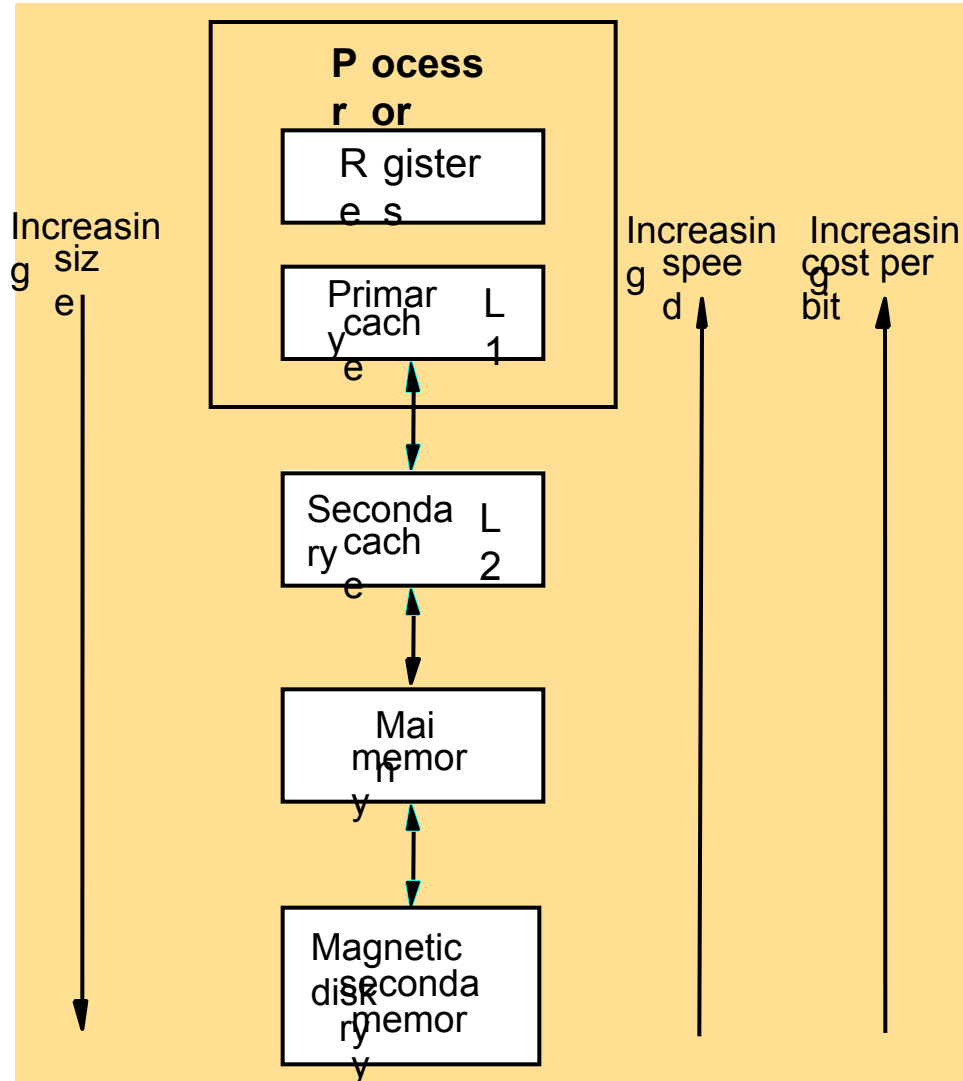


Speed, Size, and Cost

- 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



- Fastest access is to the data held in processor registers. Registers are at the top of the memory hierarchy.
- Relatively small amount of memory that can be implemented on the processor chip. This is processor cache.
- Two levels of cache. Level 1 (L1) cache is on the processor chip. Level 2 (L2) cache is in between main memory and processor.
- Next level is main memory, implemented as SIMMs. Much larger, but much slower than cache memory.
- Next level is magnetic disks. Huge amount of inexpensive storage.
- Speed of memory access is critical, the

Cache Memories

The Memory System

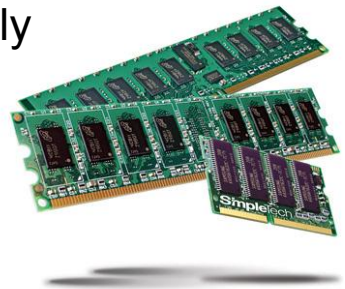
Cache Memories

- 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”.

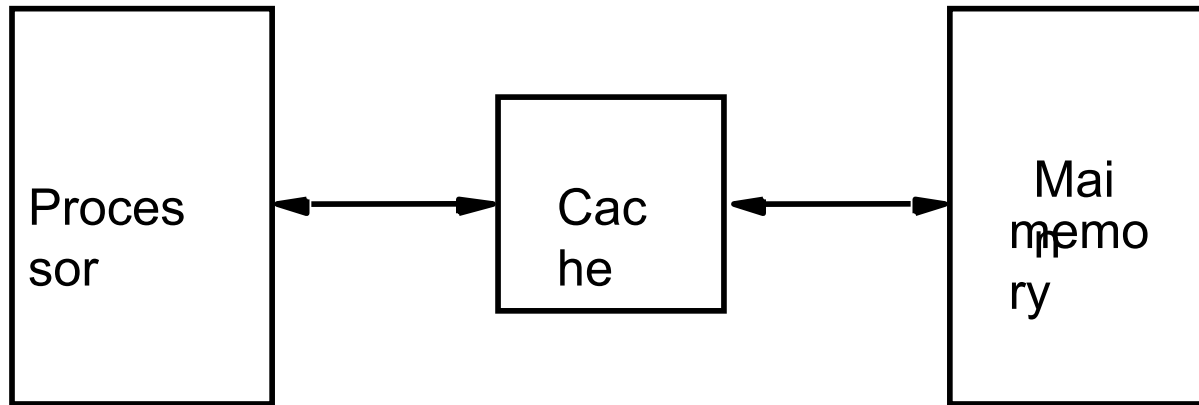


Locality of Reference

- 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”.



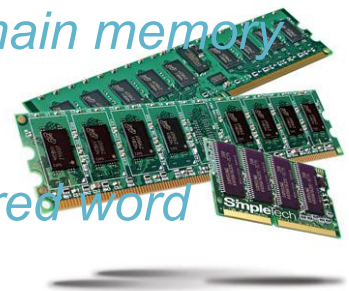
Cache hit

- 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.*



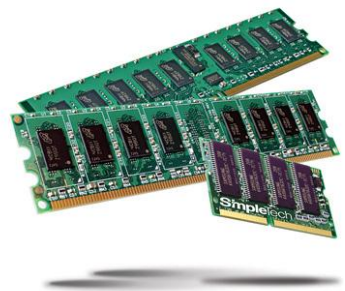
Cache Coherence Problem

- 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.

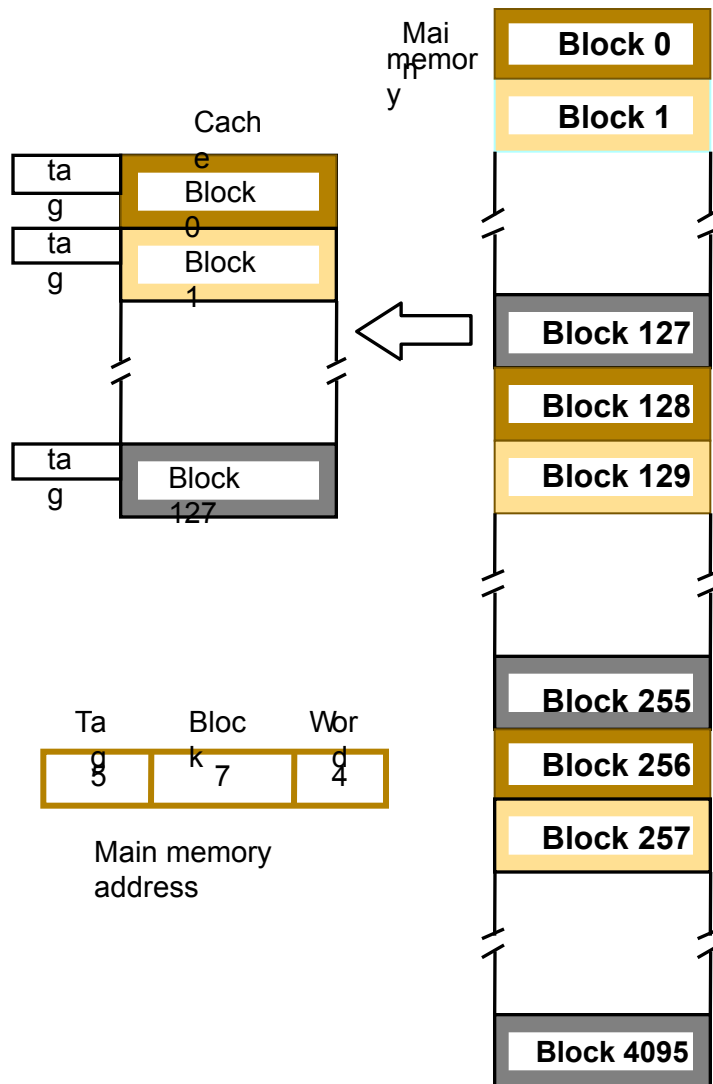


Mapping functions

- 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.

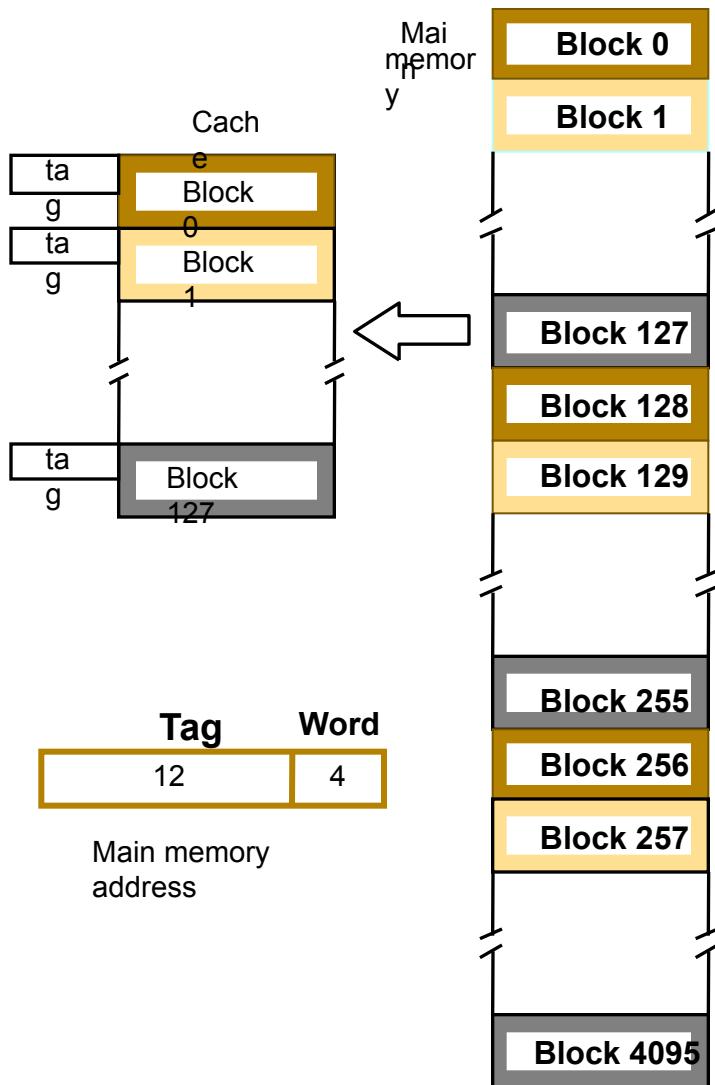


Direct mapping



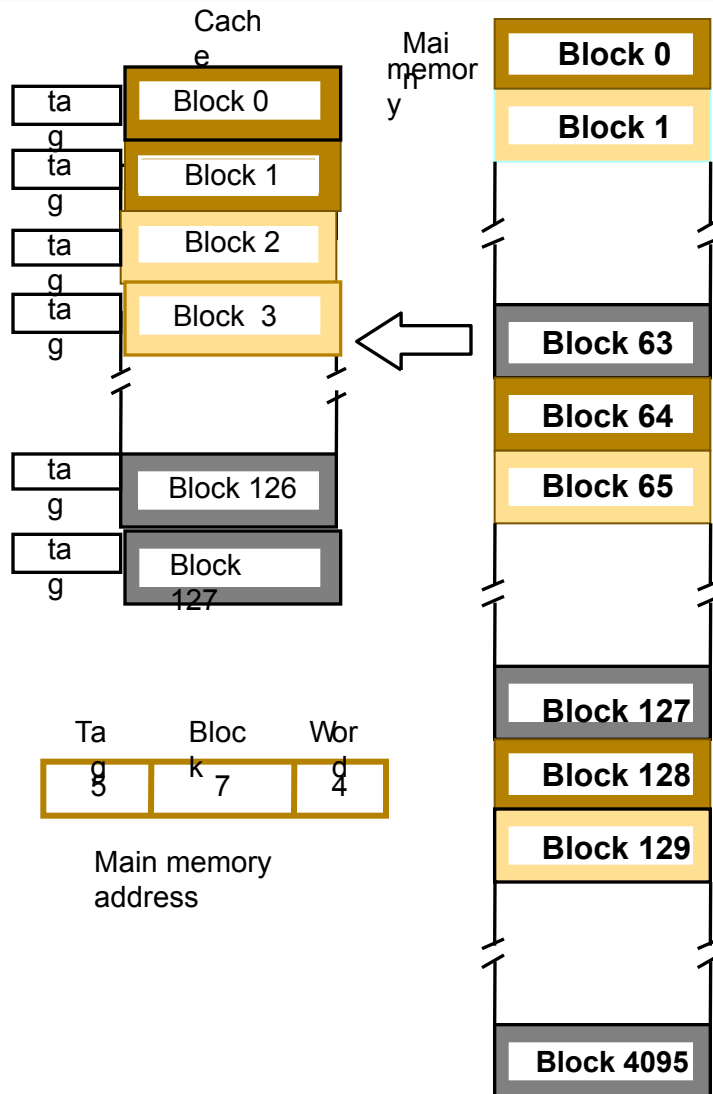
- Block j of the main memory maps to j modulo 128 of the cache. 0 maps to 0, 129 maps to 1.
- More than one memory block is mapped onto the same position in the cache.
- May lead to contention for cache blocks even if the cache is not full.
- Resolve the contention by allowing new block to replace the 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.

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).

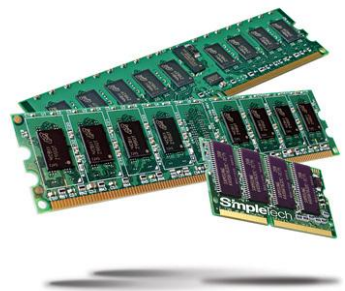
- Other extreme is to have one block per set

Performance considerations

The Memory System

Performance considerations

- A key design objective of a computer system is to achieve the best possible performance at the lowest possible cost.
 - Price/performance ratio is a common measure of success.
- Performance of a processor depends on:
 - How fast machine instructions can be brought into the processor for execution.
 - How fast the instructions can be executed.

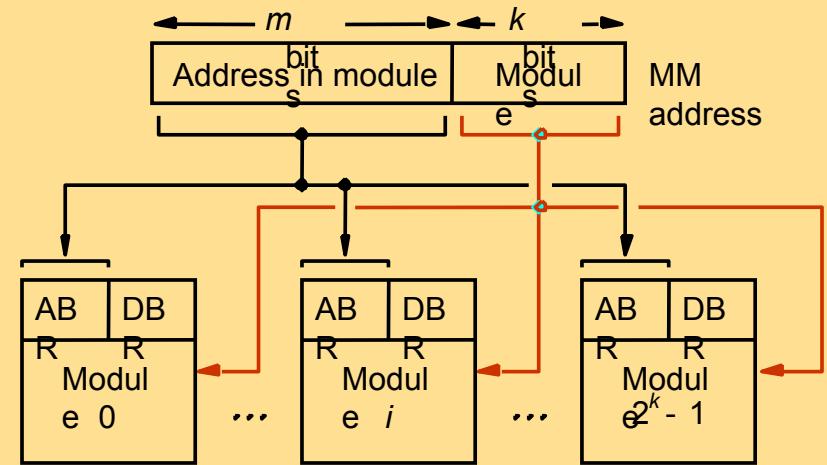
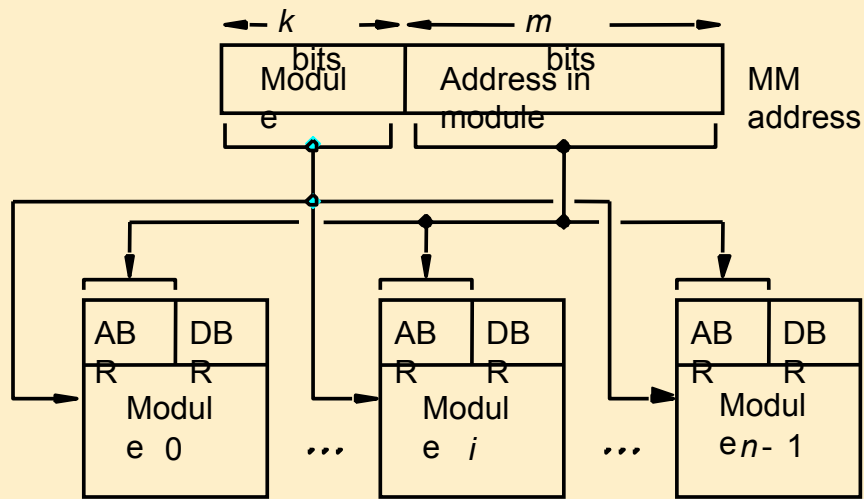


Interleaving

- Divides the memory system into a number of **memory modules**. Each module has its own address buffer register (ABR) and data buffer register (DBR).
- Arranges addressing so that successive words in the address space are placed in different modules.
- When requests for memory access involve consecutive addresses, the access will be to different modules.
- Since parallel access to these modules is possible, the average rate of fetching words from the Main Memory can be increased.



Methods of address layouts

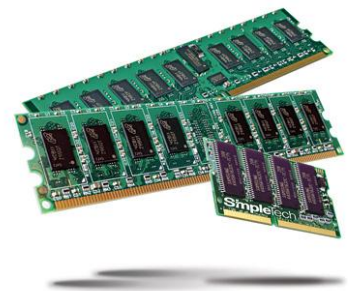


- Consecutive words are placed in a module.
- High-order k bits of a memory address determine the module.
- Low-order m bits of a memory address determine the word within a module.
- When a block of words is transferred from main memory to cache, only one module is busy at a time.
- Consecutive words are located in consecutive modules.
- Consecutive addresses can be located in consecutive modules.
- While transferring a block of data, several memory modules can be kept busy at the same time.



Hit Rate and Miss Penalty

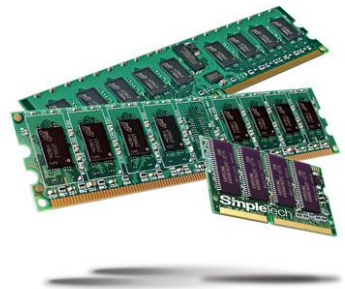
- Hit rate
- Miss penalty
- Hit rate can be improved by increasing block size, while keeping cache size constant
- Block sizes that are neither very small nor very large give best results.
- Miss penalty can be reduced if load-through approach is used when loading new blocks into cache.



Caches on the processor chip

- In high performance processors 2 levels of caches are normally used.
- Avg access time in a system with 2 levels of caches is

$$T_{ave} = h_1c_1 + (1-h_1)h_2c_2 + (1-h_1)(1-h_2)M$$



Other Performance Enhancements

Write buffer

■ Write-through:

- Each write operation involves writing to the main memory.
- If the processor has to wait for the write operation to be complete, it slows down the processor.
- Processor does not depend on the results of the write operation.
- Write buffer can be included for temporary storage of write requests.
- Processor places each write request into the buffer and continues execution.
- If a subsequent Read request references data which is still in the write buffer, then this data is referenced in the write buffer.

■ Write-back:

- Block is written back to the main memory when it is replaced.
- If the processor waits for this write to complete, before reading the new block, it is slowed down.
- Fast write buffer can hold the block to be written, and the new block can be read first.



Other Performance Enhancements (Contd.,)

Prefetching

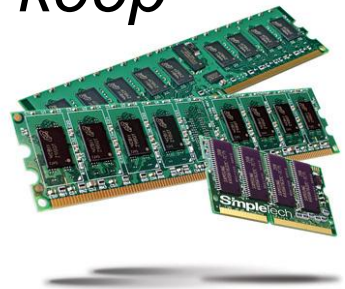
- *New data are brought into the processor when they are first needed.*
- *Processor has to wait before the data transfer is complete.*
- *Prefetch the data into the cache before they are actually needed, or a before a Read miss occurs.*
- *Prefetching can be accomplished through software by including a special instruction in the machine language of the processor.*
 - ***Inclusion of prefetch instructions increases the length of the programs.***
- *Prefetching can also be accomplished using hardware:*
 - ***Circuitry that attempts to discover patterns in memory references and then prefetches according to this pattern.***



Other Performance Enhancements (Contd.,)

Lockup-Free Cache

- *Prefetching scheme does not work if it stops other accesses to the cache until the prefetch is completed.*
- *A cache of this type is said to be “locked” while it services a miss.*
- *Cache structure which supports multiple outstanding misses is called a lockup free cache.*
- *Since only one miss can be serviced at a time, a lockup free cache must include circuits that keep track of all the outstanding misses.*
- *Special registers may hold the necessary information about these misses.*



Virtual Memory

The Memory System

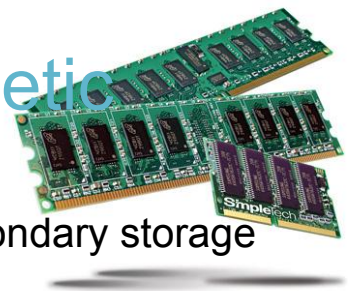
Virtual memories

- Recall that an important challenge in the design of a computer system is to provide a large, fast memory system at an affordable cost.
- Architectural solutions to increase the effective speed and size of the memory system.
- Cache memories were developed to increase the effective speed of the memory system.
- Virtual memory is an architectural solution to increase the effective size of the memory system.



Virtual memories (contd..)

- Recall that the addressable memory space depends on the number of address bits in a computer.
 - For example, if a computer issues 32-bit addresses, the addressable memory space is 4G bytes.
- Physical main memory in a computer is generally not as large as the entire possible addressable space.
 - Physical memory typically ranges from a few hundred megabytes to 1G bytes.
- Large programs that cannot fit completely into the main memory have their parts stored on secondary storage devices such as magnetic disks.
 - Pieces of programs must be transferred to the main memory from secondary storage before they can be executed.



Virtual memories (contd..)

- When a new piece of a program is to be transferred to the main memory, and the main memory is full, then some other piece in the main memory must be replaced.
 - Recall this is very similar to what we studied in case of cache memories.
- Operating system automatically transfers data between the main memory and secondary storage.
 - Application programmer need not be concerned with this transfer.
 - Also, application programmer does not need to be aware of the limitations imposed by the available physical memory.

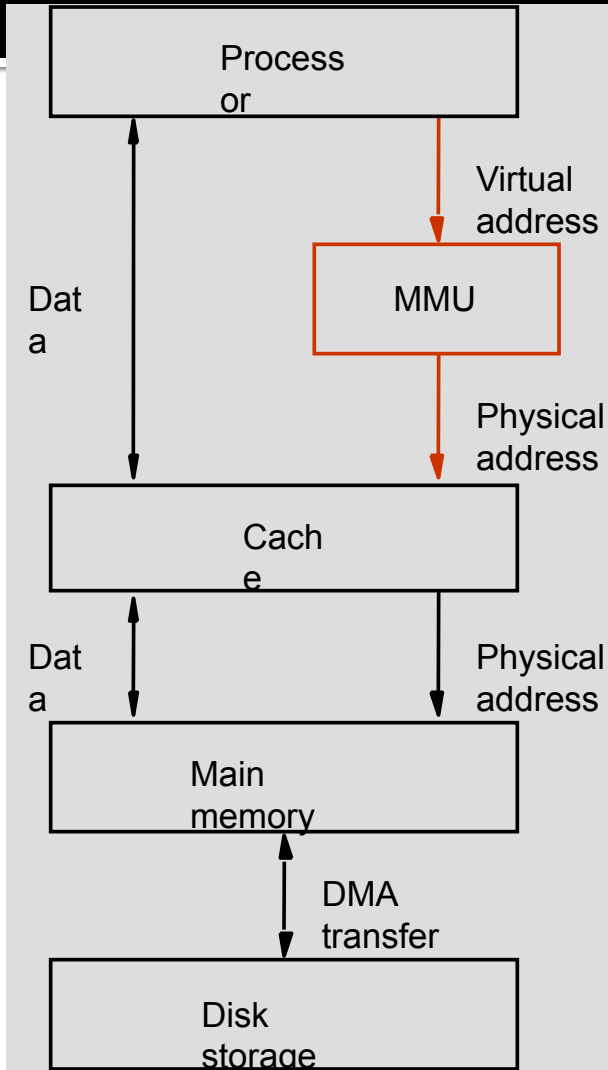


Virtual memories (contd..)

- Techniques that automatically move program and data between main memory and secondary storage when they are required for execution are called virtual-memory techniques.
- Programs and processors reference an instruction or data independent of the size of the main memory.
- Processor issues binary addresses for instructions and data.
 - These binary addresses are called logical or virtual addresses.
- Virtual addresses are translated into physical addresses by a combination of hardware and software subsystems.
 - If virtual address refers to a part of the program that is currently in the main memory, it is accessed immediately.
 - If the address refers to a part of the program that is not currently in the main memory, it is first transferred to the main memory before it can be used.



Virtual memory organization



- *Memory management unit (MMU) translates virtual addresses into physical addresses.*
- *If the desired data or instructions are in the main memory they are fetched as described previously.*
- *If the desired data or instructions are not in the main memory, they must be transferred from secondary storage to the main memory.*
- *MMU causes the operating system to bring the data from the secondary storage into the main memory.*

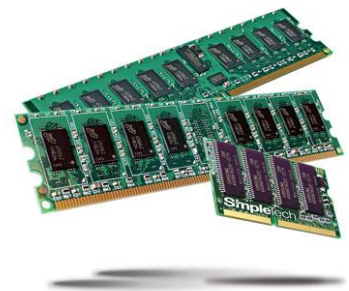
Address translation

- Assume that program and data are composed of fixed-length units called pages.
- A page consists of a block of words that occupy contiguous locations in the main memory.
- Page is a basic unit of information that is transferred between secondary storage and main memory.
- Size of a page commonly ranges from 2K to 16K bytes.
 - Pages should not be too small, because the access time of a secondary storage device is much larger than the main memory.
 - Pages should not be too large, else a large portion of the page may not be used, and it will occupy valuable space in the main memory.



Address translation (contd..)

- Concepts of virtual memory are similar to the concepts of cache memory.
- Cache memory:
 - Introduced to bridge the speed gap between the processor and the main memory.
 - Implemented in hardware.
- Virtual memory:
 - Introduced to bridge the speed gap between the main memory and secondary storage.
 - Implemented in part by software.



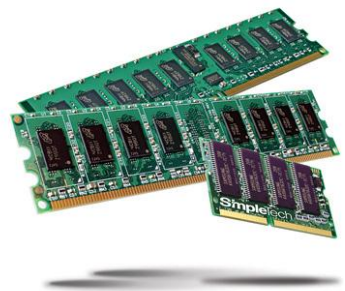
Address translation (contd..)

- Each virtual or logical address generated by a processor is interpreted as a virtual page number (high-order bits) plus an offset (low-order bits) that specifies the location of a particular byte within that page.
- Information about the main memory location of each page is kept in the page table.
 - Main memory address where the page is stored.
 - Current status of the page.
- Area of the main memory that can hold a page is called as page frame.
- Starting address of the page table is kept in a page table base register.

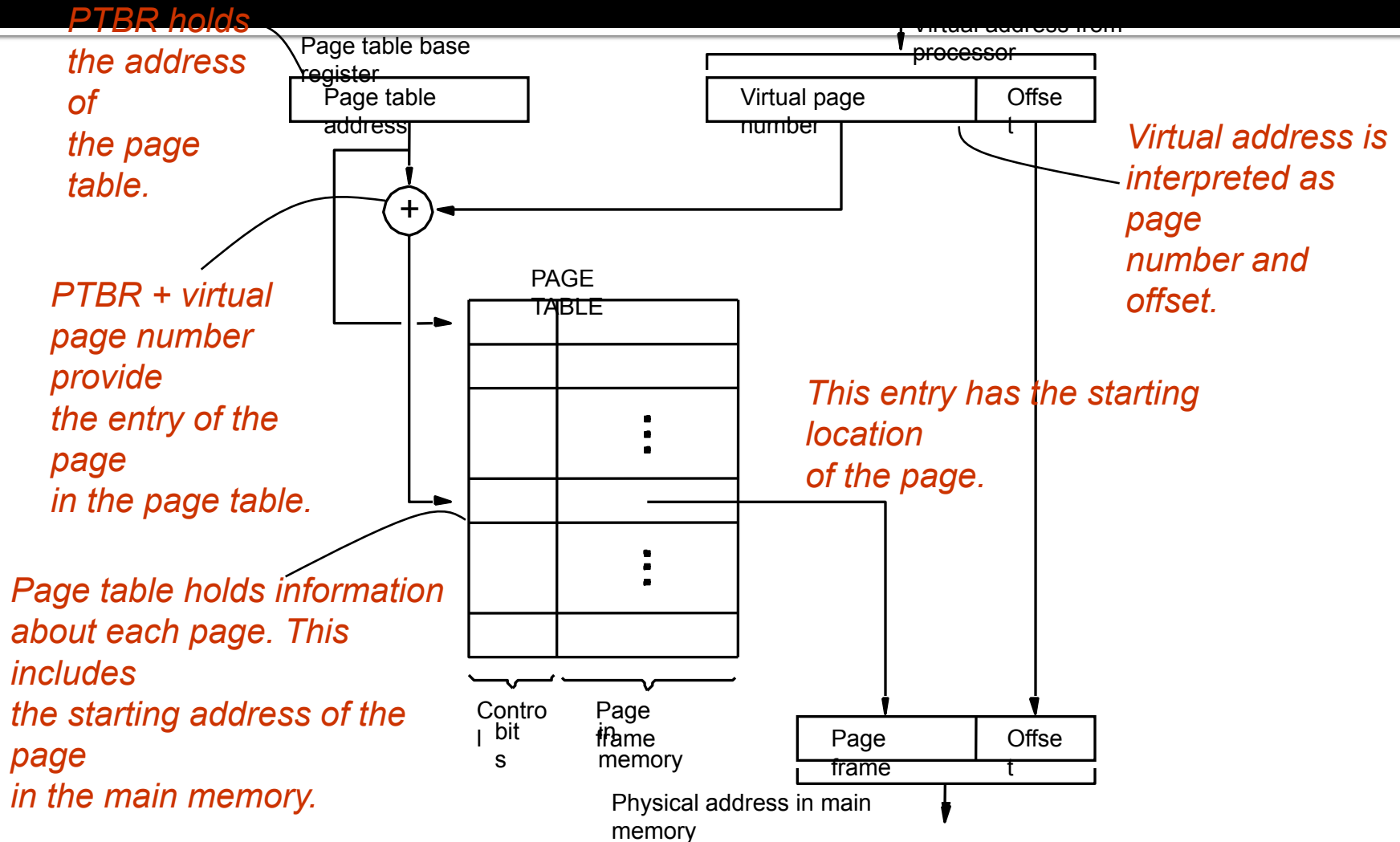


Address translation (contd..)

- Virtual page number generated by the processor is added to the contents of the page table base register.
 - This provides the address of the corresponding entry in the page table.
- The contents of this location in the page table give the starting address of the page if the page is currently in the main memory.



Address translation (contd..)



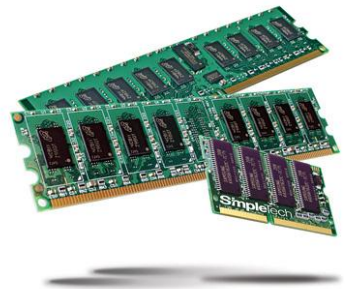
Address translation (contd..)

- Page table entry for a page also includes some control bits which describe the status of the page while it is in the main memory.
- One bit indicates the validity of the page.
 - Indicates whether the page is actually loaded into the main memory.
 - Allows the operating system to invalidate the page without actually removing it.
- One bit indicates whether the page has been modified during its residency in the main memory.
 - This bit determines whether the page should be written back to the disk when it is removed from the main memory.
 - Similar to the dirty or modified bit in case of cache memory.



Address translation (contd..)

- Other control bits for various other types of restrictions that may be imposed.
 - For example, a program may only have read permission for a page, but not write or modify permissions.



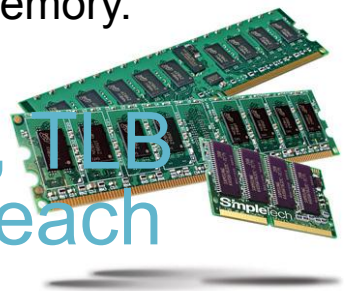
Address translation (contd..)

- Where should the page table be located?
- Recall that the page table is used by the MMU for every read and write access to the memory.
 - Ideal location for the page table is within the MMU.
- Page table is quite large.
- MMU is implemented as part of the processor chip.
- Impossible to include a complete page table on the chip.
- Page table is kept in the main memory.
- A copy of a small portion of the page table can be accommodated within the MMU.
 - Portion consists of page table entries that correspond to the most recently accessed pages.

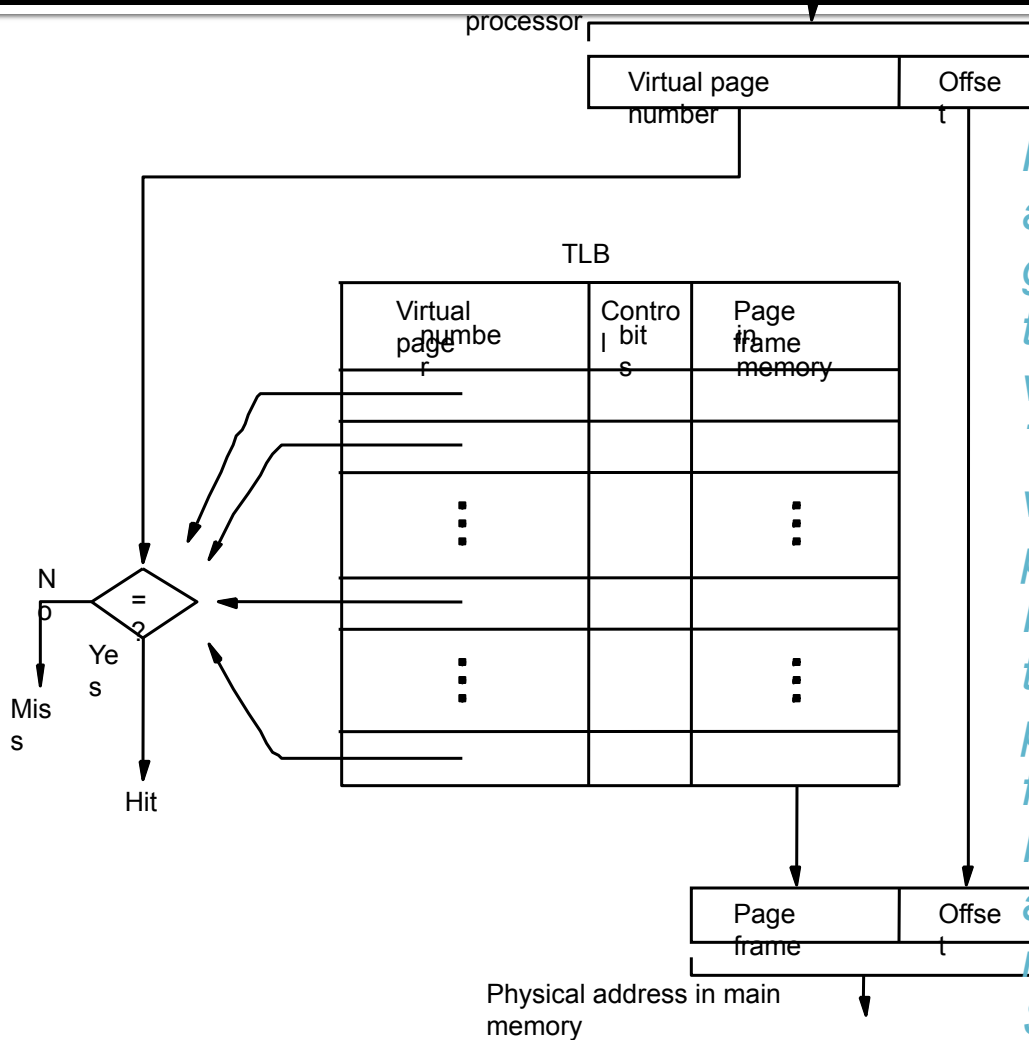


Address translation (contd..)

- A small cache called as Translation Lookaside Buffer (TLB) is included in the MMU.
 - TLB holds page table entries of the most recently accessed pages.
- Recall that cache memory holds most recently accessed blocks from the main memory.
 - Operation of the TLB and page table in the main memory is similar to the operation of the cache and main memory.
- Page table entry for a page includes:
 - Address of the page frame where the page resides in the main memory.
 - Some control bits.
- In addition to the above for each page, TLB must hold the virtual page number for each page.



Address translation (contd..)



Associative-mapped TLB

High-order bits of the virtual address generated by the processor select the virtual page.

These bits are compared to the virtual page numbers in the TLB. If there is a match, a hit occurs and the corresponding address of the page frame is read.

If there is no match, a miss occurs and the page table within the main memory must be consulted.

Set-associative mapped TLBs are found in commercial processors.

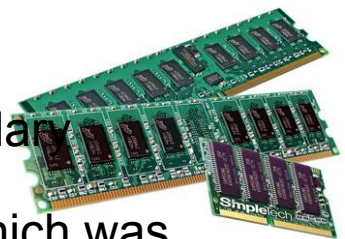
Address translation (contd..)

- How to keep the entries of the TLB coherent with the contents of the page table in the main memory?
- Operating system may change the contents of the page table in the main memory.
 - Simultaneously it must also invalidate the corresponding entries in the TLB.
- A control bit is provided in the TLB to invalidate an entry.
- If an entry is invalidated, then the TLB gets the information for that entry from the page table.
 - Follows the same process that it would follow if the entry is not found in the TLB or if a “miss” occurs.



Address translation (contd..)

- What happens if a program generates an access to a page that is not in the main memory?
- In this case, a page fault is said to occur.
 - Whole page must be brought into the main memory from the disk, before the execution can proceed.
- Upon detecting a page fault by the MMU, following actions occur:
 - MMU asks the operating system to intervene by raising an exception.
 - Processing of the active task which caused the page fault is interrupted.
 - Control is transferred to the operating system.
 - Operating system copies the requested page from secondary storage to the main memory.
 - Once the page is copied, control is returned to the task which was interrupted.



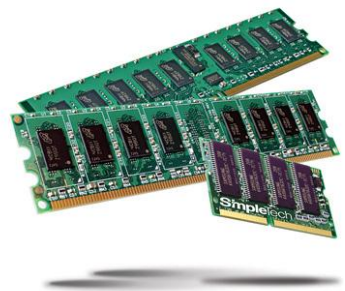
Address translation (contd..)

- Servicing of a page fault requires transferring the requested page from secondary storage to the main memory.
- This transfer may incur a long delay.
- While the page is being transferred the operating system may:
 - Suspend the execution of the task that caused the page fault.
 - Begin execution of another task whose pages are in the main memory.
- Enables efficient use of the processor



Address translation (contd..)

- How to ensure that the interrupted task can continue correctly when it resumes execution?
- There are two possibilities:
 - Execution of the interrupted task must continue from the point where it was interrupted.
 - The instruction must be restarted.
- Which specific option is followed depends on the design of the processor.



Address translation (contd..)

- When a new page is to be brought into the main memory from secondary storage, the main memory may be full.
 - Some page from the main memory must be replaced with this new page.
- How to choose which page to replace?
 - This is similar to the replacement that occurs when the cache is full.
 - The principle of locality of reference (?) can also be applied here.
 - A replacement strategy similar to LRU can be applied.
- Since the size of the main memory is relatively larger compared to cache, a relatively large amount of programs and data can be held in the main memory.
 - Minimizes the frequency of transfers between secondary storage and main memory.



Address translation (contd..)

- A page may be modified during its residency in the main memory.
- When should the page be written back to the secondary storage?
- Recall that we encountered a similar problem in the context of cache and main memory:
 - Write-through protocol(?)
 - Write-back protocol(?)
- Write-through protocol cannot be used, since it will incur a long delay each time a small amount of data is written to the disk.

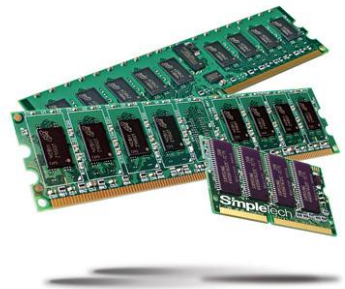


Memory Management

The Memory System

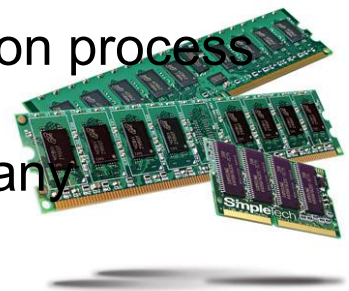
Memory management

- Operating system is concerned with transferring programs and data between secondary storage and main memory.
- Operating system needs memory routines in addition to the other routines.
- Operating system routines are assembled into a virtual address space called system space.
- System space is separate from the space in which user application programs reside.
 - This is user space.
- Virtual address space is divided into one system space + several user spaces.



Memory management (contd..)

- Recall that the Memory Management Unit (MMU) translates logical or virtual addresses into physical addresses.
- MMU uses the contents of the page table base register to determine the address of the page table to be used in the translation.
 - Changing the contents of the page table base register can enable us to use a different page table, and switch from one space to another.
- At any given time, the page table base register can point to one page table.
 - Thus, only one page table can be used in the translation process at a given time.
 - Pages belonging to only one space are accessible at any given time.



Memory management (contd..)

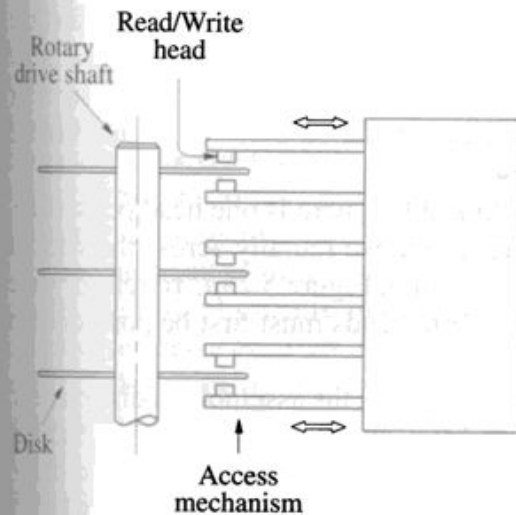
- When multiple, independent user programs coexist in the main memory, how to ensure that one program does not modify/destroy the contents of the other?
- Processor usually has two states of operation:
 - Supervisor state.
 - User state.
- Supervisor state:
 - Operating system routines are executed.
- User state:
 - User programs are executed.
 - Certain privileged instructions cannot be executed in user state.
 - These privileged instructions include the ones which change page table base register.
 - Prevents one user from accessing the space of other users.



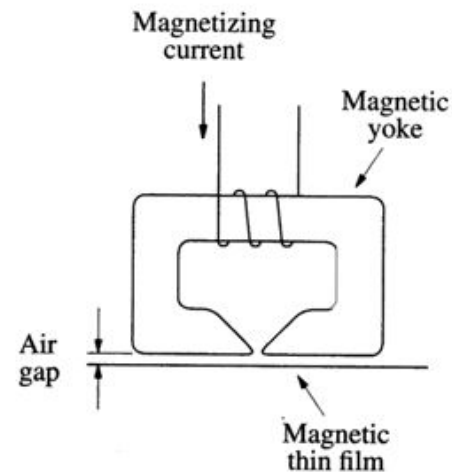
Secondary Storage

The Memory System

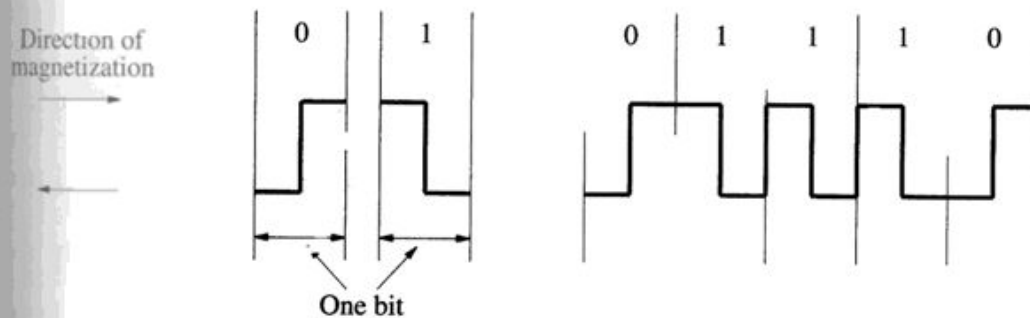
Magnetic Hard Disks



(a) Mechanical structure



(b) Read/Write head detail



(c) Bit representation by phase encoding

Disk

Disk drive

Disk controller



Organization of Data on a Disk

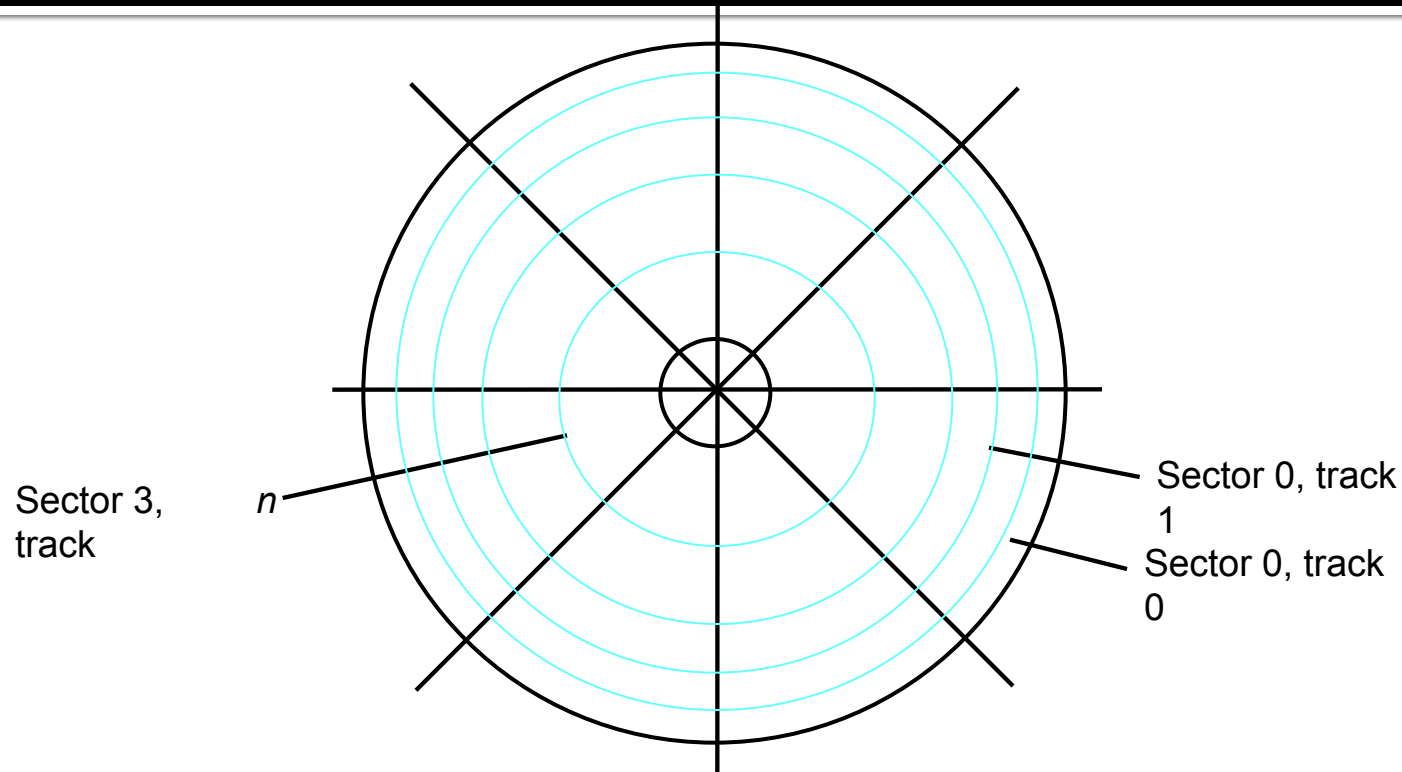
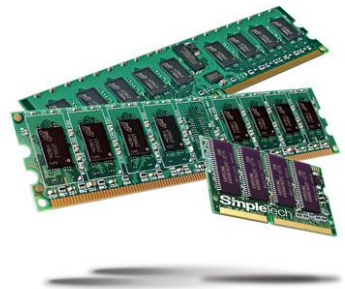
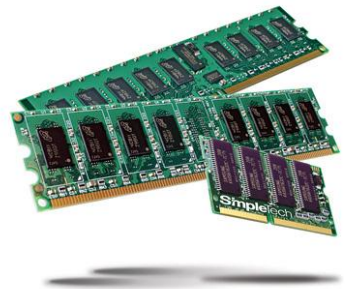


Figure 5.30. Organization of one surface of a disk.



Access Data on a Disk

- Sector header
- Following the data, there is an error-correction code (ECC).
- Formatting process
- Difference between inner tracks and outer tracks
- Access time – seek time / rotational delay (latency time)
- Data buffer/cache



Disk Controller

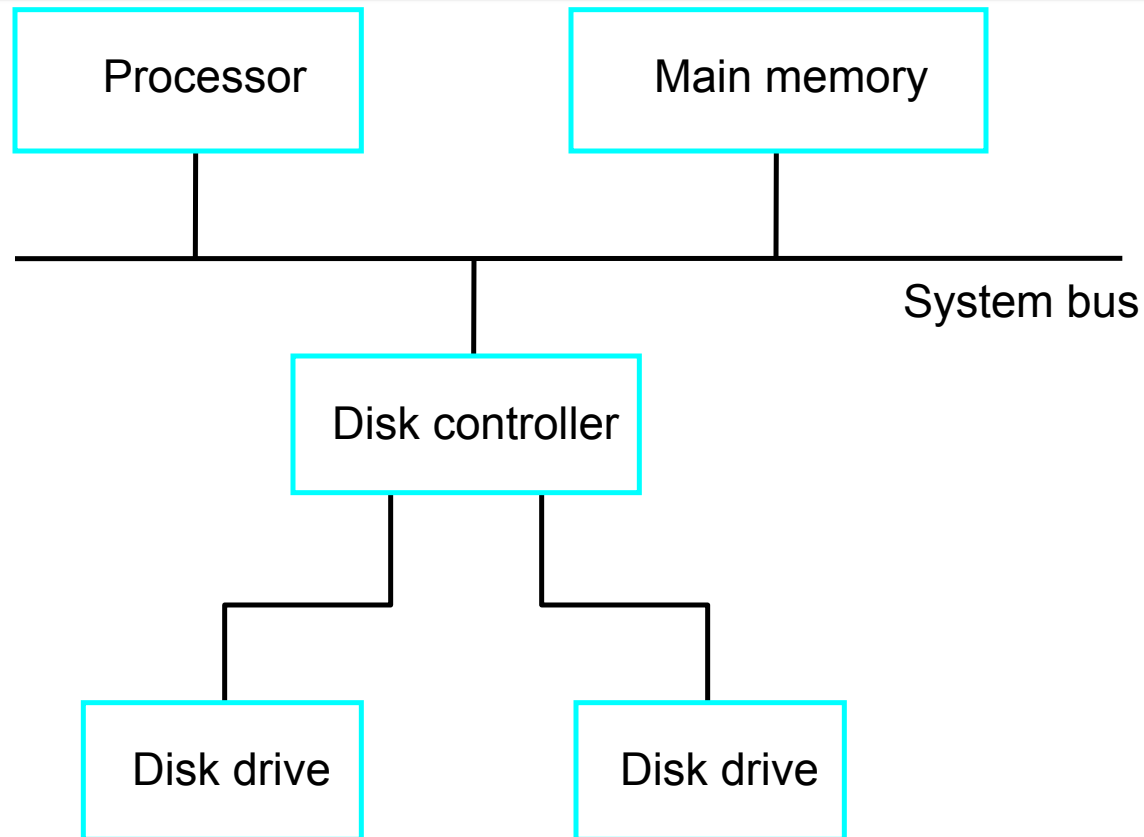
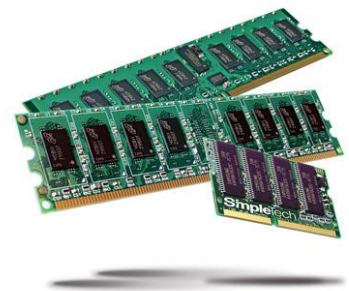
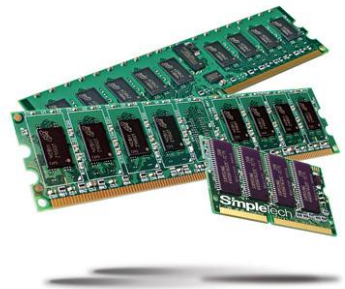


Figure 5.31. Disks connected to the system bus.



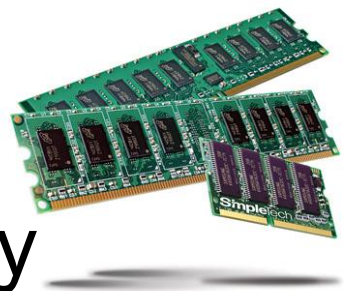
Disk Controller

- Seek
- Read
- Write
- Error checking

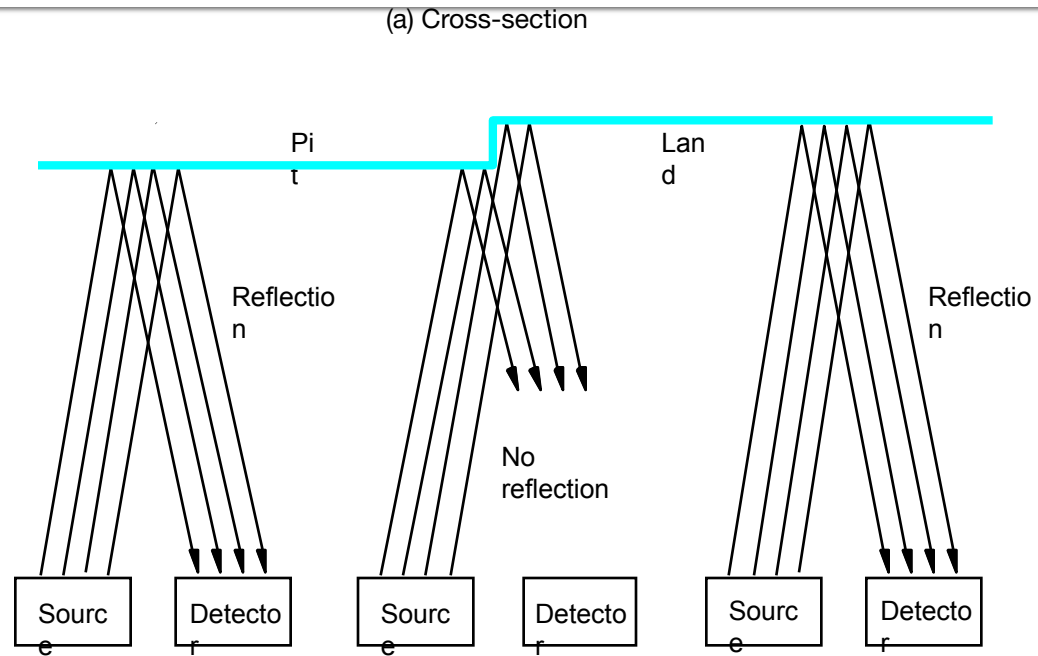


RAID Disk Arrays

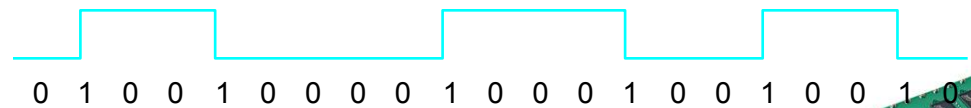
- Redundant Array of Inexpensive Disks
- Using multiple disks makes it cheaper for huge storage, and also possible to improve the reliability of the overall system.
- RAID0 – data striping
- RAID1 – identical copies of data on two disks
- RAID2, 3, 4 – increased reliability
- RAID5 – parity-based error-recovery



Optical Disks

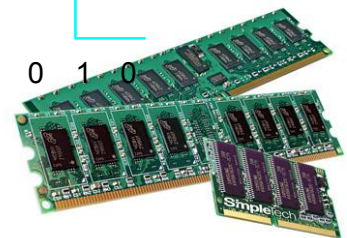


(b) Transition from pit to land



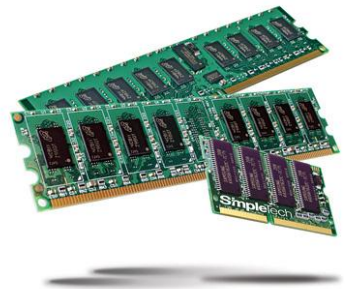
(c) Stored binary pattern

Figure 5.32. Optical disk.



Optical Disks

- CD-ROM
- CD-Recordable (CD-R)
- CD-ReWritable (CD-RW)
- DVD
- DVD-RAM



Magnetic Tape Systems

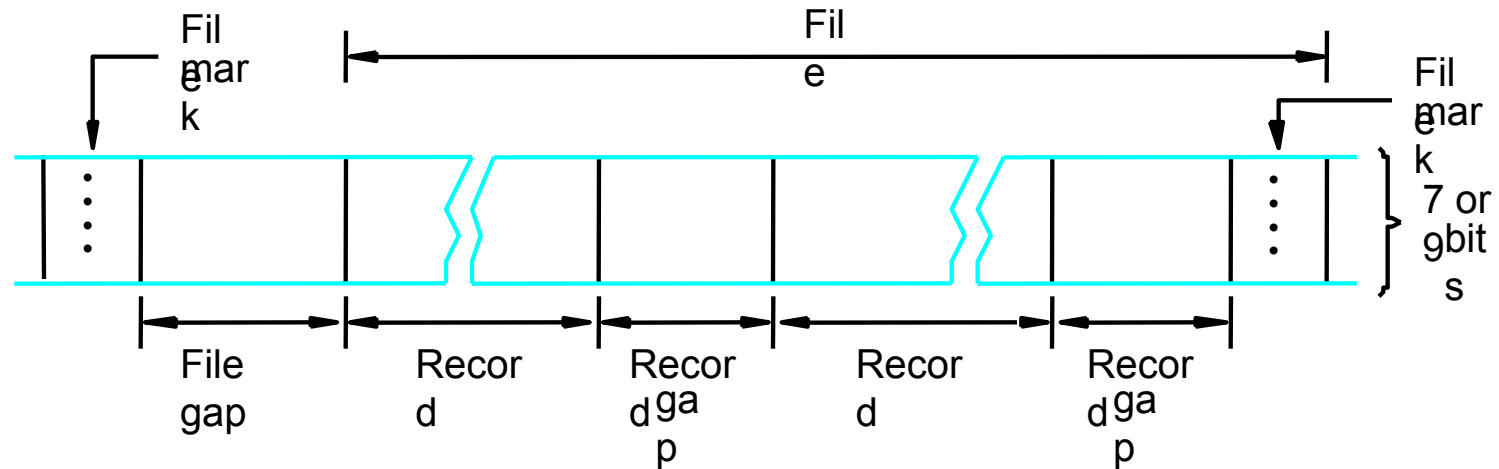


Figure 5.33. Organization of data on magnetic tape.

