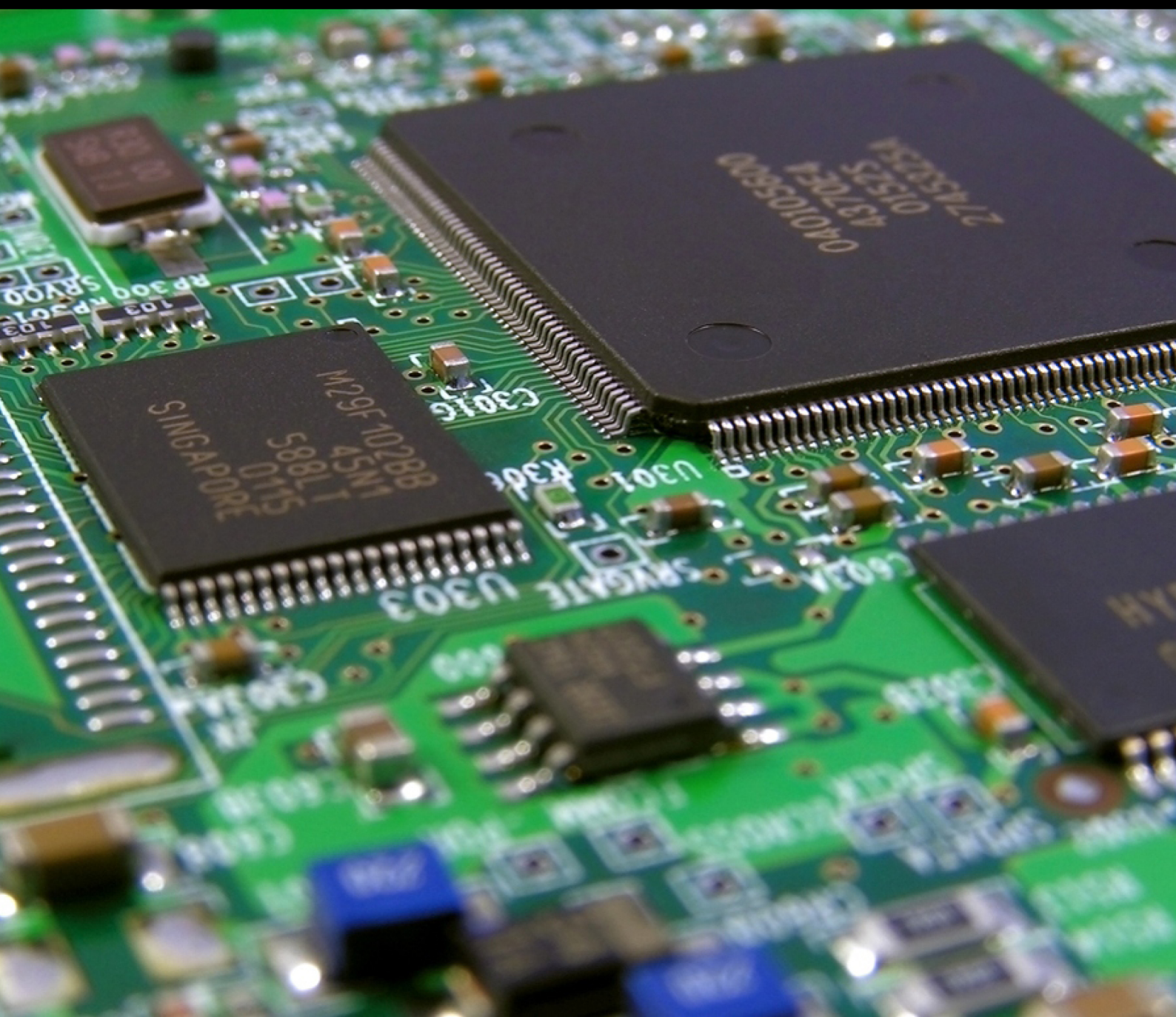


Computer Memory System

Derrick Shepherd



COMPUTER MEMORY SYSTEM

COMPUTER MEMORY SYSTEM

Derrick Shepherd



Computer Memory System
by Derrick Shepherd

Copyright© 2022 BIBLIOTEX

www.bibliotex.com

All rights reserved. No part of this book may be reproduced or used in any manner without the prior written permission of the copyright owner, except for the use brief quotations in a book review.

To request permissions, contact the publisher at info@bibliotex.com

Ebook ISBN: 9781984663962



Published by:

Bibliotex

Canada

Website: www.bibliotex.com

Contents

Chapter 1	Introduction	1
Chapter 2	Memory Management System	21
Chapter 3	Main Units: Central Processing Unit	44
Chapter 4	Computer Hard Disk Drive	89
Chapter 5	Computer File System	125
Chapter 6	Computer External Memory System	146

1

Introduction

Computer memory can refer to many types of memory within a computer, but, typically, it refers to random access memory (RAM). It is physically found on computer chips that are inserted onto the computer's motherboard. RAM is electronic, rather than mechanical; that is, it does not have moving parts and therefore data access to it is very fast. Modern computers often have somewhere between 256 MB (megabytes) and 2 GB (gigabytes) of RAM, although there are, of course, computers with more or less RAM.

RAM is also volatile, meaning that it gets lost when the computer is switched off. The expensive nature of RAM spurred the creation of another type of computer memory called virtual memory. With virtual memory, a slow down in performance is observed only when you try to operate a programme whose files are in the virtual memory. In essence, this slow down is only observed when shifting between

programmes. In this way, virtual memory often provides a cheaper alternative to RAM.

A third type of computer memory is cache. There are two types of cache. Primary cache, or level 1 cache, is built right into the central processing unit (CPU) and ensures instant availability of data that the CPU frequently needs. Secondary cache, or level 2 cache, is usually built on a memory chip, is located very close to the CPU, and has a direct connection to the CPU through a dedicated circuit. Secondary cache is bigger in capacity than primary cache.

Cache basically speeds up the rate at which data moves from the main memory to the CPU. The registers form a fourth type of computer memory. These are units within the CPU that contain specific types of data, especially for the Arithmetic and Logic Unit (ALU). A final group of computer memory is called flash. This is a solid-state, rewritable type of memory. Examples of flash memory include BIOS and memory cards. Just like the RAM, they are electronic and not mechanical. They are also non-volatile, and are therefore suitable for digital cameras, mobile phones and other miniaturized computers

OVERVIEW

The system memory is the place where the computer holds current programmes and data that are in use. There are various levels of computer memory (memory), including ROM, RAM, cache, page and graphics, each with specific objectives for system operation. This section focusses on the role of computer memory, and the technology behind it. Although memory is used in many different forms around modern PC

systems, it can be divided into two essential types: RAM and ROM. ROM, or Read Only Memory, is relatively small, but essential to how a computer works. ROM is always found on motherboards, but is increasingly found on graphics cards and some other expansion cards and peripherals. Generally speaking, ROM does not change. It forms the basic instruction set for operating the hardware in the system, and the data within remains intact even when the computer is shut down. It is possible to update ROM, but it's only done rarely, and at need. If ROM is damaged, the computer system simply cannot function.

RAM, or Random Access Memory, is "volatile." This means that it only holds data while power is present. RAM changes constantly as the system operates, providing the storage for all data required by the operating system and software. Because of the demands made by increasingly powerful operating systems and software, system RAM requirements have accelerated dramatically over time.

For instance, at the turn of the millennium a typical computer may have only 128Mb of RAM in total, but in 2007 computers commonly ship with 2Gb of RAM installed, and may include graphics cards with their own additional 512Mb of RAM and more. Clearly, modern computers have significantly more memory than the first PCs of the early 1980s, and this has had an effect on development of the PC's architecture. The trouble is, storing and retrieving data from a large block of memory is more time-consuming than from a small block. With a large amount of memory, the difference in time between a register access and a memory access is very great, and this has resulted in extra layers of

cache in the storage hierarchy. When accessing memory, a fast processor will demand a great deal from RAM. At worst, the CPU may have to waste clock cycles while it waits for data to be retrieved. Faster memory designs and motherboard buses can help, but since the 1990s "cache memory" has been employed as standard between the main memory and the processor.

Not only this, CPU architecture has also evolved to include ever larger internal caches. The organization of data this way is immensely complex, and the system uses ingenious electronic controls to ensure that the data the processor needs next is already in cache, physically closer to the processor and ready for fast retrieval and manipulation. Read on for a closer look at the technology behind computer memory, and how developments in RAM and ROM have enabled systems to function with seemingly exponentially increasing power.

Memory System

Computer memory is used to store two things: i) instructions to execute a programme and ii) data. When the computer is doing any job, the data that have to be processed are stored in the primary memory. This data may come from an input device like keyboard or from a secondary storage device like a floppy disk.

As programme or the set of instructions is kept in primary memory, the computer is able to follow instantly the set of instructions. For example, when you book ticket from railway reservation counter, the computer has to follow the same steps: take the request, check the availability of seats, calculate fare, wait for money to be paid, store the reservation

and get the ticket printed out. The programme containing these steps is kept in memory of the computer and is followed for each request.

But inside the computer, the steps followed are quite different from what we see on the monitor or screen. In computer's memory both programmes and data are stored in the binary form. You have already been introduced with decimal number system, that is the numbers 1 to 9 and 0. The binary system has only two values 0 and 1. These are called *bits*. As human beings we all understand decimal system but the computer can only understand binary system.

It is because a large number of integrated circuits inside the computer can be considered as switches, which can be made ON, or OFF. If a switch is ON it is considered 1 and if it is OFF it is 0. A number of switches in different states will give you a message like this: 110101....10. So the computer takes input in the form of 0 and 1 and gives output in the form 0 and 1 only. Is it not absurd if the computer gives outputs as 0's and 1's only? But you do not have to worry about.

Every number in binary system can be converted to decimal system and vice versa; for example, 1010 meaning decimal 10. Therefore it is the computer that takes information or data in decimal form from you, convert it in to binary form, process it producing output in binary form and again convert the output to decimal form.

The primary memory as you know in the computer is in the form of IC's (Integrated Circuits). These circuits are called Random Access Memory (RAM). Each of RAM's locations stores one *byte* of information. (One *byte* is equal to 8 *bits*).

A bit is an acronym for *binary digit*, which stands for one binary piece of information. This can be either 0 or 1. You will know more about RAM later. The Primary or internal storage section is made up of several small storage locations (ICs) called cells. Each of these cells can store a fixed number of bits called *word length*.

Each cell has a unique number assigned to it called the address of the cell and it is used to identify the cells. The address starts at 0 and goes up to (N-1). You should know that the memory is like a large cabinet containing as many drawers as there are addresses on memory. Each drawer contains a word and the address is written on outside of the drawer.

CLASSIFICATION OF MEMORY: RAM

RAM is perhaps the most important of the input/output devices. When we talk about computer memory, we are mainly talking about RAM. In this class, when we think about the banks of light switches that can be manipulated, we are thinking of RAM memory. The term Random Access is pretty unfortunate. There is nothing random about how memory is accessed. Random access memory or RAM most commonly refers to computer chips that temporarily store dynamic data to enhance computer performance.

By storing frequently used or active files in random access memory, the computer can access the data faster than if it to retrieve it from the far-larger hard drive. Random access memory is also used in printers and other devices. Random access memory is volatile memory, meaning it loses its contents once power is cut. This is different from non-volatile

memory such as hard disks and flash memory, which do not require a power source to retain data.

When a computer shuts down properly, all data located in random access memory is committed to permanent storage on the hard drive or flash drive. At the next boot-up, RAM begins to fill with programmes automatically loaded at startup, and with files opened by the user.

There are several different types of random access memory chips which come several to a "stick." A stick of RAM is a small circuit board shaped like a large stick of gum. Sticks of RAM fit into "banks" on the motherboard.

Adding one or more sticks increases RAM storage and performance. Random access memory is categorized by architecture and speed. As technology progresses, RAM chips become faster and employ new standards so that RAM must be matched to a compatible motherboard. The motherboard will only support certain types of random access memory, and it will also have a limit as to the amount of RAM it can support.

For example, one motherboard may support dual-channel Synchronous Dynamic Random Access Memory (SDRAM), while an older motherboard might only support Single In-line Memory Modules (SIMMS) or Dual In-line Memory Modules (DIMMS). Since random access memory can improve performance, the type and amount of RAM a motherboard will support becomes a major factor when considering a new computer.

If there is a faster, better random access memory chip on the market, the buyer will want to consider purchasing a motherboard capable of using it. A year down the road, that

'new' RAM might be standard, while the buyer may be stuck with an old style motherboard. A new variety of non-volatile random access memory made with nanotubes or other technologies will likely be forthcoming in the near future.

These RAM chips would retain data when powered down. Memory, which is commonly referred to as RAM (Random Access Memory), is a temporary (Volatile) storage area utilized by the CPU. Before a programme can be run, the programme is loaded into the memory which allows the CPU direct access to the programme.

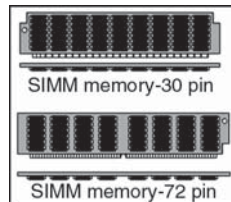
Memory is a necessity for any computer and it is recommend that you have at least 64MB of memory for your IBM or Macintosh. Memory is commonly confused with Hard Drive Space. There are two types of memory; the first type of memory is the memory explained in the above paragraph, this memory is available in computer chips; the other type of memory is actually Hard Drive Space which is stored on the computer Hard Disk Drive.

The Hard drive is actually a physical drive which contains several parts and is generally larger than the amount of memory found in your computer. The below explanation will help in describing the most commonly found types of RAM in computers.

Single In-line memory module (SIMM)

A SIMM, or single in-line memory module, is a type of memory module containing random access memory used in computers from the early 1980s to the late 1990s. It differs from a dual in-line memory module (DIMM), the most predominant form of memory module today, in that the contacts on a SIMM are redundant on both sides of the

module. SIMMs were standardised under the JEDEC JESD-21C standard.



A slender circuit board dedicated to storing memory chips. Each chip is capable of holding 8 to 9 chips per board, the ninth chip usually an error-checking chip (parity/ non parity). The typical bus from the chip to the motherboard is 32-bits wide. When upgrading a Pentium motherboard you will be required to upgrade 2 of the same type of chips at the same time to accommodate the Pentium processor.

Speed

Memory to a computer that was designed to run slower memory. However, your system will operate at the speed of the slowest memory module.

One thing to keep in mind is that the memory does need to be the same type - for example, SDRAM cannot be mixed with DDR, and DDR cannot be mixed with DDR2 and DDR2 cannot work in a DDR3 system.

We recommend that you use the Crucial Memory Advisor™ or System Scanner tools to find the right memory for your computer. Rather than give memory modules catchy names, the industry refers to modules by their specifications. But if you don't know a lot about memory, the numbers can be confusing. Here's a short summary of the most popular types of memory and what the numbers refer to.

Can determine the size amount of the chip by looking at the part number of each chip on the SIMM board. For 2-, 8-

and 9- chip SIMMs, all the chips should have the same part numbers. Look at the number that ends with a dash and a digit such as "-7". This is the rate speed or nanoseconds of the chip. With "-7" this would indicate that the memory is 70ns.

Size

Look at the four digits to the left of this number; these often carry information about the number of bits in the chip. A 4256 indicates 256K bits arranged in sets of four, for a total of 1Mb. "1000" indicates 1MB of bits arranged in one set. With some types of memory, the last one or two digits may be changed to indicate different kinds of memory; there are 1MB chips that end with 4256, 4257, and 4258. In this case, round the last digits to an even 256 or thousand.

Three-chip SIMMs will typically have two larger chips that are four times the capacity of the third chip (because 4 plus 4 plus 1 makes 9, which is the number of bits needed per byte including parity).

Parity/ Non- Parity

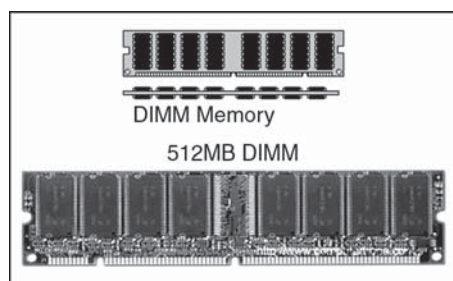
To determine if the SIMM is Parity/ Non- Parity, look for x36/ x9 which indicate that the chip is parity (Error checking). x36 is used with 72-pin SIMMs and x9 is for 30-pin SIMMs. If x32/ x8 this would be an indication that the chip is Non- Parity (Non- Error checking) x32 is used with 72-pin SIMMs and x8 is used with 30-pin SIMMs. Another method: count the chips. If you see three or nine discrete chips, the SIMM probably includes a parity. If there are two or eight chips, the SIMM probably does not include parity bit. In this case, divide the number of bits by 8 to determine bytes.

SIMM Memory and Pentium Computers

When updating the computer's memory with a Pentium processor, ensure that you purchase two SIMMs rather than one. Such as if planning to upgrade to 32 MB of RAM that you use two 16 MB SIMMs rather than one. This must be done to allow the memory to work properly with the Pentium processor.

Dual In-line Memory Module (DIMM)

A DIMM or dual in-line memory module comprises a series of dynamic random-access memory integrated circuits. These modules are mounted on a printed circuit board and designed for use in personal computers, workstations and servers. DIMMs began to replace SIMMs (single in-line memory modules) as the predominant type of memory module as Intel P5-based Pentium processors began to gain market share.



Type of Circuit board that holds memory chips. DIMMS have a 64-bit path because of the Pentium Processor requirements. Because of the new bit path, DIMMS can be installed one at a time un-like SIMMs which on a Pentium would require two be added.

Some of the Advantages DIMMs Have Over SIMMs

DIMMs have separate contacts on each side of the board, thereby providing twice as much data as a single SIMM. The

command address and control signals are buffered on the DIMMs. With heavy memory requirements this will reduce the loading effort of the memory.

Static RAM (SRAM)

Static random-access memory (SRAM) is a type of semiconductor memory that uses bistable latching circuitry to store each bit. The term static differentiates it from dynamic RAM (DRAM) which must be periodically refreshed. SRAM exhibits data remanence, but it is still volatile in the conventional sense that data is eventually lost when the memory is not powered.

Static RAM is a type of RAM that holds its data without external refresh, for as long as power is supplied to the circuit. This is contrasted to dynamic RAM (DRAM), which must be refreshed many times per second in order to hold its data contents. SRAMs are used for specific applications within the PC, where their strengths outweigh their weaknesses compared to DRAM:

- *Simplicity*: SRAMs don't require external refresh circuitry or other work in order for them to keep their data intact.
- *Speed*: SRAM is faster than DRAM.

In contrast, SRAMs have the following weaknesses, compared to DRAMs:

- *Cost*: SRAM is, byte for byte, several times more expensive than DRAM.
- *Size*: SRAMs take up much more space than DRAMs (which is part of why the cost is higher).

These advantages and disadvantages taken together obviously show that performance-wise, SRAM is superior to

DRAM, and we would use it exclusively if only we could do so economically.

Unfortunately, 32 MB of SRAM would be prohibitively large and costly, which is why DRAM is used for system memory. SRAMs are used instead for level 1 cache and level 2 cache memory, for which it is perfectly suited; cache memory needs to be very fast, and not very large. SRAM is manufactured in a way rather similar to how processors are: highly-integrated transistor patterns photo-etched into silicon. Each SRAM bit is comprised of between four and six transistors, which is why SRAM takes up much more space compared to DRAM, which uses only one (plus a capacitor). Because an SRAM chip is comprised of thousands or millions of identical cells, it is much easier to make than a CPU, which is a large die with a non-repetitive structure.

Dynamic RAM (DRAM)

Dynamic random-access memory (DRAM) is a type of random-access memory that stores each bit of data in a separate capacitor within an integrated circuit. The capacitor can be either charged or discharged; these two states are taken to represent the two values of a bit, conventionally called 0 and 1. Since capacitors leak charge, the information eventually fades unless the capacitor charge is refreshed periodically. Because of this refresh requirement, it is a dynamic memory as opposed to SRAM and other static memory.

The main memory (the "RAM") in personal computers is dynamic RAM (DRAM). It is the RAM in desktops, laptops and workstation computers as well as some of the RAM of video game consoles.

The advantage of DRAM is its structural simplicity: only one transistor and a capacitor are required per bit, compared to four or six transistors in SRAM. This allows DRAM to reach very high densities. Unlike flash memory, DRAM is volatile memory (vs. non-volatile memory), since it loses its data quickly when power is removed. The transistors and capacitors used are extremely small; billions can fit on a single memory chip.

Dynamic RAM is a type of RAM that only holds its data if it is continuously accessed by special logic called a refresh circuit. Many hundreds of times each second, this circuitry reads the contents of each memory cell, whether the memory cell is being used at that time by the computers or not.

Due to the way in which the cells are constructed, the reading action itself refreshes the contents of the memory. If this is not done regularly, then the DRAM will lose its contents, even if it continues to have power supplied to it. This refreshing action is why the memory is called dynamic. All PCs use DRAM for their main system memory, instead of SRAM, even though DRAMs are slower than SRAMs and require the overhead of the refresh circuitry. It may seem weird to want to make the computer's memory out of something that can only hold a value for a fraction of a second. In fact, DRAMs are both more complicated and slower than SRAMs.

The reason that DRAMs are used is simple: they are much cheaper and take up much less space, typically 1/4 the silicon area of SRAMs or less. To build a 64 MB core memory from SRAMs would be very expensive. The overhead of the refresh circuit is tolerated in order to allow the use of large amounts

of inexpensive, compact memory. The refresh circuitry itself is almost never a problem; many years of using DRAM has caused the design of these circuits to be all but perfected. DRAMs are smaller and less expensive than SRAMs because SRAMs are made from four to six transistors (or more) per bit, DRAMs use only one, plus a capacitor.

The capacitor, when energized, holds an electrical charge if the bit contains a "1" or no charge if it contains a "0". The transistor is used to read the contents of the capacitor. The problem with capacitors is that they only hold a charge for a short period of time, and then it fades away. These capacitors are tiny, so their charges fade particularly quickly.

This is why the refresh circuitry is needed: to read the contents of every cell and refresh them with a fresh "charge" before the contents fade away and are lost. Refreshing is done by reading every "row" in the memory chip one row at a time; the process of reading the contents of each capacitor re-establishes the charge.

DRAM is manufactured using a similar process to how processors are: a silicon substrate is etched with the patterns that make the transistors and capacitors (and support structures) that comprise each bit. DRAM costs much less than a processor because it is a series of simple, repeated structures, so there isn't the complexity of making a single chip with several million individually-located transistors.

DDR-SDRAM

Double data rate synchronous dynamic random-access memory (DDR SDRAM) is a class of memory integrated circuits used in computers. DDR SDRAM, also called DDR1 SDRAM, has been superseded by DDR2 SDRAM and DDR3

SDRAM, neither of which is either forward or backward compatible with DDR1 SDRAM -meaning that DDR2 or DDR3 memory modules will not work in DDR1-equipped motherboards, and vice versa.

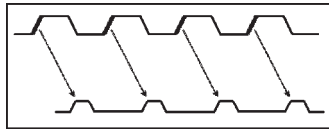
Compared to single data rate (SDR) SDRAM, the DDR SDRAM interface makes higher transfer rates possible by more strict control of the timing of the electrical data and clock signals. Implementations often have to use schemes such as phase-locked loops and self-calibration to reach the required timing accuracy. The interface uses double pumping (transferring data on both the rising and falling edges of the clock signal) to lower the clock frequency. One advantage of keeping the clock frequency down is that it reduces the signal integrity requirements on the circuit board connecting the memory to the controller. The name "double data rate" refers to the fact that a DDR SDRAM with a certain clock frequency achieves nearly twice the bandwidth of a SDR SDRAM running at the same clock frequency, due to this double pumping. With data being transferred 64 bits at a time, DDR SDRAM gives a transfer rate of (memory bus clock rate) \times 2 (for dual rate) \times 64 (number of bits transferred)/ 8 (number of bits/byte). Thus, with a bus frequency of 100 MHz, DDR SDRAM gives a maximum transfer rate of 1600 MB/s.

"Beginning in 1996 and concluding in June 2000, JEDEC developed the DDR (Double Data Rate) SDRAM specification (JESD79)." JEDEC has set standards for data rates of DDR SDRAM, divided into two parts. The first specification is for memory chips, and the second is for memory modules.

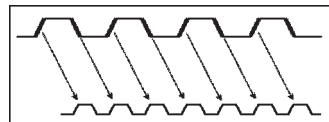
The DDR-SDRAM (Double Data Rate SDRAM) is a memory, based on the SDRAM technology, which doubles the transfer

rate of the SDRAM using the same frequency. Data are read or written into memory based on a clock.

Standard DRAM memories use a method known as SDR (Single Data Rate) involving reading or writing a piece of data at each leading edge.



The DDR doubles the frequency of reading/writing, with a clock at the same frequency, by sending data to each leading edge and to each trailing edge.



DDR memories generally have a product name such as PCXXXX where "XXXX" represents the speed in Mo/s.

DDR2-SDRAM

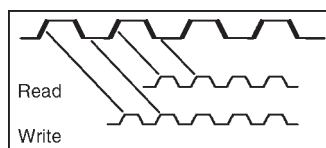
DDR2 SDRAM is a double data rate synchronous dynamic random-access memory interface. It superseded the original DDR SDRAM specification, and has since been superseded by DDR3 SDRAM. DDR2 DIMMs are neither forward compatible with DDR3 nor backward compatible with DDR.

In addition to double pumping the data bus as in DDR SDRAM (transferring data on the rising and falling edges of the bus clock signal), DDR2 allows higher bus speed and requires lower power by running the internal clock at half the speed of the data bus. The two factors combine to require a total of four data transfers per internal clock cycle. With data being transferred 64 bits at a time, DDR2 SDRAM gives a transfer rate of (memory clock rate) \times 2 (for bus clock

multiplier) $\times 2$ (for dual rate) $\times 64$ (number of bits transferred)/
8 (number of bits/byte). Thus with a memory clock frequency
of 100 MHz, DDR2 SDRAM gives a maximum transfer rate of
3200 MB/s.

Since the DDR2 internal clock runs at half the DDR
external clock rate, DDR2 memory operating at the same
external data bus clock rate as DDR results in DDR2 being
able to provide the same bandwidth but with higher latency.
Alternatively, DDR2 memory operating at twice the external
data bus clock rate as DDR may provide twice the bandwidth
with the same latency. The best-rated DDR2 memory modules
are at least twice as fast as the best-rated DDR memory
modules.

DDR2 (or DDR-II) memory achieves speeds that are twice
as high as those of the DDR with the same external frequency.
QDR (Quadruple Data Rate or quad-pumped) designates the
reading and writing method used. DDR2 memory in fact uses
two separate channels for reading and writing, so that it is
able to send or receive twice as much data as the DDR.



DDR2 also has more connectors than the classic DDR (240
for DDR2 compared with 184 for DDR).

Operation of Random Access Memory

Alternatively referred to as main memory, primary memory,
or system memory, Random Access Memory (RAM) is a
computer storage location that allows information to be stored
and accessed quickly from random locations within DRAM

on a memory module. Because information is accessed randomly instead of sequentially like a CD or hard drive the computer is able to access the data much faster than it would if it was only reading the hard drive.

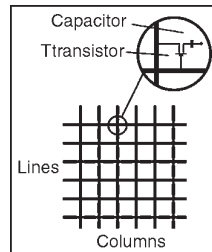
However, unlike ROM and the hard drive RAM is a volatile memory and requires power in order to keep the data accessible, if power is lost all data contained in memory lost.

As the computer loads parts of the operating system and drivers are loaded into memory, which allows the CPU to process the instructions much faster and your computer to load faster. After the operating system has loaded, each programme you open such as the browser you're using to view this page is loaded into memory while it is running. If too many programmes are open the computer will swap the data in the memory between the RAM and the hard disk drive.

The random access memory comprises hundreds of thousands of small capacitors that store loads. When loaded, the logical state of the capacitor is equal to 1, otherwise it is 0, meaning that each capacitor represents one memory bit.

Given that the capacitors become discharged they must be constantly recharged (the exact term is refresh) at regular intervals, known as the refresh cycle. DRAM memories for example require refresh cycles of around 15 nanoseconds (ns). Each capacitor is coupled with a transistor (MOS-type) enabling "recovery" or amendment of the status of the capacitor. These transistors are arranged in the form of a table (matrix) thus we access a memory box (also called memory point) via a line and a column.

Computer Memory System



Each memory point is thus characterised by an address which corresponds to a row number and a column number. This access is not instant and the access time period is known as latency time.

Consequently, time required for access to data in the memory is equal to cycle time plus latency time. Thus, for a DRAM memory, access time is 60 nanoseconds (35ns cycle time and 25ns latency time). On a computer, the cycle time corresponds to the opposite of the clock frequency; for example, for a computer with frequency of 200 MHz, cycle time is 5 ns ($1/200 \cdot 10^6$).

Consequently a computer with high frequency using memories with access time much longer than the processor cycle time must perform wait states to access the memory. For a computer with frequency of 200 MHz using DRAM memories (and access time of 60ns), there are 11 wait states for a transfer cycle. The computer's performance decreases as the number of wait states increases, therefore we recommend the use of faster memories.

2

Memory Management System

The memory management subsystem is one of the most important parts of the operating system. Since the early days of computing, there has been a need for more memory than exists physically in a system. Strategies have been developed to overcome this limitation and the most successful of these is virtual memory. Virtual memory makes the system appear to have more memory than it actually has by sharing it between competing processes as they need it.

Large Address Spaces

The operating system makes the system appear as if it has a larger amount of memory than it actually has. The virtual memory can be many times larger than the physical memory in the system,

Protection

Each process in the system has its own virtual address

space. These virtual address spaces are completely separate from each other and so a process running one application cannot affect another. Also, the hardware virtual memory mechanisms allow areas of memory to be protected against writing. This protects code and data from being overwritten by rogue applications.

Memory Mapping

Memory mapping is used to map image and data files into a processes address space. In memory mapping, the contents of a file are linked directly into the virtual address space of a process.

Fair Physical Memory Allocation

The memory management subsystem allows each running process in the system a fair share of the physical memory of the system,

Shared Virtual Memory

Although virtual memory allows processes to have separate (virtual) address spaces, there are times when you need processes to share memory. For example there could be several processes in the system running the bash command shell. Rather than have several copies of bash, one in each processes virtual address space, it is better to have only one copy in physical memory and all of the processes running bash share it. Dynamic libraries are another common example of executing code shared between several processes.

Abstract Model of Virtual Memory

Before considering the methods that Linux uses to support

virtual memory it is useful to consider an abstract model that is not cluttered by too much detail. As the processor executes a programme it reads an instruction from memory and decodes it. In decoding the instruction it may need to fetch or store the contents of a location in memory. The processor then executes the instruction and moves onto the next instruction in the programme. In this way the processor is always accessing memory either to fetch instructions or to fetch and store data.

Demand Paging

As there is much less physical memory than virtual memory the operating system must be careful that it does not use the physical memory inefficiently. One way to save physical memory is to only load virtual pages that are currently being used by the executing programme.

For example, a database programme may be run to query a database. In this case not all of the database needs to be loaded into memory, just those data records that are being examined. If the database query is a search query then it does not make sense to load the code from the database programme that deals with adding new records. This technique of only loading virtual pages into memory as they are accessed is known as demand paging.

Swapping

If a process needs to bring a virtual page into physical memory and there are no free physical pages available, the operating system must make room for this page by discarding another page from physical memory. If the page to be discarded from physical memory came from an image or data

file and has not been written to then the page does not need to be saved. Instead it can be discarded and if the process needs that page again it can be brought back into memory from the image or data file.

However, if the page has been modified, the operating system must preserve the contents of that page so that it can be accessed at a later time. This type of page is known as a *dirty* page and when it is removed from memory it is saved in a special sort of file called the swap file. Accesses to the swap file are very long relative to the speed of the processor and physical memory and the operating system must juggle the need to write pages to disk with the need to retain them in memory to be used again.

Shared Virtual Memory

Virtual memory makes it easy for several processes to share memory. All memory access are made via page tables and each process has its own separate page table.

For two processes sharing a physical page of memory, its physical page frame number must appear in a page table entry in both of their page tables shows two processes that each share physical page frame number 4. For process *X* this is virtual page frame number 4 whereas for process *Y* this is virtual page frame number 6. This illustrates an interesting point about sharing pages: the shared physical page does not have to exist at the same place in virtual memory for any or all of the processes sharing it.

Processors and Memory

There are two main components which have been part of nearly all computer systems ever designed and built. The

first is called the processor (known also as the Central Processing Unit or CPU) and the second is called the memory. The processor is the part of a computer system which does the actual computing.

That is, the part which adds, subtracts, multiplies and divides. Most processors can also compare values and perform conditional actions as a result of such comparisons. Many processors have instructions which perform various types of conversions between different representations of data. The processor itself is divided into three components that carry out the various functions that the computer is capable of performing. The first component of the processor is the controller. The controller acts as a foreman that oversees the tasks of the processor. The controller looks at the next instruction to be executed and assigns the sub-tasks that must be accomplished to carry out that instruction to the other components of the processor.

Another component of the processor is the Arithmetic-Logic Unit (ALU). This unit is the part of the processor which performs the mathematical computations and logical tasks that we expect a computer to be able to do. Addition, subtraction, comparison, etc., are all carried out by the ALU.

The last component of the processor is a collection of one or more registers. Registers are special named memory cells in the processor where information is temporarily stored during various stages of a computation. The currently executing instruction, for example, resides in a register called the Instruction Register. Since modern processors can execute millions of instructions per second it is expected that information would not stay in a register for more than a

few millionths of a second. A computer memory system is accessed (or read) by specifying the location (called an address) of the memory cell.

The memory system then responds by producing a copy of the contents of that cell. The original value of the cell is not changed by this process. This is sometimes called a non-destructive read.

A computer memory system is changed (or written) by specifying the location of the memory cell together with the new value for that cell. The previous value stored in the cell is replaced by the new value.

This is sometimes called a destructive write. The values stored in a computer memory are simply numbers. Numbers are used to represent both data and instructions. In fact, one cannot distinguish instructions from data when examining the contents of a memory system.

It is up to the programmer or operating system to keep track of which memory cells hold data and which are programme instructions. Programs have been written which manipulate and produce programs. Such programs treat instructions as data.

The processor and memory unit are wired together wiring connections called buses. A *bus* is a low resistance connection consisting of 1 or more wires.

There are three such connections. The first, called the address bus, is used by the processor to tell the memory system the number or location of the memory cell the processor wishes to access. The second bus is used to send data out to the memory. The third bus is used to transmit data and instructions to the processor.

A typical computer might be organized as indicated in the following diagram:

The processor contains a few memory cells, called registers, which are used for efficient temporary storage:

Processor Registers

The Contents of a memory cell or a register is simply a number. For purposes of illustration, suppose that each memory cell is large enough to hold numbers up to 4 decimal digits in size. If a memory cell holds an instruction, use the first two decimal digits represent the operation and the remaining digits to represent the address or location of the instruction operand.

Instruction Format

Using this scheme, we could set up the following operation codes:

Operation	Code	Function
add	01	$c(\text{acc}) = c(\text{acc}) + c(\text{addr})$
sub	02	$c(\text{acc}) = c(\text{acc}) - c(\text{addr})$
load	03	$c(\text{acc}) = c(\text{addr})$
store	04	$c(\text{addr}) = c(\text{acc})$

These codes do not correspond to any known real computer, but rather, they are the operation codes for a hypothetical model computer which we will use to illustrate important aspects of computer organization.

Computer System

In describing computer system, a distinction is often made between computer architecture and computer organization. Computer architecture refers to those attributes of a system visible to a programmer, or put another way, those attributes that have a direct impact on the logical execution of a

program. Computer organization refers to the operational units and their interconnection that realise the architecture specification. Examples of architecture attributes include the instruction set, the number of bit to represent various data types (*e.g.*, numbers, and characters), I/O mechanisms, and technique for addressing memory. Examples of organization attributes include those hardware details transparent to the programmer, such as control signals, interfaces between the computer and peripherals, and the memory technology used. As an example, it is an architectural design issue whether a computer will have a multiply instruction. It is an organizational issue whether that instruction will be implemented by a special multiply unit or by a mechanism that makes repeated use of the add unit of the system.

The organization decision may be bases on the anticipated frequency of use of the multiply instruction, the relative speed of the two approaches, and the cost and physical size of a special multiply unit.

Historically, and still today, the distinction between architecture and organization has been an important one. Many computer manufacturers offer a family of computer model, all with the same architecture but with differences in organization. Consequently, the different models in the family have different price and performance characteristics. Furthermore, an architecture may survive many years, but its organization changes with changing technology.

Development of Computers

First Generation: Vacuum Tubes

ENIAC: The ENIAC (Electronic Numerical Integrator And

Computer), designed by and constructed under the supervision of Jonh Mauchly and John Presper Eckert at the University of Pennsylvania, was the world's first general-purpose electronic digital computer.

The project was a response to U.S. wartime needs. Mauchly, a professor of electrical engineering at the University of Pennsylvania and Eckert, one of his graduate students, proposed to build a general-purpose computer using vacuum tubes. In 1943, this proposal was accepted by the Army, and work began on the ENIAC.

The resulting machine was enormous, weighting 30 tons, occupying 15,000 square feet of floor space, and containing more than 18,000 vacuum tubes. When operating, it consumed 140 kilowatts of power.

It was alos substantially faster than any electronic-mechanical computer, being capable of 5000 additions per second. The ENIAC was decimal rather than a binary machine. That is, numbers were represented in decimal form and arithmetic was performed in the decimal system. Its memory consisted of 20 "accumulators", each capable of holding a 10-digit decimal number.

Each digit was represented by a ring of 10 vacuum tubes. At any time, only one vacuum tube was in the ON state, representing one of the 10 digits.

The major drawback of the ENIAC was that it had to be programmed manually by setting switches and plugging and unplugging cables.

The ENIAC was completed in 1946, too late to be used in the war effort. Instead, its first task was to perform a series of complex calculations that were used to help determine

the feasibility of the H-bomb. The ENIAC continued to be used until 1955.

Von Neumann Machine

The programming process could be facilitated if the program could be represented in a form suitable for storing in memory alongside the data. Then, a computer could get its instructions by reading them from memory, and a program could be set or altered by setting the values of a portion of memory.

This idea, known as the Stored-program concept, is usually attributed to the ENIAC designers, most notably the mathematician John von Neumann, who was a consultant on the ENIAC project.

The idea was also developed at about the same time by Turing. The first publication of the idea was in a 1945 proposal by von Neumann for a new computer, the EDVAC (Electronic Discrete Variable Computer).

In 1946, von Neumann and his colleagues began the design of a new stored-program computer, referred to as the IAS computer, at the Princeton Institute for Advanced Studies.

The IAS computer, although not completed until 1952, is the prototype of all subsequent general-purpose computers. Figure shows the general structure of the IAS computer.

It consists of:

- A main memory, which stores both data and instructions.
- An arithmetic-logical unit (ALU) capable of operating on binary data.

- A control unit, which interprets the instructions in memory and causes them to be executed.
- Input and output (I/O) equipment operated by the control unit.

Commercial Computers

The 1950s saw the birth of the computer industry with two companies, Sperry and IBM, dominating the marketplace. In 1947, Eckert and Mauchly formed the Eckert-Mauchly computer Corporation to manufacture computers commercially. Their first successful machine was the UNIVAC I (Universal Automatic Computer), which was commissioned by the Bureau of the Census for the 1950 calculations. The Eckert-Mauchly Computer Corporation became part of the UNIVAC division of Sperry-Rand Corporation, which went on to build a series of successor machines.

The UNIVAC II, which had greater memory capacity and higher performance than the UNIVAC I, was delivered in the late 1950s and illustrates several trends that have remained characteristic of the computer industry. First, advances in technology allow companies to continue to build larger, more powerful computers. Second, each company tries to make its new machines upward compatible with the older machines. This means that the programs written for the older machines can be executed on the new machine. This strategy is adopted in the hopes of retaining the customer base; that is, when a customer decides to buy a newer machine, he is likely to get it from the same company to avoid losing the investment in programs.

The UNIVAC division also began development of the 1100 series of computers, which was to be its bread and butter. This series illustrates a distinction that existed at one time. The first model, the UNIVAC 1103, and its successors for many years were primarily intended for scientific applications, involving long and complex calculations. Other companies concentrated on business applications, which involved processing large amounts of text data. This split has largely disappeared but it was evident for a number of years.

IBM, which was then the major manufacturer of punched-card processing equipment, delivered its first electronic stored-program computer, the 701, in 1953. The 701 was intended primarily for scientific applications. In 1955, IBM introduced the companion 702 product, which had a number of hardware features that suited it to business applications. These were the first of a long series of 700/7000 computers that established IBM as the overwhelmingly dominant computer manufacturer.

Second Generation: Transistors

The first major change in the electronic computer came with the replacement of the vacuum tube by the transistor. The transistor is smaller, cheaper, and dissipates less heat than a vacuum tube but can be used in the same way as a vacuum tube to construct computers.

Unlike the vacuum tube, which requires wires, metal plates, a glass capsule, and a vacuum, the transistor is a solid-state device, made from silicon. The transistor was invented at Bell Labs in 1947 and by the 1950s had launched an electronic revolution. It was not until the late 1950s,

however, that fully transistorized computers were commercially available. IBM again was not the first company to deliver the new technology. NCR and, more successfully, RCA were the front-runners with some small transistor machines. IBM followed shortly with the 7000 series.

The use of the transistor defines the second generation of computers. It has become widely accepted to classify computers into generations based on the fundamental hardware technology employed. Each new generation is characterized by greater processing performance, larger memory capacity, and smaller size than the previous one.

Third Generation: Integrated Circuits

A single, self-contained transistor is called a discrete component. Throughout the 1950s and early 1960s, electronic equipment was composed largely of discrete components—transistors, resistors, capacitors, and so on.

Discrete components were manufactured separately, packaged in their own containers, and soldered or wired together onto circuit boards, which were then installed in computers, oscilloscopes, and other electronic equipment. Whenever an electronic device called for a transistor, a little tube of metal containing a pinhead-sized piece of silicon had to be soldered to a circuit board. The entire manufacturing process, from transistor to circuit board, was expensive and cumbersome.

These facts of life were beginning to create problems in the computer industry. Early second-generation computers contained about 10,000 transistors.

This figure grew to the hundreds of thousands, making the manufacture of newer, more powerful machines increasingly difficult. In 1958 came the achievement that revolutionized electronics and started the era of microelectronics: the invention of the integrated circuit. It is the integrated circuit that defines the third generation of computers. Perhaps the two most important members of the third generation are the IBM System/360 and the DEC PDP-8.

Later Generations

Beyond the third generation there is less general agreement on defining generations of computers. There have been a fourth and a fifth generation, based on advances in integrated circuit technology. With the introduction of large-scale integration (LSI), more than 1000,000 components can be placed on a single integrated circuit chip. Very-large-scale integration (VLSI) achieved more than 1000,000,000 components per chip, and current VLSI chips can contain more than 1000.000 components.

Computers Classification and Data Processing

Computers are classified according to their data processing speed, amount of data that they can hold and price. Generally, a computer with high processing speed and large internal storage is called a big computer. Due to rapidly improving technology, we are always confused among the categories of computers.

Depending upon their speed and memory size, computers are classified into following four main groups:

- Supercomputer.

- Mainframe computer.
- Mini computer.
- Microcomputer.

Supercomputer

Supercomputer is the most powerful and fastest, and also very expensive. It was developed in 1980s. It is used to process large amount of data and to solve the complicated scientific problems. It can perform more than one trillions calculations per second. It has large number of processors connected parallel.

So parallel processing is done in this computer. In a single supercomputer thousands of users can be connected at the same time and the supercomputer handles the work of each user separately.

Supercomputer are mainly used for:

- Weather forecasting.
- Nuclear energy research.
- Aircraft design.
- Automotive design.
- Online banking.
- To control industrial units.

The supercomputers are used in large organizations, research laboratories, aerospace centers, large industrial units etc. Nuclear scientists use supercomputers to create and analyse models of nuclear fission and fusions, predicting the actions and reactions of millions of atoms as they interact. The examples of supercomputers are CRAY-1, CRAY-2, Control Data CYBER 205 and ETA A-10 etc.

Mainframe Computers

Mainframe computers are also large-scale computers but supercomputers are larger than mainframe. These are also very expensive. The mainframe computer specially requires a very large clean room with air-conditioner.

This makes it very expensive to buy and operate. It can support a large number of various equipments. It also has multiple processors. Large mainframe systems can handle the input and output requirements of several thousand of users. For example, IBM, S/390 mainframe can support 50,000 users simultaneously. The users often access then mainframe with terminals or personal computers. Tere are basically two types of terminals used with mainframe systems.

Dumb Terminal

Dumb terminal does not have its own CPU and storage devices. This type of terminal uses the CPU and storage devices of mainframe system. Typically, a dumb terminal consists of monitor and a keyboard (or mouse).

Intelligent Terminal

Intelligent terminal has its own processor and can perform some processing operations. Usually, this type of terminal does not have its own storage. Typically, personal computrers are used as intelligent terminals. A personal computer as an intelligent terminal gives facility to access data and other services from mainframe system. It also enables to store and process data locally. The mainframe computers are specially used as servers on the World Wide Web. The mainframe computers are used in large organizations such as Banks,

Airlines and Universities etc. where many people (users) need frequent access to the same data, which is usually organized into one or more huge databases. IBM is the major manufacturer of mainframe computers. The examples of mainframes are IBM S/390, Control Data CYBER 176 and Amdahl 580 etc.

Minicomputers

These are smaller in size, have lower processing speed and also have lower cost than mainframe. These computers are known as minicomputers because of their small size as compared to other computers at that time. The capabilities of a minicomputer are between mainframe and personal computer. These computers are also known as midrange computers.

The minicomputers are used in business, education and many other government departments. Although some minicomputers are designed for a single user but most are designed to handle multiple terminals. Minicomputers are commonly used as servers in network environment and hundreds of personal computers can be connected to the network with a minicomputer acting as server like mainframes, minicomputers are used as web servers. Single user minicomputers are used for sophisticated design tasks.

The first minicomputer was introduced in the mid-1960s by Digital Equipment Corporation (DEC). After this IBM Corporation (AS/400 computers) Data General Corporation and Prime Computer also designed the mini computers.

Microcomputer

The microcomputers are also known as personal computers

or simply PCs. Microprocessor is used in this type of computer. These are very small in size and cost. The IBM's first microcomputer was designed in 1981 and was named as IBM-PC. After this many computer hardware companies copied the design of IBM-PC. The term "PC-compatible" refers any personal computer based on the original IBM personal computer design.

The most popular types of personal computers are the PC and the Apple. PC and PC-compatible computers have processors with different architectures than processors in Apple computers. These two types of computers also use different operating systems. PC and PC-compatible computers use the Windows operating system while Apple computers use the Macintosh operating system (MacOS). The majority of microcomputers sold today are part of IBM-compatible. However the Apple computer is neither an IBM nor a compatible. It is another family of computers made by Apple computer. Personal computers are available in two models.

These are:

- Desktop PCs
- Tower PCs

A desktop personal computer is most popular model of personal computer. The system unit of the desktop personal computer can lie flat on the desk or table. In desktop personal computer, the monitor is usually placed on the system unit. Another model of the personal computer is known as tower personal computer. The system unit of the tower PC is vertically placed on the desk of table. Usually the system unit of the tower model is placed on the floor to make desk

space free and user can place other devices such as printer, scanner etc. on the desktop. Today computer tables are available which are specially designed for this purpose. The tower models are mostly used at homes and offices.

Microcomputer are further divided into following categories:

- Laptop computer
- Workstation
- Network computer
- Handheld computer

Laptop Computer

Laptop computer is also known as notebook computer. It is small size (85-by-11 inch notebook computer and can fit inside a briefcase. The laptop computer is operated on a special battery and it does not have to be plugged in like desktop computer. The laptop computer is portable and fully functional microcomputer.

It is mostly used during journey. It can be used on your lap in an airplane. It is because it is referred to as laptop computer. The memory and storage capacity of laptop computer is almost equivalent to the PC or desktop computer. It also has the hard dist, floppy disk drive, Zip disk drive, CD-ROM drive, CD-writer etc. it has built-in keyboard and built-in trackball as pointing device.

Laptop computer is also available with the same processing speed as the most powerful personal computer. It means that laptop computer has same features as personal computer. Laptop computers are more expensive than desktop computers. Normally these computers are frequently used in business travellers.

Workstations

Workstations are special single user computers having the same features as personal computer but have the processing speed equivalent to minicomputer or mainframe computer. A workstation computer can be fitted on a desktop. Scientists, engineers, architects and graphic designers mostly use these computers.

Workstation computers are expensive and powerful computers. These have advanced processors, more RAM and storage capacity than personal computers. These are usually used as single-user applications but these are used as servers on computer network and web servers as well.

Network Computers

Network computers are also version of personal computers having less processing power, memory and storage. These are specially designed as terminals for network environment. Some types of network computers have no storage. The network computers are designed for network, Internet or Intranet for data entry or to access data on the network. The network computers depend upon the network's server for data storage and to use software. These computers also use the network's server to perform some processing tasks. In the mid-1990s the concept of network computers became popular among some PC manufacturers. As a result several variations of the network computers quickly became available.

In business, variations of the network computer are Windows terminals, NetPCs and diskless workstations. Some network computers are designed to access only the Internet

or to an Intranet. These devices are sometimes called Internet PCs, Internet boxes etc. In home some network computers do not include monitor. These are connected to home television, which serves as the output devices. A popular example of a home-based network computer is Web TV, which enables the user to connect a television to the Internet.

The Web TV has a special set-top box used to connect to the Internet and also provides a set of simple controls which enable the user to navigate the Internet, send and receive e-mails and to perform other tasks on the network while watching television. Network computers are cheaper to purchase and to maintain than personal computers.

Handheld Computer

In the mid 1990s, many new types of small personal computing devices have been introduced and these are referred to as handheld computers. These computers are also referred to as Palmtop Computers. The handheld computers sometimes called Mini-Notebook Computers. The type of computer is named as handheld computer because it can fit in one hand while you can operate it with the other hand. Because of its reduced size, the screen of handheld computer is quite small. Similarly it also has small keyboard. The handheld computers are preferred by business traveller. Some handheld computers have a specialized keyboard. These computers are used by mobile employees, such as meter readers and parcel delivery people, whose jobs require them to move from place to place.

The examples of handheld computers are:

- Personal digital assistance

- Cellular telephones
- H/PC pro devices

Personal Digital Assistance (PDAs)

The PDA is one of the more popular lightweight mobile devices in use today. A PDA provides special functions such as taking notes, organizing telephone numbers and addresses. Most PDAs also offer a variety of other application software such as word processing, spreadsheet and games etc. Some PDAs include electronic books that enable users to read a book on the PDA's screen. Many PDAs are web-based and users can send/receive e-mails and access the Internet. Similarly, some PDAs also provide telephone capabilities.

The primary input device of a PDA is the stylus. A stylus is an electronic pen and looks like a small ballpoint pen. This input device is used to write notes and store in the PDA by touching the screen. Some PDAs also support voice input.

Cellular Phones

A cellular phone is a web-based telephone having features of analog and digital devices. It is also referred to as Smart Phone. In addition to basic phone capabilities, a cellular phone also provides the functions to receive and send e-mails & faxes and to access the Internet.

H/PC Pro Devices

H/PC Pro dive is new development in handheld technology. These systems are larger than PDAs but they are not quite as large as typical notebook PCs. These devices have features between PDAs and notebook PCs. The H/PC Pro device

Computer Memory System

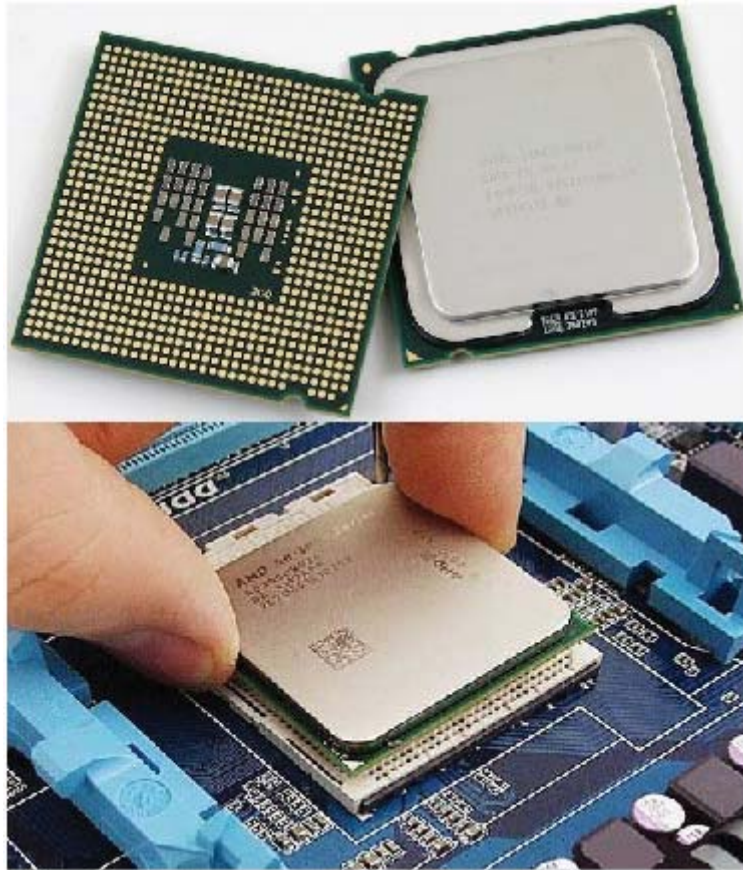
includes a full-size keyboard but it does not include disk. These systems also have RAM with very low storage capacity and slow speed of processor.

3

Main Units: Central Processing Unit

A central processing unit (CPU) or processor is an electronic circuit that can execute computer programmes. This broad definition can easily be applied to many early computers that existed long before the term “CPU” ever came into widespread usage. The term itself and its initialism have been in use in the computer industry at least since the early 1960s (Weik 1961). The form, design and implementation of CPUs have changed dramatically since the earliest examples, but their fundamental operation has remained much the same. Early CPUs were custom-designed as a part of a larger, sometimes one-of-a-kind, computer.

Computer Memory System



However, this costly method of designing custom CPUs for a particular application has largely given way to the development of mass-produced processors that are made for one or many purposes.

This standardization trend generally began in the era of discrete transistor mainframes and minicomputers and has rapidly accelerated with the popularization of the integrated circuit (IC). The IC has allowed increasingly complex CPUs to be designed and manufactured to tolerances on the order of nanometers. Both the miniaturization and standardization of CPUs have increased the presence of these digital devices in modern life far beyond the limited application of dedicated

computing machines. Modern microprocessors appear in everything from automobiles to cell phones to children's toys.

History of CPU

Prior to the advent of machines that resemble today's CPUs, computers such as the ENIAC had to be physically rewired in order to perform different tasks. These machines are often referred to as "fixed-programme computers," since they had to be physically reconfigured in order to run a different programme. Since the term "CPU" is generally defined as a software (computer programme) execution device, the earliest devices that could rightly be called CPUs came with the advent of the stored-programme computer.

The idea of a stored-programme computer was already present during ENIAC's design, but was initially omitted so the machine could be finished sooner. On June 30, 1945, before ENIAC was even completed, mathematician John von Neumann distributed the paper entitled "First Draft of a Report on the EDVAC." It outlined the design of a stored-programme computer that would eventually be completed in August 1949 (von Neumann 1945). EDVAC was designed to perform a certain number of instructions (or operations) of various types. These instructions could be combined to create useful programmes for the EDVAC to run. Significantly, the programmes written for EDVAC were stored in high-speed computer memory rather than specified by the physical wiring of the computer. This overcame a severe limitation of ENIAC, which was the large amount of time and effort it took to reconfigure the computer to perform a new task. With von Neumann's design, the programme, or

software, that EDVAC ran could be changed simply by changing the contents of the computer's memory.

While von Neumann is most often credited with the design of the stored-programme computer because of his design of EDVAC, others before him such as Konrad Zuse had suggested similar ideas. Additionally, the so-called Harvard architecture of the Harvard Mark I, which was completed before EDVAC, also utilized a stored-programme design using punched paper tape rather than electronic memory. The key difference between the von Neumann and Harvard architectures is that the latter separates the storage and treatment of CPU instructions and data, while the former uses the same memory space for both. Most modern CPUs are primarily von Neumann in design, but elements of the Harvard architecture are commonly seen as well.

Being digital devices, all CPUs deal with discrete states and therefore require some kind of switching elements to differentiate between and change these states.

Prior to commercial acceptance of the transistor, electrical relays and vacuum tubes (thermionic valves) were commonly used as switching elements. Although these had distinct speed advantages over earlier, purely mechanical designs, they were unreliable for various reasons.

For example, building direct current sequential logic circuits out of relays requires additional hardware to cope with the problem of contact bounce. While vacuum tubes do not suffer from contact bounce, they must heat up before becoming fully operational and eventually stop functioning altogether. Usually, when a tube failed, the CPU would have to be diagnosed to locate the failing component so it could

be replaced. Therefore, early electronic (vacuum tube based) computers were generally faster but less reliable than electromechanical (relay based) computers.

Tube computers like EDVAC tended to average eight hours between failures, whereas relay computers like the (slower, but earlier) Harvard Mark I failed very rarely (Weik 1961:238). In the end, tube based CPUs became dominant because the significant speed advantages afforded generally outweighed the reliability problems. Most of these early synchronous CPUs ran at low clock rates compared to modern microelectronic designs (see below for a discussion of clock rate). Clock signal frequencies ranging from 100 kHz to 4 MHz were very common at this time, limited largely by the speed of the switching devices they were built with.

Discrete Transistor and IC CPUs

The design complexity of CPUs increased as various technologies facilitated building smaller and more reliable electronic devices. The first such improvement came with the advent of the transistor. Transistorized CPUs during the 1950s and 1960s no longer had to be built out of bulky, unreliable, and fragile switching elements like vacuum tubes and electrical relays. With this improvement more complex and reliable CPUs were built onto one or several printed circuit boards containing discrete (individual) components.

During this period, a method of manufacturing many transistors in a compact space gained popularity. The integrated circuit (IC) allowed a large number of transistors to be manufactured on a single semiconductor-based die, or “chip.”

At first only very basic non-specialized digital circuits such as NOR gates were miniaturized into ICs. CPUs based upon these “building block” ICs are generally referred to as “small-scale integration” (SSI) devices.

SSI ICs, such as the ones used in the Apollo guidance computer, usually contained transistor counts numbering in multiples of ten. To build an entire CPU out of SSI ICs required thousands of individual chips, but still consumed much less space and power than earlier discrete transistor designs. As microelectronic technology advanced, an increasing number of transistors were placed on ICs, thus decreasing the quantity of individual ICs needed for a complete CPU. MSI and LSI (medium- and large-scale integration) ICs increased transistor counts to hundreds, and then thousands. In 1964 IBM introduced its System/360 computer architecture which was used in a series of computers that could run the same programmes with different speed and performance.

This was significant at a time when most electronic computers were incompatible with one another, even those made by the same manufacturer. To facilitate this improvement, IBM utilized the concept of a microprogram (often called “microcode”), which still sees widespread usage in modern CPUs (Amdahl *et al.* 1964). The System/360 architecture was so popular that it dominated the mainframe computer market for the decades and left a legacy that is still continued by similar modern computers like the IBM zSeries. In the same year (1964), Digital Equipment Corporation (DEC) introduced another influential computer aimed at the scientific and research markets, the PDP-8.

DEC would later introduce the extremely popular PDP-11 line that originally was built with SSI ICs but was eventually implemented with LSI components once these became practical. In stark contrast with its SSI and MSI predecessors, the first LSI implementation of the PDP-11 contained a CPU composed of only four LSI integrated circuits (Digital Equipment Corporation 1975).

Transistor-based computers had several distinct advantages over their predecessors. Aside from facilitating increased reliability and lower power consumption, transistors also allowed CPUs to operate at much higher speeds because of the short switching time of a transistor in comparison to a tube or relay. Thanks to both the increased reliability as well as the dramatically increased speed of the switching elements (which were almost exclusively transistors by this time), CPU clock rates in the tens of megahertz were obtained during this period. Additionally while discrete transistor and IC CPUs were in heavy usage, new high-performance designs like SIMD (Single Instruction Multiple Data) vector processors began to appear. These early experimental designs later gave rise to the era of specialized supercomputers like those made by Cray Inc.

Microprocessors

The introduction of the microprocessor in the 1970s significantly affected the design and implementation of CPUs. Since the introduction of the first microprocessor (the Intel 4004) in 1970 and the first widely used microprocessor (the Intel 8080) in 1974, this class of CPUs has almost completely

overtaken all other central processing unit implementation methods. Mainframe and minicomputer manufacturers of the time launched proprietary IC development programmes to upgrade their older computer architectures, and eventually produced instruction set compatible microprocessors that were backward-compatible with their older hardware and software. Combined with the advent and eventual vast success of the now ubiquitous personal computer, the term “CPU” is now applied almost exclusively to microprocessors.

Previous generations of CPUs were implemented as discrete components and numerous small integrated circuits (ICs) on one or more circuit boards. Microprocessors, on the other hand, are CPUs manufactured on a very small number of ICs; usually just one. The overall smaller CPU size as a result of being implemented on a single die means faster switching time because of physical factors like decreased gate parasitic capacitance. This has allowed synchronous microprocessors to have clock rates ranging from tens of megahertz to several gigahertz. Additionally, as the ability to construct exceedingly small transistors on an IC has increased, the complexity and number of transistors in a single CPU has increased dramatically. This widely observed trend is described by Moore’s law, which has proven to be a fairly accurate predictor of the growth of CPU (and other IC) complexity to date.

While the complexity, size, construction, and general form of CPUs have changed drastically over the past sixty years, it is notable that the basic design and function has not changed much at all. Almost all common CPUs today can

be very accurately described as von Neumann stored-programme machines. As the aforementioned Moore's law continues to hold true, concerns have arisen about the limits of integrated circuit transistor technology. Extreme miniaturization of electronic gates is causing the effects of phenomena like electromigration and subthreshold leakage to become much more significant. These newer concerns are among the many factors causing researchers to investigate new methods of computing such as the quantum computer, as well as to expand the usage of parallelism and other methods that extend the usefulness of the classical von Neumann model.

CPU Operation

The fundamental operation of most CPUs, regardless of the physical form they take, is to execute a sequence of stored instructions called a programme. The programme is represented by a series of numbers that are kept in some kind of computer memory. There are four steps that nearly all CPUs use in their operation: fetch, decode, execute, and writeback.

The first step, fetch, involves retrieving an instruction (which is represented by a number or sequence of numbers) from programme memory. The location in programme memory is determined by a programme counter (PC), which stores a number that identifies the current position in the programme. In other words, the programme counter keeps track of the CPU's place in the current programme. After an instruction is fetched, the PC is incremented by the length of the instruction word in terms of memory units. Often the

instruction to be fetched must be retrieved from relatively slow memory, causing the CPU to stall while waiting for the instruction to be returned. This issue is largely addressed in modern processors by caches and pipeline architectures.

The instruction that the CPU fetches from memory is used to determine what the CPU is to do. In the decode step, the instruction is broken up into parts that have significance to other portions of the CPU. The way in which the numerical instruction value is interpreted is defined by the CPU's instruction set architecture (ISA). Often, one group of numbers in the instruction, called the opcode, indicates which operation to perform. The remaining parts of the number usually provide information required for that instruction, such as operands for an addition operation. Such operands may be given as a constant value (called an immediate value), or as a place to locate a value: a register or a memory address, as determined by some addressing mode. In older designs the portions of the CPU responsible for instruction decoding were unchangeable hardware devices. However, in more abstract and complicated CPUs and ISAs, a microprogram is often used to assist in translating instructions into various configuration signals for the CPU. This microprogram is sometimes rewritable so that it can be modified to change the way the CPU decodes instructions even after it has been manufactured.

After the fetch and decode steps, the execute step is performed. During this step, various portions of the CPU are connected so they can perform the desired operation. If, for instance, an addition operation was requested, an arithmetic logic unit (ALU) will be connected to a set of inputs

and a set of outputs. The inputs provide the numbers to be added, and the outputs will contain the final sum. The ALU contains the circuitry to perform simple arithmetic and logical operations on the inputs (like addition and bitwise operations). If the addition operation produces a result too large for the CPU to handle, an arithmetic overflow flag in a flags register may also be set.

The final step, writeback, simply “writes back” the results of the execute step to some form of memory. Very often the results are written to some internal CPU register for quick access by subsequent instructions. In other cases results may be written to slower, but cheaper and larger, main memory. Some types of instructions manipulate the programme counter rather than directly produce result data. These are generally called “jumps” and facilitate behavior like loops, conditional programme execution (through the use of a conditional jump), and functions in programmes. Many instructions will also change the state of digits in a “flags” register. These flags can be used to influence how a programme behaves, since they often indicate the outcome of various operations. For example, one type of “compare” instruction considers two values and sets a number in the flags register according to which one is greater. This flag could then be used by a later jump instruction to determine programme flow.

After the execution of the instruction and writeback of the resulting data, the entire process repeats, with the next instruction cycle normally fetching the next-in-sequence instruction because of the incremented value in the programme counter. If the completed instruction was a jump,

the programme counter will be modified to contain the address of the instruction that was jumped to, and programme execution continues normally. In more complex CPUs than the one described here, multiple instructions can be fetched, decoded, and executed simultaneously. This section describes what is generally referred to as the “Classic RISC pipeline,” which in fact is quite common among the simple CPUs used in many electronic devices (often called microcontroller). It largely ignores the important role of CPU cache, and therefore the access stage of the pipeline.

ADMINISTRATIVE AND COMPUTER USE

The computing facilities of the Chabot-Las Positas Community College District are provided for the use of students, faculty, and staff in support of the programmes of the Colleges and District. In order to facilitate proper and responsible use of computers, the following administrative rules and procedures are established for all users. Instructors, managers, departments, or colleges may elect to impose additional requirements or restrictions. Beyond the consequences listed herein, rule violations may have consequences determined by District Board policy and applicable law.

PROPER USE

- a. Board Policy 2311 specifies that the computer systems of the District are provided solely for the following purposes:
 - 1) use by authorized employees and agents of the Chabot-Las Positas Community College District for District business;

Computer Memory System

- 2) use by authorized employees of the Chabot-Las Positas Community College District for professional activities related to the employee's job function, or
 - 3) use by registered students or authorized employees of the Chabot-Las Positas Community College District for instructional activities; or
 - 4) public access to approved District or College information resources via the public telephone and data networks.
- b. Use of District computer resources for personal or recreational purposes is prohibited. Prohibited activities include, but are not limited to, the following examples:
- storing personal recipes
 - balancing your personal checkbook
 - preparing a homeowner's association newsletter
 - playing any sort of computer games unless the games are a specific component of an instructional activity or assignment.
- c. Use of District computer resources for personal gain, profit, or commercial purposes is prohibited. Prohibited activities include, but are not limited to, the following examples:
- consulting for profit
 - typing services for profit
 - maintaining commercial business records

Computer Memory System

- developing software for sale, except as permitted in Board Policy pertaining to intellectual property rights
 - any activity which is not District business or a professional activity related to the employee's job function.
- d. Use of District computer resources for unauthorized activities is prohibited. Unauthorized activities include, but shall not be limited to, the following examples:
- use of passwords or accounts of another user
 - attempts to capture or "crack" passwords
 - attempts to break encryption protocols
 - attempts to use loopholes in computer security or special passwords to gain access to systems, obtain extra resources, or make unauthorized use of systems
 - destruction or unauthorized alteration of data belonging to the District or to another user
 - creation or communication of "viruses", "worms", or "Trojan horses"
 - acts that restrict access to the system or damage the system
 - acts that deliberately misrepresent the identity of the source of a message
 - acts that harass, threaten, or defame other persons
 - acts that violate any law

Copyrights and Licenses

- a. The District acquires a substantial portion of its computer software from vendors under license agreements which restrict the use of the software to specific computer systems and which require the District to limit the use and copying of the software. Board Policy 2311 requires compliance with the terms of these licenses and with copyright law.

Use of District computer resources in violation of copyright restrictions or software license terms is prohibited under Board Policy 2311. Prohibited activities include, but shall not be limited to, the following examples:

- copying District-licensed software in violation of the license terms or copyright law
 - installing software on District computers in violation of the license terms or copyright law
 - "giving" District software to students or colleagues
- b. Each major organization shall be responsible for implementation of this policy: For computer software used on College computers, the College Presidents shall be responsible for establishing implementation procedures.

For computer software used on District-organizational unit computers, the Chief Management Information Officer shall be responsible for establishing implementation procedures.

System Access

- a. *Administrative Systems:* The District's administrative

Computer Memory System

systems are operated by MIS. Access to these systems requires MIS approval of a written request prepared by the employee's supervisor or manager. In addition, other administrative review is sometimes required. For example, the Controller will review the need for Finance System access. Usually, requests will be approved for staff who have specific administrative responsibilities requiring system access. Administrative responsibilities that require system access include, but are not limited to, the following examples:

- management or overseeing of department or area budgets
- management of financial records of special projects or grants
- data entry of information pertaining to students, personnel, or finance records
- student information inquiry by counselors or A&R staff

Administrative system users shall access only those system accounts authorized by MIS.

All other access to administrative systems is prohibited. Administrative system users may not, under any circumstances, transfer or confer their system access privileges to another individual or permit use of their assigned system accounts by another individual. Users will be held responsible for all administrative system transactions conducted under their login passwords.

Administrative system users will be granted access privileges only if they agree in writing to adhere to the rules

and procedures presented in this section. System access privileges may be revoked without notice in response to violations of these rules and procedures or in response to legitimate requests from the employee's supervisor or manager.

- b. *Instructional Systems:* The District's instructional systems are owned and operated by the Colleges. (The sole exception is the instructional Sequent computer, which is operated by MIS.) Access and privileges for these systems are assigned by the systems administrators of specific individual systems. Eligible individuals may become authorized users of a system and be granted appropriate access and privileges by following the approval steps prescribed by the College for that system.

Passwords

- a. Passwords are the keys to system security, and they provide the most important defense against unauthorized use of District systems.

Each system user is responsible to:

- follow certain rules when creating passwords
- select passwords that are secure
- change login passwords periodically
- keep passwords secret

Users shall fulfill these responsibilities in conformity with established CLPCCD Password Guidelines.

- b. Users of a terminal or PC that is logged in to an administrative system must not leave it unattended. Users will be held responsible for all system transactions conducted under their login passwords.

Ownership

- a. The District's computer systems, including hardware, software, and all computerized information and data are owned by the District or are licensed from vendors under license agreements. Except as provided in Board Policy pertaining to intellectual property rights, employees and students have no rights of ownership to these systems or to the information they contain, even if the employee or student entered the information into these systems. Employees may use this information only as directed in the legitimate business of the Colleges and District and only as prescribed by Board Policy 5511.

Electronic Mail Privacy

- a. Under Board Policy 2311, the District's electronic mail system and messages are owned by the District and provided for legitimate business use by its employees.
- b. The District's E-mail system uses encrypted messages and is relatively secure. It is contrary to MIS department policy for MIS staff to snoop or routinely examine the contents of employee E-mail, and most E-mail messages will enjoy private status.

Nevertheless, E-mail messages are not guaranteed to be private or confidential, and the District accepts no responsibility for consequences that might arise from disclosure of an E-mail message. Please remember that control over a message is lost once it

is sent, and future events may have unanticipated results:

- The recipient of the message might forward it to others on the system.
- The recipient of the message might print it and hand it to another reader or might even post it on the wall.
- The message might accidentally be sent to an unintended recipient, especially when using named Groups for the "TO" address.
- In unusual circumstances, MIS staff might need to examine mail in order to resolve a system problem.
- Conceivably, one or more messages might be subpoenaed in a legal proceeding, and then MIS would be required to provide the subpoenaed material.

Etiquette

Users are expected to use the system in a manner that reflects respect for other users.

- a. It is a violation of system etiquette to transmit material which is offensive, harassing, or needlessly affects the work of other users.
- b. Please carefully consider the appropriateness of any E-mail message being sent to EVERYONE; notification of the arrival of such a message will interrupt every user on the system and consume a portion of their system resources. Such messages are sometimes perceived as the electronic equivalent of "junk mail".

- c. Mail messages composed in all capitals are difficult to read and are often perceived as the electronic equivalent of "SHOUTING". Please use such messages sparingly.

Nondiscrimination

- a. All users have the right to be free from any conduct associated with the use of District computer systems which discriminates against any person on the basis of race, colour, national origin, gender, or disability. Users of District systems shall refrain from such discriminatory acts.
- b. Discriminatory conduct includes, but is not limited to, written or graphic conduct that satisfies both of the following conditions: (a) harasses, denigrates, or shows hostility to or aversion toward an individual or group based on race, colour, national origin, gender, or disability, and (b) has the purpose or effect of creating a hostile, intimidating, or offensive educational environment.
 - 1) "Harassing conduct" includes, but is not limited to, epithets, slurs, negative stereotyping, or threatening, intimidating, or hostile acts, that relate to race, colour, national origin, gender, or disability. This includes acts that purport to be "jokes" or "pranks" but that are hostile or demeaning.
 - 2) A "hostile educational environment" is established when harassing conduct is sufficiently severe, pervasive, or persistent so

as to interfere with or limit the ability of an individual to participate in or benefit from District computing systems.

- c. Any user who believes he or she has been subject to conduct associated with the use of District computer systems which discriminates on the basis of race, colour, national origin, gender, or disability may report the incident to the College or District Affirmative Action/ Harassment Officer.

MIS Staff Rights and Responsibilities

- a. In the normal course of systems administration, the MIS staff occasionally may need to examine files, electronic mail, and printer output in order to gather sufficient information to diagnose and correct system problems or perform technical maintenance. In the course of this work, the staff reserves the right to inspect, copy, remove, or otherwise alter any data, file, or system resources which may adversely affect the system without notice to the user. In addition, the MIS staff reserves the right to restrict system access of any user who violates the rules/procedures presented in this section.
- b. Although MIS staff have the right to examine any system files, they also have a responsibility to maintain users' privacy to the maximum extent possible.

User Rights and Responsibilities

- a. As described herein, users of District systems have the right to

Computer Memory System

- use District systems as authorized
 - own information stored on District systems solely as provided in Board policy pertaining to intellectual property rights
 - be free of routine intrusions on privacy
 - be free of discrimination in use of District systems.
- b. As described herein, users of District systems have the responsibility to:
- use the systems in compliance with the rules and procedures presented in this section
 - make proper use of District systems
 - comply with copyright law
 - access systems only as authorized
 - keep passwords secret and maintain password security
 - use the system with proper etiquette and respect for other users
 - refrain from acts that are discriminatory, defamatory, harassing, or illegal
 - agree that the District is not responsible for the content of external networks and for actions by individual users of the systems in violation of these rules.

Management Rights and Responsibilities

- a. Administrative Systems Managers shall make written requests for employees' access to the District's administrative systems. In addition, to maintain

system security, managers shall notify MIS in writing immediately when system access is no longer required or authorized for an employee. Managers shall be responsible to provide general supervision of departmental employees' adherence to the rules and procedures presented herein, and managers shall have the right to impose additional departmental rules or procedures. In the event of conflict, the rules and procedures presented herein shall take precedence over departmental rules and procedures.

IMPORTANCE OF THE COMPUTER SCIENCES

Study of Computer Sciences has been given due place in our school education programme by being made as a compulsory subject. Not only that more and more emphasis is now being paid over the scientific and technical education. By doing so indeed a right step has been taken to push our country forward and to enable us to compete with other progressive nations. It has necessitated to lay due emphasis on the Computer Science Education & Teaching right from the primary stage. Realizing such need Kothari Commission has very rightly remarked in their recommendations as follows:

“Science and Mathematics should be taught on compulsory basis of all pupils as a part of general education during the first ten years of schooling”.

PLACE IN THE SCHOOL CURRICULUM

The question here may arise that why the subject science including Computer Science has been given so much weight as to be taught compulsorily during the ten years of one's

school life. The answer may be satisfactorily discovered in its various uses, applications and advantages in the following pages.

UTILITARIAN VALUE OF DAY-TO-DAY USE

Modern age is science age. We see a network of scientific gadgets based on latest scientific inventions all around us in the field of Computer Science. There was a time in olden age when we used to burn an earth ware oil lamp for light but today, we have electrical bulbs, tube lights illuminating as bright as sun. Olden ways of transport on horses and camels have given way to electric trains, cars and aeroplanes. We have newspapers, radios televisions, cinemas applications in their places. On looking at the Moon, we used to wonder what that old woman is doing with that spinning wheel. Today man has reached the moon, stepped on it and is busy in exploring it for further details. There is a miraculous change labour saving gadgets and machines. They have made all our complicated work not only easy but also these are done much faster. Science has revolutionized our way of living. Now our lives depend on scientifically invented gadgets so much that we cannot do without them. It is now imperative for everyone not only to understand science but also to master it from all angles Herbert Spencer has rightly said, “The knowledge gained through science is much useful in guiding our life style than gained through other sources”.

INTELLECTUAL VALUE

The study of Computer Science provides us the opportunity of developing our mental faculties of reasoning, imagination, memory, observation, concentration, analysis,

and originality and of systematic thinking. Science gives us the insight, which enables us to search the truth and the reality of nature around us. We are able to do useful except anything which we cannot prove by actual observation, reasoning and experimentations. The queries like 'why' and 'how' of all problems and phenomena can be satisfactorily answered only by the wisdom of science. A mature scientist can solve all problems with his shape intelligence and wisdom.

DISCIPLINARY VALUE

Study of Computer Sciences develops our personality as a whole. It inculcate spirit of enquiry, seriousness and systematic thinking. It brings about total transformation of one's viewpoint and makes thought process more organized. These sciences make us think seriously and help to observe the real nature of the problem. It helps us to judge all the good and bad points together, with the gain and loss likely to be incurred in the plan of action contemplated.

These sciences also check one to take a hasty action prompted by one's sentiments and influenced by the biased opinions of others, which are against the principles of science. Study of Computer Sciences promotes interest in study, concentration on habit of hard and systematic work. It also inculcates the habit of viewing a problem impartially with an alert mind. This helps to lead one's daily life successfully in well-organized and systematic way.

CULTURAL VALUE

Before we discuss the cultural value of Computer Sciences let us try to understand what is culture.

According to the great Indian poet Dinkar, “Culture is the way of life which has been handed over to society from one generation to another in the form of accumulated customs, habits and mode of living. The mode and style of living is different from one society to another and therefore their culture is also not the same”.

According to the above definition of culture, our activities like way of life, our food, customs, our views, our art and craft, scientific interests, social and economic conditions, all determine our standard of life. From time immemorial man has been trying to maintain and preserve their way of life and standard through the use of science. But somehow our way of life has been changing with the passage of time and progress of science. This change in our life-style is due to the inventions of computer science. The development of culture is the history of science. We can judge the progress of civilization and culture of a nation by its progress in science. Therefore study of Computer Sciences not only develops our culture but also helps in preserving it. We are constantly adjusting and modifying our style of living according to the latest scientific inventions and discoveries. Thus study of Computer Sciences develop culture of a society and nation.

AESTHETIC VALUE

Science is beauty, art, a source of entertainment and a successful means of attaining physical comforts. Can you imagine the joy and bliss of Graham Bell inventor of ‘telephone’, when his word were clearly heard on the other side by Watson? The words were, “Watson, come here, I want you”. When Watson heard these words, he rushed to the

room of Graham Bell shouting,” “I have clearly heard your message.” The joy and satisfaction they got from their successful invention was similar to the joy got by Archimedes from his discovery that he became totally oblivious of his naked body while coming out of the public bath and ran on the streets shouting, “Eureka, Eureka”. One gets extreme joy on one’s successful discovery or invention. Even the study of science is a source of great pleasure when one gets answers to his questions about the mysteries of nature. A child even experiences great joy when he gets answers to his queries like—‘Where do stars and moon go away during the day time?’ ‘Why does a hydrogen gas filled balloon rise up like a ball?’ “How is light produced in the room on pressing a switch”? etc. The child’s joy takes the form of this interest in studying Computer Sciences. Man derives boundless bliss on the revelation of the mysteries of nature, and the mysteries of greatest ‘Artist’ and His art.

Study of Computer Science help us to utilize our leisure purposefully. In addition to makes our lives happy by providing scientific gadgets for our entertainment in our spare time, cinema, television, radio, newspapers are gifts of Computer Sciences for our entertainment. Scientific hobbies can well be a source of joy in our leisure time. The useful and gainful hobbies related to Computer Sciences like photography, construction of scientific toys, manufacturing tooth paste, varnishes, inks etc. are not only interesting but are also financially useful.

MORAL VALUE

Some people believe that Science is responsible for lack of faith in God, but in reality the situation is reverse. No

doubt that science does not permit blind faith. It also does not admit faith in idol worship nor follows many useless customs and rituals. The pursuit of knowledge of nature of study of science cannot be called contrary to religion and faithlessness. Therefore the study of Computer Sciences and its pursuit not only include all the traits of morality but also develop them. The qualities of honesty of purpose, truth, justice, punctuality, determination, patience, self-control, self-respect, self-confidence and tolerance are automatically developed in man if he follows scientific methods in his pursuit of knowledge. The habit of lawfulness, the faculty of distinction between right and wrong and respect for other's point of view are some other desirable traits of character developed through the study of Computer Sciences. The person busy in pursuit of truth and reality imbibes in himself the qualities of morality. In Computer Sciences there is no place for self-opinionated person and his sentiments or democracies. Here every conclusion depends upon tests and actual observations and not by cheat and deceit. The ideal of "Truth is beauty" is always kept in view.

SOCIAL VALUE

Study of Computer Sciences provides great impetus and opportunities for the development and progress of the society as well as useful experiences of the society as well as useful experiences of the inculcation of social values and virtues in the students. It helps the students to know the material and objects surrounded in their physical and social environment. Harnessing the forces of nature for the welfare and progress of the society is the every essence of any invention and development put forward in the subject

Computer Sciences. Due to the all round progress and development put forward by the study of Computer Sciences the World Society has become now a global society. What we do in our present society is totally dependent upon the progress and development brought out by the study of Computer Sciences. It has given eyes to the society in the form of telescopes, microscopes, cameras and televisions, transport in the shape of fast electric trains and a variety of transport vehicles and wings in the form of supersonic Jets and rockets. The globalization and proximity brought out by the progress in Computer Sciences has necessitated developing social, democratic and humanistic values in the individuals for the welfare and their own adjustment in the society as well as for the broader cause of a globalize society.

TRAINING IN PROBLEM SOLVING ABILITY

Students of Computer Sciences get enough opportunities for being trained in the use of scientific method and scientific approach for the solution of the problems. In general, as we know, the following thought process is involved in solving a problem through scientific method in Computer Sciences.

1. The nature and purpose of the problem.
2. Analysis of the problem.
3. Testing the validity of the various solutions of the problem.
4. Testing of the results through applications in other problems of similar nature.

Acceptance of the result thus obtained as scientific principle and law for the future use. The above thought processes and methods are also needed in solving all

problems of daily life. First we try to understand the nature of the problem. Then after thinking of various solutions, we test them and finally accept one solution after examining its validity and correctness. As a result one gets desirable training in problem solving ability through the study of Computer Sciences.

The values and benefits so derived with the study of Computer Sciences in the school classes are thus very much responsible for making its study compulsory in all the ten years of schooling in our country. Utility of a subject in one's personal and social life is enough to justify its place in school curriculum. The study of Computer Sciences, whether we consider it from personal or social angle, has a quite unique importance in our own as well as in life of the nation. Therefore, the decision of providing a compulsory status for the study of sciences in corporation Computer Science is quite justified and valid.

VOCATIONAL VALUE

Study of Computer Science is helpful in opening vast vista of useful professions and vocations for being adopted in the future life. The most prestigious professions and vocations, related to engineering, computer technology, physical and chemical analysis of objects and environmental surrounding, physical and chemical inventions related to progress and development of the society etc. are keenly attached and dependent upon the study of Computer Sciences. The key of entering into a wide number of useful professions and vocations in our modern society (the list of which we have already given in chapter first of this text under the title "Professions in the area of Computer Sciences"). Very well

lies in the study of Computer Sciences. In this way, we can safely claim that study of Computer Sciences have a great “bread and butter value” enabling or preparing the would be citizens for earning their livelihood in their future life.

GLOBALIZATION AND COMPUTER SCIENCES

The term globalization derived from the words ‘globe’ and ‘signifies’ the removal of barriers of distances or of other nature for bringing people of the world together in terms of their relationships of events. For the understanding of its meaning let us think over on a few definitions given below :

1. The sociologist, Anthony Giddens, defines globalization as a decoupling of space and time, emphasizing that with instantaneous communications, knowledge and culture can be shared around the world simultaneously.
2. The Dutch academician Rudd Lubbers defines it as process in which geographical distances become of diminishing importance in the establishment and maintenance of cross economic, political and socio-cultural relations.
3. David Held and Anthony McGrew say that globalization can be conceived as a process (or set of processes) which embodies a transformation in the spatial organization of social relations and transactions, expressed in trans-continental or inter-regional flows and network of activity, interaction, and power.

These definitions point out the following things about the meaning and nature of the term globalization :

1. Globalization is the result of a rapid and efficient system of communication network. It aims towards the sharing of knowledge and culture among the people and communities of the world.
2. Globalization makes possible nearness between the individual and organization of the world community belonging to different regions, countries or continents in terms of their social relationships, interactions and sharing of comforts and power.

In the light of the above meaning and characteristics, the term globalization may be defined as the process of interlinking the individuals and organizations of the world community economically, politically and socio-culturally by barriers of the space and time.

ROLE OF COMPUTER SCIENCE IN THE PROCESS OF GLOBALIZATION

The global society or globalization is something, which is an accepted reality of the modern age. The development of science and technology especially related to the field of Computer Sciences is responsible to a great extent for accelerating the process of globalization. A few of such inventions and discoveries related to the subject of Computer Science may be named as below:

- Development of printing material and communication through written verbal material.
- Development of transport system with the discoveries and inventions related to travel on land surface, air, sea and space.
- Development of communication network with the inventions and discoveries like telephone, wireless,

telegraph, teleprompter, radio, television, video, camera, photography still and movie pictures, satellite communication, mobile telephony, tele-conferencing, video conferencing etc.

- Development of modern sophisticated communication system with the help of computer technology, internet access facilitating, e-mail and audio-visual transmission through web camera and web sites coupled with audio listening and speaking devices.
- The great invention and discoveries mentioned as above in the field of transport and communication services brought out by Computer Sciences have helped in the process of globalization in the ways narrated as below :
 - (i) Travelling from one corner of the globe to the other has become too convenient and less time consuming. It has provided needed mobility to the individuals or group as a whole for the interaction with the rest of the world.
 - (ii) It is quite easy and speedy as well as to transport commodities from one corner of the globe to the other. That is why, it is no wonder to get benefit of the productions of the commodities grown or manufactured in one part of the globe to the other parts in no time with a reduced cost and labour.
 - (iii) We can be in touch with the individuals, organizations and events belonging to the farthest corners of the globe with the help of the much development means and tools of communication. The distance

has to barriers at all for any type of communication— audio as well as visual.

The facilities available as above in terms of removing barriers of distances have invariably brought together the individuals, groups and nations of the world in the form of a single knitted unit truly in the spirit of “Basudheve Kutumbkam” i.e. “the whole world is my family”. As a result the people of the world have a close interaction, inter-dependence and inter-relationship with each other irrespective of their living and working distances. What happens with an individual, community, region or a country now affects in one way or the other the rest of the world. Poverty of a region is now the concern of the other events and political turmoil. As a result you may find that terrorism has become global, trafficking in drugs or human beings has become global, sexually transmitted and also other types of viral and infectious diseases have become global. On the other hand we now have a tremendous increase in the speed of globalization at the economic as well as socio-cultural exchange fronts. It has resulted in the process of denationalization of market, politics and legal systems for paving the way of the so-called global economy.

The emergence of the phenomenon of global trade and commerce has now provided enormous opportunities for an organization or company belonging to a particular part of the globe to establish itself in the foreign market. It is easy for her now to adopt first her product or services to the final users linguistic and cultural requirements and then take advantages of the internet revolution for establishing a virtual presence on the international market place with a multilingual corporate website or even resorting to an

electronic business, trade and commerce according to the mutual convenience of the business/trade partners. In this way globalization has resulted in intertwining the fates of the individuals, community and nations of the world and this dream has become true only on account of the services provided by the numerous inventions and discoveries made in the field of Computer Sciences.

COMPUTER SCIENCE RESEARCH INSTITUTIONS

Three categories of institutions conduct CS research in India: the major teaching and research institutes devoted to science and technology, government-sponsored laboratories and industry-sponsored laboratories.

- *The Six Major Research and Teaching Institutes Devoted to Science and Technology:* These six institutes are the IITs (of which there are five, with one more coming up in Assam) and IISc, located in Bangalore. These institutions form a select group in the minds of the government as well as the citizens.

The next tier of institutions is made up primarily of the Regional Engineering Colleges (RECs), with one located in each state. Also, there are several other universities where computer science research is being conducted quietly. (One such is the University of Hyderabad; researchers here are very active in (collaborative) AI research, keeping close contacts with overseas colleagues.) But a considerable gap does exist between the six top tier institutions and the next because of high teaching load imposed on the faculty, students being, on average, of a lower quality, and finally

poorer infrastructure, namely library and computing facilities.

- *Government-sponsored Institutions:* These institutions are funded by different government ministries and departments.

TIFR and the Institute for Mathematical Sciences (MatScience) perform research which is predominantly of a theoretical nature. These are funded by DAE. Defence-related work takes place in a number of labs around the country, many located in Bangalore and Hyderabad, both in Southern India. A good example is CAIR which can be described as a “think-tank” serving the AI and robotics needs of Indian Ministry of Defence. It is a component of the Defence Research and Development Organization (DRDO).

The Ministry of Planning funds ISI, with its primary location in Calcutta. (It is worth noting that the first indigenous digital computer—fabricated using discrete transistor units—was commissioned by ISI in 1966 in collaboration with Jadavpur University.) NCST also carries out research in several areas of computer science besides having education and training among its functions. NCST is a successor of the erstwhile National Centre for Software Development and Computing Techniques (NCSDCT) which was a component of TIFR.

The Indian Space Research Organization (ISRO), is also involved in computer science work, but most of its work is of an applied nature, in the context of satellites and launch vehicles. ISRO has been building satellites for remote sensing as well as for communication. Its most recent success involved the launch of the Polar Synchronous Launch Vehicle

capable of launching 1000-Kg class satellites into sun-synchronous orbits. National Aerospace Laboratories (NAL), Bhabha Atomic Research Centre (BARC), and Centre for the Development of Advanced Computation (CDAC) have had the development of parallel processing platforms for solving computational science problems as the main focus of their computer science research.

- *Industrial Labs:* Tata Research Development and Design Centre (TRDDC), supported by the Tata group of companies, and the SPIC science foundation (SSF), sponsored by SPIC, a petrochemical corporation, are good examples. While the latter is primarily involved in theoretical computer science research, TRDDC is geared up to “result-oriented research” to meet the needs of Tata Consultancy Services (TCS) and its clients, and more generally, the Tata group of companies. The uniqueness of TRDDC comes from its self-supporting R&D effort. Even though most of the projects are done for TCS, TRDDC also has funds from DST, MoD, and other government organizations.

In addition, there are several labs that are sponsored by many multinationals in India, such as Texas Instruments (which led the way, as early as in 1987), Motorola, and Oracle. But, at this point, these labs are mostly involved in developmental activities defined and subcontracted by their parent organizations.

STUDENTS AND FACULTY AT THE EDUCATIONAL INSTITUTIONS

US and Europe have a very high regard for graduate students trained at the IITs and in IISc. The reasons are

easy to find. Admission to the Bachelors' programme (called B.Tech) at the IITs is through a fiercely competitive entrance examination called the Joint Entrance Examination (JEE). It is written by over 100,000 students every year, with less than 1500 selected—based purely on their ranking in the JEE.

The curriculum at the IITs is on par with top institutions in developed countries. A laudable feature of the B.Tech programme is that during the final year, students are required to do a project, which in many instances are quite ambitious with students going on to publish their work in conferences and journals.

All the academic research institutions have a Masters' programme (called M.Tech), students selected once again after a competitive examination called GATE, perhaps not quite as competitive as JEE. M.Tech's come (more often than not) from a non-IIT background and their preparedness in computer science is also considered to be lacking, but they make up for it through their hard work. While some of the IITs require computer science background (for example, IIT Bombay does not admit students who are not among the top 7 per cent in computer science GATE), others do not.

M.Tech's also do a project, but since the programme itself is only three semesters long and the project is done during the third, there is less time to execute a project of an involved nature. Also, if a student arrives with no computer science background, it does not allow him or her to do a substantial project.

Some institutions have a programme called Master of Science, which is purely research project oriented. The

number in this category is quite small and not all academic institutions have an M.S. programme.

As part of the Quality Improvement Programme (QIP), an effort to modernize the faculty at teaching institutions, faculty from smaller colleges and universities are encouraged to apply to do their Ph.Ds at IITs and other institutions. For instance, four of the computer science faculty at the Regional Engineering College (REC) in Trichy (out of just about a dozen) are currently doing their Ph.D.'s under this programme. Since most of the faculty at smaller educational institutions offering programmes in CS do not have Ph.Ds, the QIP is a worthy enterprise.

All graduate students are supported by the Ministry of Human Resource Development (MHRD) (even though there is a move, especially for Ph.D. students to be paid from research grants). If one enters the Ph.D. programme after a Bachelors' degree, assistantship is guaranteed for five years; for those joining after a Masters degree, four and a half years. Because of this and the fact that all faculty salaries are on a 12-month basis, there is very little compulsion for an academic researcher to seek external funding for his or her research, except perhaps to support travel and for books and equipment.

IISc has the largest contingent of Ph.D. students, having produced 20 Ph.Ds during 94-95 with over 40 in the pipeline. Each IIT has between 10 and 15 students pursuing their Ph.Ds., a number smaller on average compared to the recent past. I often heard the remark that the number of students interested in a Ph.D degree and of Ph.D. caliber has been decreasing. This does not augur well for the future of Indian computer science. It is hence important to understand the

reasons for this trend. We will discuss some of these here. Faculty at the top educational institutions belong to three ranks, Assistant Professor, Associate Professor and Professor. The nomenclature varies in other institutions. While fresh Ph.Ds have joined as Assistant Professors in some of the institutions, most places insist on a few years' experience. Even though Assistant Professorship is tenured, the requirements to move into a higher rank are not very transparent. In fact, IITs don't appear to have the notion of promotions. If one desires a higher rank, he or she must respond to an advertisement for that rank. The lack of transparency as well as norms for such promotions causes, as a junior faculty put it, "the dilemma of whether to focus on 'local developmental work' or 'internationally publishable work'."

With the "graying" of the research-oriented academic institutions, there is a need to infuse fresh blood. But with very few retirements in the offing, the number of openings in established institutions is small. Even if one finds a position, a fresh Assistant Professor should be content with a take-home-pay that could be substantially less than the earnings of a fresh Master's student entering the private sector.

A large proportion of the well-trained IIT and IISc students go abroad, mainly to the US, for their further studies. So India does not reap the benefits of having trained them. Among the IIT students who stay in India after their B.Tech's, a sizable number opt for a career in management, and hence join the MBA programme in one of the four prestigious Indian Institutes of Management (IIM). As a result, the number of

IIT students who pursue a career in science or engineering within India is dishearteningly small.

Whereas in the pure sciences, the percentage of women in India (both among students and faculty) is much higher than in the US, in engineering disciplines, the situation is the opposite. Being associated with engineering, the computer science research roster in India has a very small number of women. To its credit, IIT Madras has two women faculty in the computer science department, with one serving as its current head. IISc and IIT Bombay have one each in their faculty.

Inbreeding, often associated with stagnation, is a common phenomenon at the top educational institutions, with the ratio—of those with a Ph.D. from the same institution—sometimes approaching one-half. Even though an institution tries not to hire one of its own immediately after his or her Ph.D., most have no hesitation if the person has spent a year or two elsewhere before returning to his or her Ph.D. institution.

Among the educational institutions, according to many computer science researchers around the country, IISc has a greater appreciation for research in pure sciences and this is true for computer science as well. Many of the IIT faculty I talked to feel that training and educating students is their foremost task, in contrast, for example, with IISc, ISI, and TIFR, where research is. There is also a general feeling that the IITs can do a better job if they set their performance thresholds higher. The directors of the IITs have a bigger say in setting the directions of their respective institutions than the heads of the above mentioned research-oriented

institutions. Hence, any changes or redirections, giving more emphasis and recognition for research and thereby tapping the available research potential, awaits their directives and blessings.

COMPUTER SCIENCES ON MODERN COMMUNITIES

Development and progress in the field of Computer Sciences has influenced the life and livings of the modern Communities and Society in so many ways like below:

CONSTRUCTION OF BUILDINGS AND RESIDENTIAL COLONIES

On Account of the availability of a variety of load bearable and non-load bearable material through the development in Computer Sciences, the modern communities look modern in terms of the construction of their buildings in the form of business centers, offices and residential colonies.

Transportation and Communication Systems

Development in Computer Sciences has been able to provide the latest available transportation and communication systems to the modern communities. The distances are no more a barrier for the people living at the farthest distances of the globe.

Modernization of the Systems

The development in the Computer Sciences have modernized the sources of availability of food stuff in the shape of modernization of the method of farming, poultry farming, cattle beading, fisheries, bee keeping etc. It has resulted in multiplying the production of the food stuff as well as reducing the complexities or manual labour. The food

stuff cannot be better preserved through the modern techniques available as a result of invention and discoveries in Computer Sciences.

Development in Computer Sciences is helping the modern communities to take care of their water resources. It has provided artificial irrigation means as well as availability of drinking water with the construction of big water reservoirs, dams and sophisticated distribution system. It has provided big plants and simple household gadgets for the availability of pure drinking water to the modern communities. It has also provided means to have artificial rains and cultivation of water for providing additional sources of water to the modern communities.

Modern means for the entertainment and leisure time hobbies

Development in Computer Science has provided modern and methods for the entertainment and uplifting of leisure to the modern communities. Radio, Television Video, Films and Computer services have taken a total command of providing entertainment and leisure time hobbies to the modern communities.

Health care and treatment of diseases

Development in Computer Sciences have helped much in taking care of the health including treatment of illness and diseases of the members of the modern communities. It has provided better knowledge and information for the prevention and care of the diseases as well maintenance of good physical and mental health through its wider network of information technology. The dreaded diseases are now no more so dreaded

as happened to be in the past. With the vast discoveries and invention in the field of health and medical sciences as well as tremendous progress in chemical sciences. Modern communities can avail the latest treatment of the diseases and look after of their health.

Development of inter-relationship and dependence

Development of Computer Sciences are responsible for making the modern communities too much inter-related and inter-dependent. It has given birth to the phenomenon of globalization in every aspect—physical, mental, emotional social, cultural and ethical of the behaviour of teach and every person belonging to the modern communities of this globe.

What happens to a community of the world, residing at any farthest corner of the world equally affects the working and behaviour of the other segment of the world communities. We will discuss in detail such globalization impact of Computer Sciences soon in this very chapter.

With all what has been said above, we should not conclude the developments in Computer Sciences are always bound to cast positive and desirable impact well being and progress of the modern communities. If handled improperly and utilized destructively these can yield bitter and horrifying results. Such negative impact of the development Computer Sciences on modern communities may be summarized as below:

- Too much urbanization of the communities.
- Causing heavy pollution of every sort like air

Computer Memory System

pollution, water pollution, noise pollution, cultural pollution etc.

- Inequitable distribution of wealth and other material comforts in the population.
- Abolishment of the concept and existence of the harming of health and welfare of the people.
- Development of the weapons of destruction and their unmindful application.
- Side effects of the fertilizers, chemicals, pesticides insecticides used in growing foodstuff and killing harmful bio-stuff.
- Neglect of moral values and social responsibilities at the cost of material development and individualism.

4

Computer Hard Disk Drive

A hard disk drive (HDD) is a non-volatile, random access device for digital data. It features rotating rigid platters on a motor-driven spindle within a protective enclosure. Data is magnetically read from and written to the platter by read/write heads that float on a film of air above the platters.

Introduced by IBM in 1956, hard disk drives have fallen in cost and physical size over the years while dramatically increasing in capacity. Hard disk drives have been the dominant device for secondary storage of data in general purpose computers since the early 1960s. They have maintained this position because advances in their areal recording density have kept pace with the requirements for secondary storage. Today's HDDs operate on high-speed serial interfaces; i.e., serial ATA (SATA) or serial attached SCSI (SAS).

HISTORY OF HARD DISK DRIVE

Hard disk drives were introduced in 1956 as data storage for an IBM accounting computer and were developed for use with general purpose mainframe and mini computers.

Driven by areal density doubling every two to four years since their invention, HDDs have changed in many ways, a few highlights include:

- Capacity per HDD increasing from 3.75 megabytes to greater than 1 terabyte, a greater than 270 thousand to 1 improvement.
- Size of HDD decreasing from 87.9 cubic feet (a double wide refrigerator) to 0.002 cubic feet (2½-inch form factor, a pack of cards), a greater than 44 thousand to 1 improvement.
- Price decreasing from about \$15,000 per megabyte to less than \$0.0001 per megabyte (\$100/1 terabyte), a greater than 150 million to 1 improvement.
- Average access time decreasing from greater than 0.1 second to a few thousandths of a second, a greater than 40 to 1 improvement.
- Market application expanding from general purpose computers to most computing applications including consumer applications.

PERFORMANCE CHARACTERISTICS

Access Time

The factors that limit the time to access the data on a hard disk drive (Access time) are mostly related to the mechanical nature of the rotating disks and moving heads.

Seek time is a measure of how long it takes the head assembly to travel to the track of the disk that contains data. Latency is rotational delay incurred because the desired disk sector may not be directly under the head when data transfer is requested. These two delays are on the order of milliseconds each. The bit rate or data transfer rate once the head is in the right position creates delay which is a function of the number of blocks transferred; typically relatively small, but can be quite long with the transfer of large contiguous files. Delay may also occur if the drive disks are stopped to save energy, see Power management.

An HDD's Average Access Time is its average Seek time which technically is the time to do all possible seeks divided by the number of all possible seeks, but in practice is determined by statistical methods or simply approximated as the time of a seek over one-third of the number of tracks

Defragmentation is a procedure used to minimize delay in retrieving data by moving related items to physically proximate areas on the disk. Some computer operating systems perform defragmentation automatically. Although automatic defragmentation is intended to reduce access delays, the procedure can slow response when performed while the computer is in use.

Access time can be improved by increasing rotational speed, thus reducing latency and/or by decreasing seek time. Increasing areal density increases throughput by increasing data rate and by increasing the amount of data under a set of heads, thereby potentially reducing seek activity for a given amount of data. Based on historic trends,

analysts predict a future growth in HDD areal density (and therefore capacity) of about 40% per year. Access times have not kept up with throughput increases, which themselves have not kept up with growth in storage capacity.

Seek Time

Average Seek time ranges from 3 ms for high-end server drives, to 15 ms for mobile drives, with the most common mobile drives at about 12 ms and the most common desktop type typically being around 9 ms. The first HDD had an average seek time of about 600 ms and by the middle 1970s HDDs were available with seek times of about 25 ms. Some early PC drives used a stepper motor to move the heads, and as a result had seek times as slow as 80–120 ms, but this was quickly improved by voice coil type actuation in the 1980s, reducing seek times to around 20 ms. Seek time has continued to improve slowly over time.

Latency

Latency is the delay for the rotation of the disk to bring the required disk sector under the read-write mechanism. It depends on rotational speed of a disk, measured in revolutions per minute (RPM). Average rotational delay is shown in the table below, based on the empirical relation that the average latency in milliseconds for such a drive is one-half the rotational period:

<i>Spindle [rpm]</i>	<i>Average latency [ms]</i>
4200	7.14
5400	5.56
7200	4.17
10000	3
15000	2

Data Transfer Rate

As of 2010, a typical 7200 rpm desktop hard drive has a sustained “disk-to-buffer” data transfer rate up to 1030 Mbits/sec. This rate depends on the track location, so it will be higher for data on the outer tracks (where there are more data sectors) and lower toward the inner tracks (where there are fewer data sectors); and is generally somewhat higher for 10,000 rpm drives. A current widely used standard for the “buffer-to-computer” interface is 3.0 Gbit/s SATA, which can send about 300 megabyte/s from the buffer to the computer, and thus is still comfortably ahead of today’s disk-to-buffer transfer rates. Data transfer rate (read/write) can be measured by writing a large file to disk using special file generator tools, then reading back the file. Transfer rate can be influenced by file system fragmentation and the layout of the files.

HDD data transfer rate depends upon the rotational speed of the platters and the data recording density. Because heat and vibration limit rotational speed, advancing density becomes the main method to improve sequential transfer rates. Areal density advances by increasing both the number of tracks across the disk and the number of sectors per track, the later will increase the data transfer rate (for a given RPM). Since data transfer rate performance only tracks one of the two components of areal density, its performance improves at lower rate,

Power Consumption

Power consumption has become increasingly important, not only in mobile devices such as laptops but also in server

and desktop markets. Increasing data center machine density has led to problems delivering sufficient power to devices (especially for spin up), and getting rid of the waste heat subsequently produced, as well as environmental and electrical cost concerns. Heat dissipation directly tied to power consumption, and as drive age, disk failure rates increase at higher drive temperatures. Similar issues exist for large companies with thousands of desktop PCs. Smaller form factor drives often use less power than larger drives. One interesting development in this area is actively controlling the seek speed so that the head arrives at its destination only just in time to read the sector, rather than arriving as quickly as possible and then having to wait for the sector to come around (i.e. the rotational latency). Many of the hard drive companies are now producing Green Drives that require much less power and cooling. Many of these Green Drives spin slower (<5,400 rpm compared to 7,200, 10,000 or 15,000 rpm) and also generate less waste heat. Power consumption can also be reduced by parking the drive heads when the disk is not in use reducing friction, adjusting spin speeds according to transfer rates, and disabling internal components when not in use.

Also in systems where there might be multiple hard disk drives, there are various ways of controlling when the hard drives spin up since the highest current is drawn at that time.

- On SCSI hard disk drives, the SCSI controller can directly control spin up and spin down of the drives.
- On Parallel ATA (aka PATA) and Serial ATA (SATA)

hard disk drives, some support power-up in standby or PUIS. The hard disk drive will not spin up until the controller or system BIOS issues a specific command to do so. This limits the power draw or consumption upon power on.

- Some SATA II hard disk drives support staggered spin-up, allowing the computer to spin up the drives in sequence to reduce load on the power supply when booting.

Power Management

Most hard disk drives today support some form of power management which uses a number of specific power modes that save energy by reducing performance. When implemented an HDD will change between a full power mode to one or more power saving modes as a function of drive usage. Recovery from the deepest mode, typically called Sleep, may take as long as several seconds.

Audible Noise

Measured in dBA, audible noise is significant for certain applications, such as DVRs, digital audio recording and quiet computers. Low noise disks typically use fluid bearings, slower rotational speeds (usually 5,400 rpm) and reduce the seek speed under load (AAM) to reduce audible clicks and crunching sounds. Drives in smaller form factors (e.g. 2.5 inch) are often quieter than larger drives.

Shock Resistance

Shock resistance is especially important for mobile devices. Some laptops now include active hard drive

protection that parks the disk heads if the machine is dropped, hopefully before impact, to offer the greatest possible chance of survival in such an event. Maximum shock tolerance to date is 350 g for operating and 1000 g for non-operating.

ACCESS AND INTERFACES

Hard disk drives are accessed over one of a number of bus types, including parallel ATA (P-ATA, also called IDE or EIDE), Serial ATA (SATA), SCSI, Serial Attached SCSI (SAS), and Fibre Channel. Bridge circuitry is sometimes used to connect hard disk drives to buses that they cannot communicate with natively, such as IEEE 1394, USB and SCSI.

For the ST-506 interface, the data encoding scheme as written to the disk surface was also important. The first ST-506 disks used Modified Frequency Modulation (MFM) encoding, and transferred data at a rate of 5 megabits per second. Later controllers using 2,7 RLL (or just "RLL") encoding caused 50% more data to appear under the heads compared to one rotation of an MFM drive, increasing data storage and data transfer rate by 50%, to 7.5 megabits per second.

Many ST-506 interface disk drives were only specified by the manufacturer to run at the 1/3 lower MFM data transfer rate compared to RLL, while other drive models (usually more expensive versions of the same drive) were specified to run at the higher RLL data transfer rate. In some cases, a drive had sufficient margin to allow the MFM

specified model to run at the denser/faster RLL data transfer rate (not recommended nor guaranteed by manufacturers). Also, any RLL-certified drive could run on any MFM controller, but with 1/3 less data capacity and as much as 1/3 less data transfer rate compared to its RLL specifications.

Enhanced Small Disk Interface (ESDI) also supported multiple data rates (ESDI disks always used 2,7 RLL, but at 10, 15 or 20 megabits per second), but this was usually negotiated automatically by the disk drive and controller; most of the time, however, 15 or 20 megabit ESDI disk drives were not downward compatible (i.e. a 15 or 20 megabit disk drive would not run on a 10 megabit controller). ESDI disk drives typically also had jumpers to set the number of sectors per track and (in some cases) sector size.

Modern hard drives present a consistent interface to the rest of the computer, no matter what data encoding scheme is used internally. Typically a DSP in the electronics inside the hard drive takes the raw analog voltages from the read head and uses PRML and Reed-Solomon error correction to decode the sector boundaries and sector data, then sends that data out the standard interface. That DSP also watches the error rate detected by error detection and correction, and performs bad sector remapping, data collection for Self-Monitoring, Analysis, and Reporting Technology, and other internal tasks.

SCSI originally had just one signaling frequency of 5 MHz for a maximum data rate of 5 megabytes/second over 8 parallel conductors, but later this was increased dramatically. The SCSI bus speed had no bearing on the disk's internal

speed because of buffering between the SCSI bus and the disk drive's internal data bus; however, many early disk drives had very small buffers, and thus had to be reformatted to a different interleave (just like ST-506 disks) when used on slow computers, such as early Commodore Amiga, IBM PC compatibles and Apple Macintoshes.

ATA disks have typically had no problems with interleave or data rate, due to their controller design, but many early models were incompatible with each other and could not run with two devices on the same physical cable in a master/slave setup. This was mostly remedied by the mid-1990s, when ATA's specification was standardized and the details began to be cleaned up, but still causes problems occasionally (especially with CD-ROM and DVD-ROM disks, and when mixing Ultra DMA and non-UDMA devices).

Serial ATA does away with master/slave setups entirely, placing each disk on its own channel (with its own set of I/O ports) instead. FireWire/IEEE 1394 and USB(1.0/2.0) HDDs are external units containing generally ATA or SCSI disks with ports on the back allowing very simple and effective expansion and mobility.

Most FireWire/IEEE 1394 models are able to daisy-chain in order to continue adding peripherals without requiring additional ports on the computer itself. USB however, is a point to point network and does not allow for daisy-chaining. USB hubs are used to increase the number of available ports and are used for devices that do not require charging since the current supplied by hubs is typically lower than what's available from the built-in USB ports.

Disk Interface Families used in Personal Computers

Notable families of disk interfaces include:

- Historical bit serial interfaces connect a hard disk drive (HDD) to a hard disk controller (HDC) with two cables, one for control and one for data. (Each drive also has an additional cable for power, usually connecting it directly to the power supply unit). The HDC provided significant functions such as serial/parallel conversion, data separation, and track formatting, and required matching to the drive (after formatting) in order to assure reliability. Each control cable could serve two or more drives, while a dedicated (and smaller) data cable served each drive.
 - o ST506 used MFM (Modified Frequency Modulation) for the data encoding method.
 - o ST412 was available in either MFM or RLL (Run Length Limited) encoding variants.
 - o Enhanced Small Disk Interface (ESDI) was an industry standard interface similar to ST412 supporting higher data rates between the processor and the disk drive.
- Modern bit serial interfaces connect a hard disk drive to a host bus interface adapter (today typically integrated into the “south bridge”) with one data/control cable. (As for historical *bit serial interfaces* above, each drive also has an additional power cable, usually direct to the power supply unit.)
 - o Fibre Channel (FC), is a successor to parallel SCSI interface on enterprise market. It is a serial protocol.

In disk drives usually the Fibre Channel Arbitrated Loop (FC-AL) connection topology is used. FC has much broader usage than mere disk interfaces, and it is the cornerstone of storage area networks (SANs). Recently other protocols for this field, like iSCSI and ATA over Ethernet have been developed as well. Confusingly, drives usually use *copper* twisted-pair cables for Fibre Channel, not fibre optics. The latter are traditionally reserved for larger devices, such as servers or disk array controllers.

- o Serial ATA (SATA). The SATA data cable has one data pair for differential transmission of data to the device, and one pair for differential receiving from the device, just like EIA-422. That requires that data be transmitted serially. Similar differential signaling system is used in RS485, LocalTalk, USB, Firewire, and differential SCSI.
- o Serial Attached SCSI (SAS). The SAS is a new generation serial communication protocol for devices designed to allow for much higher speed data transfers and is compatible with SATA. SAS uses a mechanically identical data and power connector to standard 3.5-inch SATA1/SATA2 HDDs, and many server-oriented SAS RAID controllers are also capable of addressing SATA hard drives. SAS uses serial communication instead of the parallel method found in traditional SCSI devices but still uses SCSI commands.
- Word serial interfaces connect a hard disk drive to a host bus adapter (today typically integrated into

the “south bridge”) with one cable for combined data/control. (As for all *bit serial interfaces* above, each drive also has an additional power cable, usually direct to the power supply unit.) The earliest versions of these interfaces typically had a 8 bit parallel data transfer to/from the drive, but 16-bit versions became much more common, and there are 32 bit versions. Modern variants have serial data transfer. The word nature of data transfer makes the design of a host bus adapter significantly simpler than that of the precursor HDD controller.

- o Integrated Drive Electronics (IDE), later renamed to ATA, with the alias P-ATA (“parallel ATA”) retroactively added upon introduction of the new variant Serial ATA. The original name reflected the innovative integration of HDD controller with HDD itself, which was not found in earlier disks. Moving the HDD controller from the interface card to the disk drive helped to standardize interfaces, and to reduce the cost and complexity. The 40-pin IDE/ATA connection transfers 16 bits of data at a time on the data cable. The data cable was originally 40-conductor, but later higher speed requirements for data transfer to and from the hard drive led to an “ultra DMA” mode, known as UDMA. Progressively swifter versions of this standard ultimately added the requirement for a 80-conductor variant of the same cable, where half of the conductors provides grounding necessary for enhanced high-speed signal quality by reducing cross talk. The interface

for 80-conductor only has 39 pins, the missing pin acting as a key to prevent incorrect insertion of the connector to an incompatible socket, a common cause of disk and controller damage.

- o EIDE was an unofficial update (by Western Digital) to the original IDE standard, with the key improvement being the use of direct memory access (DMA) to transfer data between the disk and the computer without the involvement of the CPU, an improvement later adopted by the official ATA standards. By directly transferring data between memory and disk, DMA eliminates the need for the CPU to copy byte per byte, therefore allowing it to process other tasks while the data transfer occurs.
- o Small Computer System Interface (SCSI), originally named SASI for Shugart Associates System Interface, was an early competitor of ESDI. SCSI disks were standard on servers, workstations, Commodore Amiga, and Apple Macintosh computers through the mid-1990s, by which time most models had been transitioned to IDE (and later, SATA) family disks. Only in 2005 did the capacity of SCSI disks fall behind IDE disk technology, though the highest-performance disks are still available in SCSI and Fibre Channel only. The range limitations of the data cable allows for external SCSI devices. Originally SCSI data cables used single ended (common mode) data transmission, but server class SCSI could use differential transmission, either low voltage

differential (LVD) or high voltage differential (HVD). (“Low” and “High” voltages for differential SCSI are relative to SCSI standards and do not meet the meaning of low voltage and high voltage as used in general electrical engineering contexts, as apply e.g. to statutory electrical codes; both LVD and HVD use low voltage signals (3.3 V and 5 V respectively) in general terminology.)

Integrity

Due to the extremely close spacing between the heads and the disk surface, hard disk drives are vulnerable to being damaged by a head crash—a failure of the disk in which the head scrapes across the platter surface, often grinding away the thin magnetic film and causing data loss. Head crashes can be caused by electronic failure, a sudden power failure, physical shock, contamination of the drive’s internal enclosure, wear and tear, corrosion, or poorly manufactured platters and heads.

The HDD’s spindle system relies on air pressure inside the disk enclosure to support the heads at their proper *flying height* while the disk rotates. Hard disk drives require a certain range of air pressures in order to operate properly. The connection to the external environment and pressure occurs through a small hole in the enclosure (about 0.5 mm in breadth), usually with a filter on the inside (the *breather filter*). If the air pressure is too low, then there is not enough lift for the flying head, so the head gets too close to the disk, and there is a risk of head crashes and data loss. Specially manufactured sealed and pressurized disks are needed for

reliable high-altitude operation, above about 3,000 m (10,000 feet). Modern disks include temperature sensors and adjust their operation to the operating environment. Breather holes can be seen on all disk drives—they usually have a sticker next to them, warning the user not to cover the holes. The air inside the operating drive is constantly moving too, being swept in motion by friction with the spinning platters. This air passes through an internal recirculation (or “recirc”) filter to remove any leftover contaminants from manufacture, any particles or chemicals that may have somehow entered the enclosure, and any particles or outgassing generated internally in normal operation. Very high humidity for extended periods can corrode the heads and platters.

For giant magnetoresistive (GMR) heads in particular, a minor head crash from contamination (that does not remove the magnetic surface of the disk) still results in the head temporarily overheating, due to friction with the disk surface, and can render the data unreadable for a short period until the head temperature stabilizes (so called “thermal asperity”, a problem which can partially be dealt with by proper electronic filtering of the read signal).

Actuation of Moving Arm

The hard drive’s electronics control the movement of the actuator and the rotation of the disk, and perform reads and writes on demand from the disk controller. Feedback of the drive electronics is accomplished by means of special segments of the disk dedicated to servo feedback. These are either complete concentric circles (in the case of dedicated

servo technology), or segments interspersed with real data (in the case of embedded servo technology). The servo feedback optimizes the signal to noise ratio of the GMR sensors by adjusting the voice-coil of the actuated arm. The spinning of the disk also uses a servo motor. Modern disk firmware is capable of scheduling reads and writes efficiently on the platter surfaces and remapping sectors of the media which have failed.

Landing Zones and Load/Unload Technology

Modern HDDs prevent power interruptions or other malfunctions from landing its heads in the data zone by parking the heads either in a landing zone or by unloading (i.e., load/unload) the heads. Some early PC HDDs did not park the heads automatically and they would land on data. In some other early units the user manually parked the heads by running a programme to park the HDD's heads.

A landing zone is an area of the platter usually near its inner diameter (ID), where no data are stored. This area is called the Contact Start/Stop (CSS) zone. Disks are designed such that either a spring or, more recently, rotational inertia in the platters is used to park the heads in the case of unexpected power loss. In this case, the spindle motor temporarily acts as a generator, providing power to the actuator.

Spring tension from the head mounting constantly pushes the heads towards the platter. While the disk is spinning, the heads are supported by an air bearing and experience no physical contact or wear. In CSS drives the sliders carrying the head sensors (often also just called *heads*) are

designed to survive a number of landings and takeoffs from the media surface, though wear and tear on these microscopic components eventually takes its toll. Most manufacturers design the sliders to survive 50,000 contact cycles before the chance of damage on startup rises above 50%. However, the decay rate is not linear: when a disk is younger and has had fewer start-stop cycles, it has a better chance of surviving the next startup than an older, higher-mileage disk (as the head literally drags along the disk's surface until the air bearing is established). For example, the Seagate Barracuda 7200.10 series of desktop hard disks are rated to 50,000 start-stop cycles, in other words no failures attributed to the head-platter interface were seen before at least 50,000 start-stop cycles during testing.

Around 1995 IBM pioneered a technology where a landing zone on the disk is made by a precision laser process (*Laser Zone Texture* = LZT) producing an array of smooth nanometer-scale "bumps" in a landing zone, thus vastly improving stiction and wear performance. This technology is still largely in use today (2008), predominantly in desktop and enterprise (3.5 inch) drives. In general, CSS technology can be prone to increased stiction (the tendency for the heads to stick to the platter surface), e.g. as a consequence of increased humidity. Excessive stiction can cause physical damage to the platter and slider or spindle motor.

Load/Unload technology relies on the heads being lifted off the platters into a safe location, thus eliminating the risks of wear and stiction altogether. The first HDD RAMAC and most early disk drives used complex mechanisms to

load and unload the heads. Modern HDDs use ramp loading, first introduced by Memorex in 1967, to load/unload onto plastic “ramps” near the outer disk edge.

All HDDs today still use one of these two technologies listed above. Each has a list of advantages and drawbacks in terms of loss of storage area on the disk, relative difficulty of mechanical tolerance control, non-operating shock robustness, cost of implementation, etc.

Addressing shock robustness, IBM also created a technology for their ThinkPad line of laptop computers called the Active Protection System. When a sudden, sharp movement is detected by the built-in accelerometer in the Thinkpad, internal hard disk heads automatically unload themselves to reduce the risk of any potential data loss or scratch defects. Apple later also utilized this technology in their PowerBook, iBook, MacBook Pro, and MacBook line, known as the Sudden Motion Sensor. Sony, HP with their HP 3D DriveGuard and Toshiba have released similar technology in their notebook computers.

This accelerometer-based shock sensor has also been used for building cheap earthquake sensor networks.

TECHNOLOGY: MAGNETIC RECORDING

HDDs record data by magnetizing ferromagnetic material directionally. Sequential changes in the direction of magnetization represent patterns of binary data bits. The data are read from the disk by detecting the transitions in magnetization and decoding the originally written data. Different encoding schemes, such as Modified Frequency

Modulation, group code recording, run-length limited encoding, and others are used.

A typical HDD design consists of a spindle that holds flat circular disks called platters, onto which the data are recorded. The platters are made from a non-magnetic material, usually aluminum alloy or glass, and are coated with a shallow layer of magnetic material typically 10–20 nm in depth, with an outer layer of carbon for protection. For reference, standard copy paper is 0.07–0.18 millimetre (70,000–180,000 nm).

The platters are spun at speeds varying from 3,000 RPM in energy-efficient portable devices, to 15,000 RPM for high performance servers. Information is written to, and read from a platter as it rotates past devices called read-and-write heads that operate very close (tens of nanometers in new drives) over the magnetic surface. The read-and-write head is used to detect and modify the magnetization of the material immediately under it. In modern drives there is one head for each magnetic platter surface on the spindle, mounted on a common arm. An actuator arm (or access arm) moves the heads on an arc (roughly radially) across the platters as they spin, allowing each head to access almost the entire surface of the platter as it spins. The arm is moved using a voice coil actuator or in some older designs a stepper motor.

The magnetic surface of each platter is conceptually divided into many small sub-micrometer-sized magnetic regions referred to as magnetic domains. In older disk designs the regions were oriented horizontally and parallel

to the disk surface, but beginning about 2005, the orientation was changed to perpendicular to allow for closer magnetic domain spacing. Due to the polycrystalline nature of the magnetic material each of these magnetic regions is composed of a few hundred magnetic grains. Magnetic grains are typically 10 nm in size and each form a single magnetic domain. Each magnetic region in total forms a magnetic dipole which generates a magnetic field.

For reliable storage of data, the recording material needs to resist self-demagnetization, which occurs when the magnetic domains repel each other. Magnetic domains written too densely together to a weakly magnetizable material will degrade over time due to physical rotation of one or more domains to cancel out these forces. The domains rotate sideways to a halfway position that weakens the readability of the domain and relieves the magnetic stresses. Older hard disks used iron(III) oxide as the magnetic material, but current disks use a cobalt-based alloy.

A write head magnetizes a region by generating a strong local magnetic field. Early HDDs used an electromagnet both to magnetize the region and to then read its magnetic field by using electromagnetic induction. Later versions of inductive heads included metal in Gap (MIG) heads and thin film heads. As data density increased, read heads using magnetoresistance (MR) came into use; the electrical resistance of the head changed according to the strength of the magnetism from the platter. Later development made use of spintronics; in these heads, the magnetoresistive effect was much greater than in earlier types, and was dubbed “giant” magnetoresistance (GMR). In today’s heads,

the read and write elements are separate, but in close proximity, on the head portion of an actuator arm. The read element is typically magneto-resistive while the write element is typically thin-film inductive.

The heads are kept from contacting the platter surface by the air that is extremely close to the platter; that air moves at or near the platter speed. The record and playback head are mounted on a block called a slider, and the surface next to the platter is shaped to keep it just barely out of contact. This forms a type of air bearing.

In modern drives, the small size of the magnetic regions creates the danger that their magnetic state might be lost because of thermal effects. To counter this, the platters are coated with two parallel magnetic layers, separated by a 3-atom layer of the non-magnetic element ruthenium, and the two layers are magnetized in opposite orientation, thus reinforcing each other. Another technology used to overcome thermal effects to allow greater recording densities is perpendicular recording, first shipped in 2005, and as of 2007 the technology was used in many HDDs.

Components

A typical hard disk drive has two electric motors; a disk motor to spin the disks and an actuator (motor) to position the read/write head assembly across the spinning disks. The disk motor has an external rotor attached to the disks; the stator windings are fixed in place. Opposite the actuator at the end of the head support arm is the read-write head (near center in photo); thin printed-circuit cables connect the read-write heads to amplifier electronics mounted at the

pivot of the actuator. A flexible, somewhat U-shaped, ribbon cable, seen edge-on below and to the left of the actuator arm continues the connection to the controller board on the opposite side.

The head support arm is very light, but also stiff; in modern drives, acceleration at the head reaches 550 Gs. The silver-colored structure at the upper left of the first image is the top plate of the actuator, a permanent-magnet and moving coil motor that swings the heads to the desired position (it is shown removed in the second image). The plate supports a squat neodymium-iron-boron (NIB) high-flux magnet. Beneath this plate is the moving coil, often referred to as the *voice coil* by analogy to the coil in loudspeakers, which is attached to the actuator hub, and beneath that is a second NIB magnet, mounted on the bottom plate of the motor (some drives only have one magnet).

The voice coil itself is shaped rather like an arrowhead, and made of doubly coated copper magnet wire. The inner layer is insulation, and the outer is thermoplastic, which bonds the coil together after it is wound on a form, making it self-supporting. The portions of the coil along the two sides of the arrowhead (which point to the actuator bearing center) interact with the magnetic field, developing a tangential force that rotates the actuator. Current flowing radially outward along one side of the arrowhead and radially inward on the other produces the tangential force. If the magnetic field were uniform, each side would generate opposing forces that would cancel each other out. Therefore the surface of the magnet is half N pole, half S pole, with the radial dividing line in the middle, causing the two sides

of the coil to see opposite magnetic fields and produce forces that add instead of canceling. Currents along the top and bottom of the coil produce radial forces that do not rotate the head.

Error Handling

Modern drives also make extensive use of Error Correcting Codes (ECCs), particularly Reed–Solomon error correction. These techniques store extra bits for each block of data that are determined by mathematical formulas. The extra bits allow many errors to be fixed. While these extra bits take up space on the hard drive, they allow higher recording densities to be employed, resulting in much larger storage capacity for user data. In 2009, in the newest drives, low-density parity-check codes (LDPC) are supplanting Reed-Solomon. LDPC codes enable performance close to the Shannon Limit and thus allow for the highest storage density available.

Typical hard drives attempt to “remap” the data in a physical sector that is going bad to a spare physical sector—hopefully while the errors in that bad sector are still few enough that the ECC can recover the data without loss. The S.M.A.R.T. system counts the total number of errors in the entire hard drive fixed by ECC, and the total number of remappings, in an attempt to predict hard drive failure.

Future Development

Because of bit-flipping errors and other issues, perpendicular recording densities may be supplanted by other magnetic recording technologies. Toshiba is promoting bit-patterned recording (BPR), while Xyratex are developing heat-assisted magnetic recording (HAMR).

CAPACITY

Capacity Measurements

Hard disk manufacturers quote disk capacity in multiples of SI-standard powers of 1000, where a *terabyte* is 1000 gigabytes and a *gigabyte* is 1000 megabytes. With file systems that report capacity in powers of 1024, available space appears somewhat less than advertised capacity. The discrepancy between the two methods of reporting sizes had serious financial consequences for at least one hard drive manufacturer when a class action suit argued the different methods effectively misled consumers.

Semiconductor memory chips are organized so that memory sizes are expressed in multiples of powers of two. Hard disks by contrast have no inherent binary size. Capacity is the product of the number of heads, number of tracks, number of sectors per track, and the size of each sector. Sector sizes are standardized for convenience at 256 or 512 and more recently 4096 bytes, which are powers of two. This can cause some confusion because operating systems may report the formatted capacity of a hard drive using binary prefix units which increment by powers of 1024. For example, Microsoft Windows reports disk capacity both in a decimal integer to 12 or more digits and in binary prefix units to three significant digits.

A one terabyte (1 TB) disk drive would be expected to hold around 1 trillion bytes (1,000,000,000,000) or 1000 GB; and indeed most 1 TB hard drives will contain slightly more than this number. However some operating system utilities

would report this as around 931 GB or 953,674 MB. (The actual number for a formatted capacity will be somewhat smaller still, depending on the file system.) Following are the several ways of reporting one Terabyte.

Addressing Data on Large Drives

The capacity of an HDD can be calculated by multiplying the number of cylinders by the number of heads by the number of sectors by the number of bytes/sector (most commonly 512). Drives with the ATA interface and a capacity of eight gigabytes or more behave as if they were structured into 16383 cylinders, 16 heads, and 63 sectors, for compatibility with older operating systems. Unlike in the 1980s, the cylinder, head, sector (C/H/S) counts reported to the CPU by a modern ATA drive are no longer actual physical parameters since the reported numbers are constrained by historic operating-system interfaces and with zone bit recording the actual number of sectors varies by zone. Disks with SCSI interface address each sector with a unique integer number; the operating system remains ignorant of their head or cylinder count.

The old C/H/S scheme has been replaced by logical block addressing. In some cases, to try to “force-fit” the C/H/S scheme to large-capacity drives, the number of heads was given as 64, although no modern drive has anywhere near 32 platters.

Not all the space on a hard drive is available for user files. The operating system file system uses some of the disk space to organize files on the disk, recording their file names and the sequence of disk areas that represent the

file. Examples of data structures stored on disk to retrieve files include the MS DOS file allocation table (FAT), and UNIX inodes, as well as other operating system data structures. This file system overhead is usually less than 1% on drives larger than 100 MB.

For RAID drives, data integrity and fault-tolerance requirements also reduce the realized capacity. For example, a RAID1 drive will be about half the total capacity as a result of data mirroring. For RAID5 drives with x drives you would lose $1/x$ of your space to parity. RAID drives are multiple drives that appear to be one drive to the user, but provides some fault-tolerance. A general rule of thumb to quickly convert the manufacturer's hard disk capacity to the standard Microsoft Windows formatted capacity is $0.93 \times \text{capacity of HDD from manufacturer}$ for HDDs less than a terabyte and $0.91 \times \text{capacity of HDD from manufacturer}$ for HDDs equal to or greater than 1 terabyte.

HDD Formatting

The presentation of an HDD to its host is determined by its controller. This may differ substantially from the drive's native interface particularly in mainframes or servers.

Modern HDDs, such as SAS and SATA drives, appear at their interfaces as a contiguous set of logical blocks; typically 512 bytes long but the industry is in the process of changing to 4,096 byte logical blocks; see Advanced Format.

The process of initializing these logical blocks on the physical disk platters is called *low level formatting* which is usually performed at the factory and is not normally changed in the field. *High level formatting* then writes the

file system structures into selected logical blocks to make the remaining logical blocks available to the host OS and its applications.

Form Factors

Mainframe and minicomputer hard disks were of widely varying dimensions, typically in free standing cabinets the size of washing machines (e.g. HP 7935 and DEC RP06 Disk Drives) or designed so that dimensions enabled placement in a 19" rack (e.g. Diablo Model 31). In 1962, IBM introduced its model 1311 disk, which used 14 inch (nominal size) platters. This became a standard size for mainframe and minicomputer drives for many years, but such large platters were never used with microprocessor-based systems.

With increasing sales of microcomputers having built in floppy-disk drives (FDDs), HDDs that would fit to the FDD mountings became desirable, and this led to the evolution of the market towards drives with certain Form factors, initially derived from the sizes of 8-inch, 5.25-inch, and 3.5-inch floppy disk drives. Smaller sizes than 3.5 inches have emerged as popular in the marketplace and/or been decided by various industry groups.

- 8 inch: 9.5 in × 4.624 in × 14.25 in (241.3 mm × 117.5 mm × 362 mm)

In 1979, Shugart Associates' SA1000 was the first form factor compatible HDD, having the same dimensions and a compatible interface to the 83 FDD.

- 5.25 inch: 5.75 in × 3.25 in × 8 in (146.1 mm × 82.55 mm × 203 mm)

This smaller form factor, first used in an HDD by Seagate in 1980, was the same size as full-height 5¹/₄-inch-diameter (130 mm) FDD, 3.25-inches high. This is twice as high as “half height”; i.e., 1.63 in (41.4 mm). Most desktop models of drives for optical 120 mm disks (DVD, CD) use the half height 5¹/₄ dimension, but it fell out of fashion for HDDs. The Quantum Bigfoot HDD was the last to use it in the late 1990s, with “low-profile” (H”25 mm) and “ultra-low-profile” (H”20 mm) high versions.

- 3.5 inch: 4 in × 1 in × 5.75 in (101.6 mm × 25.4 mm × 146 mm) = 376.77344 cm³

This smaller form factor, first used in an HDD by Rodime in 1983, was the same size as the “half height” 3¹/₂ FDD, i.e., 1.63 inches high. Today it has been largely superseded by 1-inch high “slimline” or “low-profile” versions of this form factor which is used by most desktop HDDs.

- 2.5 inch: 2.75 in × 0.275–0.59 in × 3.945 in (69.85 mm × 7–15 mm × 100 mm) = 48.895–104.775 cm³

This smaller form factor was introduced by PrairieTek in 1988; there is no corresponding FDD. It is widely used today for hard-disk drives in mobile devices (laptops, music players, etc.) and as of 2008 replacing 3.5 inch enterprise-class drives. It is also used in the Playstation 3 and Xbox 360 video game consoles. Today, the dominant height of this form factor is 9.5 mm for laptop drives (usually having two platters inside), but higher capacity drives have a height of

12.5 mm (usually having three platters). Enterprise-class drives can have a height up to 15 mm. Seagate has released a wafer-thin 7mm drive aimed at entry level laptops and high end netbooks in December 2009.

- 1.8 inch: $54 \text{ mm} \times 8 \text{ mm} \times 71 \text{ mm} = 30.672 \text{ cm}^3$
This form factor, originally introduced by Integral Peripherals in 1993, has evolved into the ATA-7 LIF with dimensions as stated. It was increasingly used in digital audio players and subnotebooks, but is rarely used today. An original variant exists for 2–5GB sized HDDs that fit directly into a PC card expansion slot. These became popular for their use in iPods and other HDD based MP3 players.
- 1 inch: $42.8 \text{ mm} \times 5 \text{ mm} \times 36.4 \text{ mm}$
This form factor was introduced in 1999 as IBM's Microdrive to fit inside a CF Type II slot. Samsung calls the same form factor "1.3 inch" drive in its product literature.
- 0.85 inch: $24 \text{ mm} \times 5 \text{ mm} \times 32 \text{ mm}$
Toshiba announced this form factor in January 2004 for use in mobile phones and similar applications, including SD/MMC slot compatible HDDs optimized for video storage on 4G handsets. Toshiba currently sells a 4 GB (MK4001MTD) and 8 GB (MK8003MTD) version and holds the Guinness World Record for the smallest hard disk drive.
- 3.5-inch and 2.5-inch hard disks currently dominate the market.

By 2009 all manufacturers had discontinued the development of new products for the 1.3-inch, 1-inch and 0.85-inch form factors due to falling prices of flash memory, which is slightly more stable and resistant to damage from impact and/or dropping.

The inch-based nickname of all these form factors usually do not indicate any actual product dimension (which are specified in millimeters for more recent form factors), but just roughly indicate a size relative to disk diameters, in the interest of historic continuity.

DISK FAILURES AND THEIR METRICS

Most major hard disk and motherboard vendors now support S.M.A.R.T. (Self-Monitoring, Analysis, and Reporting Technology), which measures drive characteristics such as operating temperature, spin-up time, data error rates, etc. Certain trends and sudden changes in these parameters are thought to be associated with increased likelihood of drive failure and data loss.

However, not all failures are predictable. Normal use eventually can lead to a breakdown in the inherently fragile device, which makes it essential for the user to periodically back up the data onto a separate storage device. Failure to do so can lead to the loss of data. While it may sometimes be possible to recover lost information, it is normally an extremely costly procedure, and it is not possible to guarantee success. A 2007 study published by Google suggested very little correlation between failure rates and either high temperature or activity level; however, the correlation between

manufacturer/model and failure rate was relatively strong. Statistics in this matter is kept highly secret by most entities. Google did not publish the manufacturer's names along with their respective failure rates, though they have since revealed that they use Hitachi Deskstar drives in some of their servers. While several S.M.A.R.T. parameters have an impact on failure probability, a large fraction of failed drives do not produce predictive S.M.A.R.T. parameters. S.M.A.R.T. parameters alone may not be useful for predicting individual drive failures.

A common misconception is that a colder hard drive will last longer than a hotter hard drive. The Google study seems to imply the reverse—"lower temperatures are associated with higher failure rates". Hard drives with S.M.A.R.T.-reported average temperatures below 27 °C (80.6 °F) had higher failure rates than hard drives with the highest reported average temperature of 50 °C (122 °F), failure rates at least twice as high as the optimum S.M.A.R.T.-reported temperature range of 36 °C (96.8 °F) to 47 °C (116.6 °F).

SCSI, SAS, and FC drives are typically more expensive and are traditionally used in servers and disk arrays, whereas inexpensive ATA and SATA drives evolved in the home computer market and were perceived to be less reliable. This distinction is now becoming blurred.

The mean time between failures (MTBF) of SATA drives is usually about 600,000 hours (some drives such as Western Digital Raptor have rated 1.4 million hours MTBF), while SCSI drives are rated for upwards of 1.5 million hours.

However, independent research indicates that MTBF is not a reliable estimate of a drive's longevity. MTBF is conducted in laboratory environments in test chambers and is an important metric to determine the quality of a disk drive before it enters high volume production. Once the drive product is in production, the more valid metric is annualized failure rate (AFR). AFR is the percentage of real-world drive failures after shipping.

SAS drives are comparable to SCSI drives, with high MTBF and high reliability.

Enterprise S-ATA drives designed and produced for enterprise markets, unlike standard S-ATA drives, have reliability comparable to other enterprise class drives.

Typically enterprise drives (all enterprise drives, including SCSI, SAS, enterprise SATA, and FC) experience between 0.70%–0.78% annual failure rates from the total installed drives.

Eventually all mechanical hard disk drives fail, so to mitigate loss of data, some form of redundancy is needed, such as RAID or a regular backup system.

External Removable Drives

External removable hard disk drives connect to the computer using a USB cable or other means. External drives are used for:

- Backup of files and information
- Data recovery
- Disk cloning
- Running virtual machines

- Scratch disk for video editing applications and video recording.

Larger models often include full-sized 3.5" PATA or SATA desktop hard drives. Features such as biometric security or multiple interfaces generally increase cost.

Market Segments

- As of July 2010, the highest capacity consumer HDDs are 3 TB.
- “Desktop HDDs” typically store between 120 GB and 2 TB and rotate at 5,400 to 10,000 rpm, and have a media transfer rate of 0.5 Gbit/s or higher. (1 GB = 10^9 bytes; 1 Gbit/s = 10^9 bit/s)
- Enterprise HDDs are typically used with multiple-user computers running enterprise software. Examples are
 - o transaction processing databases;
 - o internet infrastructure (email, webserver, e-commerce);
 - o scientific computing software;
 - o nearline storage management software.

The fastest enterprise HDDs spin at 10,000 or 15,000 rpm, and can achieve sequential media transfer speeds above 1.6 Gbit/s. and a sustained transfer rate up to 1 Gbit/s. Drives running at 10,000 or 15,000 rpm use smaller platters to mitigate increased power requirements (as they have less air drag) and therefore generally have lower capacity than the highest capacity desktop drives.

Enterprise drives commonly operate continuously (“24/7”) in demanding environments while delivering the highest

possible performance without sacrificing reliability. Maximum capacity is not the primary goal, and as a result the drives are often offered in capacities that are relatively low in relation to their cost.

- Mobile HDDs or laptop HDDs, smaller than their desktop and enterprise counterparts, tend to be slower and have lower capacity. A typical mobile HDD spins at either 4200 rpm, 5200 rpm, 5400 rpm, or 7200 rpm, with 5400 rpm being the most prominent. 7200 rpm drives tend to be more expensive and have smaller capacities, while 4200 rpm models usually have very high storage capacities. Because of smaller platter(s), mobile HDDs generally have lower capacity than their greater desktop counterparts.

The exponential increases in disk space and data access speeds of HDDs have enabled the commercial viability of consumer products that require large storage capacities, such as digital video recorders and digital audio players. In addition, the availability of vast amounts of cheap storage has made viable a variety of web-based services with extraordinary capacity requirements, such as free-of-charge web search, web archiving, and video sharing (Google, Internet Archive, YouTube, etc.).

Sales

Worldwide revenue from shipments of HDDs is expected to reach \$27.7 billion in 2010, up 18.4% from \$23.4 billion in 2009 corresponding to a 2010 unit shipment forecast of 674.6 million compared to 549.5 million units in 2009.

Icons

Hard drives are traditionally symbolized as either a stylized stack of platters (in orthographic projection) or, more abstractly, as a cylinder. This is particularly found in schematic diagrams or on indicator lights, as on laptops, to indicate hard drive access. In most modern operating systems, hard drives are instead represented by an illustration or photograph of a hard drive enclosure.

Manufacturers

More than 200 companies have manufactured hard disk drives. Today most drives are made by Seagate, Western Digital, Hitachi, Samsung, and Toshiba (though Toshiba does not manufacture 3.5 inch drives).

5

Computer File System

A file system (sometimes written as filesystem) is a method of storing and organizing computer files and their data. Essentially, it organizes these files into a database for the storage, organization, manipulation, and retrieval by the computer's operating system. File systems are used on data storage devices such as hard disks or CD-ROMs to maintain the physical location of the files. Beyond this, they might provide access to data on a file server by acting as clients for a network protocol (e.g., NFS, SMB, or 9P clients), or they may be virtual and exist only as an access method for virtual data (e.g., procfs). It is distinguished from a directory service and registry.

Aspects of File Systems

Most file systems make use of an underlying data storage device that offers access to an array of fixed-size physical

sectors, generally a power of 2 in size (512 bytes or 1, 2, or 4 KiB are most common). The file system is responsible for organizing these sectors into files and directories, and keeping track of which sectors belong to which file and which are not being used. Most file systems address data in fixed-sized units called “clusters” or “blocks” which contain a certain number of disk sectors (usually 1-64). This is the smallest amount of disk space that can be allocated to hold a file. However, file systems need not make use of a storage device at all. A file system can be used to organize and represent access to any data, whether it is stored or dynamically generated (e.g., procs).

File Names

A file name (or filename) is a name assigned to a file in order to secure storage location in the computer memory. Whether the file system has an underlying storage device or not, file systems typically have directories which associate file names with files, usually by connecting the file name to an index in a file allocation table of some sort, such as the FAT in a DOS file system, or an inode in a Unix-like file system. Directory structures may be flat, or allow hierarchies where directories may contain subdirectories. In some file systems, file names are structured, with special syntax for filename extensions and version numbers. In others, file names are simple strings, and per-file metadata is stored elsewhere.

Metadata

Other bookkeeping information is typically associated with each file within a file system. The length of the data

contained in a file may be stored as the number of blocks allocated for the file or as an exact byte count. The time that the file was last modified may be stored as the file's timestamp. Some file systems also store the file creation time, the time it was last accessed, and the time the file's meta-data was changed. (Note that many early PC operating systems did not keep track of file times.) Other information can include the file's device type (e.g., block, character, socket, subdirectory, etc.), its owner user ID and group ID, and its access permission settings (e.g., whether the file is read-only, executable, etc.).

Arbitrary attributes can be associated on advanced file systems, such as NTFS, XFS, ext2/ext3, some versions of UFS, and HFS+, using extended file attributes. This feature is implemented in the kernels of Linux, FreeBSD and Mac OS X operating systems, and allows metadata to be associated with the file at the *file system* level.

This, for example, could be the author of a document, the character encoding of a plain-text document, or a checksum.

Hierarchical File Systems

The hierarchical file system (not to be confused with Apple's HFS) was an early research interest of Dennis Ritchie of Unix fame; previous implementations were restricted to only a few levels, notably the IBM implementations, even of their early databases like IMS. After the success of Unix, Ritchie extended the file system concept to every object in his later operating system developments, such as Plan 9 and Inferno.

Facilities

Traditional file systems offer facilities to create, move and delete both files and directories. They lack facilities to create additional links to a directory (hard links in Unix), rename parent links (“..” in Unix-like OS), and create bidirectional links to files. Traditional file systems also offer facilities to truncate, append to, create, move, delete and in-place modify files. They do not offer facilities to prepend to or truncate from the beginning of a file, let alone arbitrary insertion into or deletion from a file. The operations provided are highly asymmetric and lack the generality to be useful in unexpected contexts. For example, interprocess pipes in Unix have to be implemented outside of the file system because the pipes concept does not offer truncation from the beginning of files.

Secure Access

Secure access to basic file system operations can be based on a scheme of access control lists or capabilities. Research has shown access control lists to be difficult to secure properly, which is why research operating systems tend to use capabilities. Commercial file systems still use access control lists.

Types of File Systems

File system types can be classified into disk file systems, network file systems and special purpose file systems.

Disk File Systems

A disk file system is a file system designed for the storage of files on a data storage device, most commonly a disk

drive, which might be directly or indirectly connected to the computer. Examples of disk file systems include FAT (FAT12, FAT16, FAT32, exFAT), NTFS, HFS and HFS+, HPFS, UFS, ext2, ext3, ext4, btrfs, ISO 9660, ODS-5, Veritas File System, VMFS, ZFS, ReiserFS and UDF. Some disk file systems are journaling file systems or versioning file systems.

Optical Discs

ISO 9660 and Universal Disk Format (UDF) are the two most common formats that target Compact Discs, DVDs and Blu-ray discs. Mount Rainier is a newer extension to UDF supported by Linux 2.6 series and Windows Vista that facilitates rewriting to DVDs in the same fashion as has been possible with floppy disks.

Flash File Systems

A *flash file system* is a file system designed for storing files on flash memory devices. These are becoming more prevalent as the number of mobile devices is increasing, and the capacity of flash memories increase. While a disk file system can be used on a flash device, this is suboptimal for several reasons:

- Erasing blocks: Flash memory blocks have to be explicitly erased before they can be rewritten. The time taken to erase blocks can be significant, thus it is beneficial to erase unused blocks while the device is idle.
- Random access: Disk file systems are optimized to avoid disk seeks whenever possible, due to the high cost of seeking. Flash memory devices impose no seek latency.

- **Wear levelling:** Flash memory devices tend to wear out when a single block is repeatedly overwritten; flash file systems are designed to spread out writes evenly.

Log-structured file systems have many of the desirable properties for a flash file system. Such file systems include JFFS2 and YAFFS.

No Need for Defragmenting

Because flash media effectively has zero seek time, defragmentation of flash media is unnecessary and has no performance benefit. Instead defragmentation is detrimental to the life of the media since it wears out the data storage cells with unnecessary writing.

Media is Typically pPre-formatted

Due to the problems of limited write cycles per data cell, it is generally not recommended to format flash storage devices since it reduces the life of all cells on the media. Flash devices are typically sold with a common type of file system already created on the media, to remove the need for formatting.

Tape File Systems

A *tape file system* is a file system and tape format designed to store files on tape in a self-describing form. Magnetic tapes are sequential storage media with significantly longer random data access times than disks, posing challenges to the creation and efficient management of a general-purpose file system. In a disk file system there is typically a master file directory, and a map of used and free data regions. Any

file additions, changes, or removals require updating the directory and the used/free maps. Random access to data regions is measured in milliseconds so this system works well for disks. However, tape requires linear motion to wind and unwind potentially very long reels of media, and this tape motion may take several seconds to several minutes to move the read/write head from one end of the tape to the other.

Consequently, a master file directory and usage map can be extremely slow and inefficient with tape. Writing typically involves reading the block usage map to find free blocks for writing, updating the usage map and directory to add the data, and then advancing the tape to write the data in the correct spot. Each additional file write requires updating the map and directory and writing the data, which may take several seconds to occur for each file.

Tape file systems instead typically allow for the file directory to be spread across the tape intermixed with the data, referred to as *streaming*, so that time-consuming and repeated tape motions are not required to write new data.

However a side effect of this design is that reading the file directory of a tape usually requires scanning the entire tape to read all the scattered directory entries. Most data archiving software that works with tape storage will store a local copy of the tape catalog on a disk file system, so that adding files to a tape can be done quickly without having to rescan the tape media. The local tape catalog copy is usually discarded if not used for a specified period of time, at which point the tape must be re-scanned if it is to be used in the future.

IBM has recently announced a new file system for tape called the Linear Tape File System. The IBM implementation of this file system has been released as the open-source IBM Long Term File System product.

Tape Formatting

Writing data to a tape is often a significantly time-consuming process that may take several hours. Similarly, completely erasing or formatting a tape can also take several hours. With many data tape technologies it is not necessary to format the tape before over-writing new data to the tape. This is due to the inherently destructive nature of overwriting data on sequential media. Because of the time it can take to format a tape, typically tapes are pre-formatted so that the tape user does not need to spend time preparing each new tape for use. All that is usually necessary is to write an identifying media label to the tape before use, and even this can be automatically written by software when a new tape is used for the first time.

Tape Booting

Tapes can sometimes be used to boot a computer, but the long random data access delays can make this a very slow and tedious process. The boot-file storage can be optimized so that programme data is stored in the order required for use, with minimal need for seeking. Typically tape booting is only used for infrequent special purposes such as operating system installation or for system recovery/restore.

Database File Systems

A recent concept for file management is the idea of a

database-based file system. Instead of, or in addition to, hierarchical structured management, files are identified by their characteristics, like type of file, topic, author, or similar metadata.

Transactional File Systems

Some programmes need to update multiple files “all at once.” For example, a software installation may write programme binaries, libraries, and configuration files. If the software installation fails, the programme may be unusable. If the installation is upgrading a key system utility, such as the command shell, the entire system may be left in an unusable state. Transaction processing introduces the isolation guarantee, which states that operations within a transaction are hidden from other threads on the system until the transaction commits, and that interfering operations on the system will be properly serialized with the transaction. Transactions also provide the atomicity guarantee, that operations inside of a transaction are either all committed, or the transaction can be aborted and the system discards all of its partial results. This means that if there is a crash or power failure, after recovery, the stored state will be consistent. Either the software will be completely installed or the failed installation will be completely rolled back, but an unusable partial install will not be left on the system.

Windows, beginning with Vista, added transaction support to NTFS, abbreviated TxF. TxF is the only commercial implementation of a transactional file system, as transactional file systems are difficult to implement correctly in practice. There are a number of research prototypes of

transactional file systems for UNIX systems, including the Valor file system, Amino, LFS, and a transactional ext3 file system on the TxOS kernel, as well as transactional file systems targeting embedded systems, such as TFFS.

Ensuring consistency across multiple file system operations is difficult, if not impossible, without file system transactions. File locking can be used as a concurrency control mechanism for individual files, but it typically does not protect the directory structure or file metadata. For instance, file locking cannot prevent TOCTTOU race conditions on symbolic links. File locking also cannot automatically roll back a failed operation, such as a software upgrade; this requires atomicity. Journaling file systems are one technique used to introduce transaction-level consistency to file system structures. Journal transactions are not exposed to programmes as part of the OS API; they are only used internally to ensure consistency at the granularity of a single system call.

Network File Systems

A network file system is a file system that acts as a client for a remote file access protocol, providing access to files on a server. Examples of network file systems include clients for the NFS, AFS, SMB protocols, and file-system-like clients for FTP and WebDAV.

Shared Disk File Systems

A shared disk file system is one in which a number of machines (usually servers) all have access to the same external disk subsystem (usually a SAN). The file system arbitrates access to that subsystem, preventing write

collisions. Examples include GFS from Red Hat, GPFS from IBM, and SFS from DataFlow.

Special Purpose File Systems

A special purpose file system is basically any file system that is not a disk file system or network file system. This includes systems where the files are arranged dynamically by software, intended for such purposes as communication between computer processes or temporary file space. Special purpose file systems are most commonly used by file-centric operating systems such as Unix. Examples include the `procfs (/proc)` file system used by some Unix variants, which grants access to information about processes and other operating system features.

Deep space science exploration craft, like Voyager I and II used digital tape-based special file systems. Most modern space exploration craft like Cassini-Huygens used Real-time operating system file systems or RTOS influenced file systems. The Mars Rovers are one such example of an RTOS file system, important in this case because they are implemented in flash memory.

File Systems and Operating Systems

Most operating systems provide a file system, as a file system is an integral part of any modern operating system. Early microcomputer operating systems' only real task was file management — a fact reflected in their names. Some early operating systems had a separate component for handling file systems which was called a disk operating system. On some microcomputers, the disk operating system was loaded separately from the rest of the operating system.

On early operating systems, there was usually support for only one, native, unnamed file system; for example, CP/M supports only its own file system, which might be called “CP/M file system” if needed, but which didn’t bear any official name at all. Because of this, there needs to be an interface provided by the operating system software between the user and the file system. This interface can be textual (such as provided by a command line interface, such as the Unix shell, or OpenVMS DCL) or graphical (such as provided by a graphical user interface, such as file browsers). If graphical, the metaphor of the *folder*, containing documents, other files, and nested folders is often used.

Flat File Systems

In a flat file system, there are no subdirectories—everything is stored at the same (root) level on the media, be it a hard disk, floppy disk, etc. While simple, this system rapidly becomes inefficient as the number of files grows, and makes it difficult for users to organize data into related groups. Like many small systems before it, the original Apple Macintosh featured a flat file system, called Macintosh File System. Its version of Mac OS was unusual in that the file management software (Macintosh Finder) created the illusion of a partially hierarchical filing system on top of EMFS. This structure meant that every file on a disk had to have a unique name, even if it appeared to be in a separate folder. MFS was quickly replaced with Hierarchical File System, which supported real directories. A recent addition to the flat file system family is Amazon’s S3, a remote storage service, which is intentionally simplistic to

allow users the ability to customize how their data is stored. The only constructs are buckets (imagine a disk drive of unlimited size) and objects (similar, but not identical to the standard concept of a file). Advanced file management is allowed by being able to use nearly any character (including '/') in the object's name, and the ability to select subsets of the bucket's content based on identical prefixes.

File Systems under Unix-like Operating Systems

Unix-like operating systems create a virtual file system, which makes all the files on all the devices appear to exist in a single hierarchy. This means, in those systems, there is one root directory, and every file existing on the system is located under it somewhere. Unix-like systems can use a RAM disk or network shared resource as its root directory. Unix-like systems assign a device name to each device, but this is not how the files on that device are accessed. Instead, to gain access to files on another device, the operating system must first be informed where in the directory tree those files should appear. This process is called mounting a file system. For example, to access the files on a CD-ROM, one must tell the operating system "Take the file system from this CD-ROM and make it appear under such-and-such directory". The directory given to the operating system is called the *mount point* – it might, for example, be /media. The /media directory exists on many Unix systems (as specified in the Filesystem Hierarchy Standard) and is intended specifically for use as a mount point for removable media such as CDs, DVDs, USB drives or floppy disks. It may be empty, or it may contain subdirectories for mounting

individual devices. Generally, only the administrator (i.e. root user) may authorize the mounting of file systems.

Unix-like operating systems often include software and tools that assist in the mounting process and provide it new functionality. Some of these strategies have been coined “auto-mounting” as a reflection of their purpose.

1. In many situations, file systems other than the root need to be available as soon as the operating system has booted. All Unix-like systems therefore provide a facility for mounting file systems at boot time. System administrators define these file systems in the configuration file *fstab* (*vfstab* in Solaris), which also indicates options and mount points.
2. In some situations, there is no need to mount certain file systems at boot time, although their use may be desired thereafter. There are some utilities for Unix-like systems that allow the mounting of predefined file systems upon demand.
3. Removable media have become very common with microcomputer platforms. They allow programmes and data to be transferred between machines without a physical connection. Common examples include USB flash drives, CD-ROMs, and DVDs. Utilities have therefore been developed to detect the presence and availability of a medium and then mount that medium without any user intervention.
4. Progressive Unix-like systems have also introduced a concept called supermounting; see, for example, the Linux supermount-ng project. For example, a

floppy disk that has been supermounted can be physically removed from the system. Under normal circumstances, the disk should have been synchronized and then unmounted before its removal. Provided synchronization has occurred, a different disk can be inserted into the drive. The system automatically notices that the disk has changed and updates the mount point contents to reflect the new medium. Similar functionality is found on Windows machines.

5. An automounter will automatically mount a file system when a reference is made to the directory atop which it should be mounted. This is usually used for file systems on network servers, rather than relying on events such as the insertion of media, as would be appropriate for removable media.

File Systems Under Linux

Linux supports many different file systems, but common choices for the system disk include the ext* family (such as ext2, ext3 and ext4), XFS, JFS, ReiserFS and btrfs.

File Systems Under Solaris

The Sun Microsystems Solaris operating system in earlier releases defaulted to (non-journaled or non-logging) UFS for bootable and supplementary file systems. Solaris defaulted to, supported, and extended UFS. Support for other file systems and significant enhancements were added over time, including Veritas Software Corp. (Journaling) VxFS, Sun Microsystems (Clustering) QFS, Sun Microsystems

(Journaling) UFS, and Sun Microsystems (open source, poolable, 128 bit compressible, and error-correcting) ZFS. Kernel extensions were added to Solaris to allow for bootable Veritas VxFS operation. Logging or Journaling was added to UFS in Sun's Solaris 7. Releases of Solaris 10, Solaris Express, OpenSolaris, and other open source variants of the Solaris operating system later supported bootable ZFS. Logical Volume Management allows for spanning a file system across multiple devices for the purpose of adding redundancy, capacity, and/or throughput. Legacy environments in Solaris may use Solaris Volume Manager (formerly known as Solstice DiskSuite.) Multiple operating systems (including Solaris) may use Veritas Volume Manager. Modern Solaris based operating systems eclipse the need for Volume Management through leveraging virtual storage pools in ZFS.

File Systems Under Mac OS X

Mac OS X uses a file system that it inherited from classic Mac OS called HFS Plus, sometimes called *Mac OS Extended*. HFS Plus is a metadata-rich and case preserving file system. Due to the Unix roots of Mac OS X, Unix permissions were added to HFS Plus. Later versions of HFS Plus added journaling to prevent corruption of the file system structure and introduced a number of optimizations to the allocation algorithms in an attempt to defragment files automatically without requiring an external defragmenter. Filenames can be up to 255 characters. HFS Plus uses Unicode to store filenames. On Mac OS X, the filetype can come from the type code, stored in file's metadata, or the filename. HFS Plus has three kinds of links: Unix-style hard links, Unix-

style symbolic links and aliases. Aliases are designed to maintain a link to their original file even if they are moved or renamed; they are not interpreted by the file system itself, but by the File Manager code in userland. Mac OS X also supports the UFS file system, derived from the BSD Unix Fast File System via NeXTSTEP. However, as of Mac OS X 10.5 (Leopard), Mac OS X can no longer be installed on a UFS volume, nor can a pre-Leopard system installed on a UFS volume be upgraded to Leopard. Newer versions Mac OS X are capable of reading and writing to the legacy FAT file systems(16 & 32). They are capable of reading, but not writing to the NTFS file system. Third party software is still necessary to write to the NTFS file system under Snow Leopard 10.6.4.

File Systems under Plan 9 from Bell Labs

Plan 9 from Bell Labs was originally designed to extend some of Unix's good points, and to introduce some new ideas of its own while fixing the shortcomings of Unix. With respect to file systems, the Unix system of treating things as files was continued, but in Plan 9, *everything* is treated as a file, and accessed as a file would be (i.e., no ioctl or mmap). Perhaps surprisingly, while the file interface is made universal it is also simplified considerably: symlinks, hard links and suid are made obsolete, and an atomic create/open operation is introduced. More importantly the set of file operations becomes well defined and subversions of this like ioctl are eliminated. Secondly, the underlying 9P protocol was used to remove the difference between local and remote files (except for a possible difference in latency

or in throughput). This has the advantage that a device or devices, represented by files, on a remote computer could be used as though it were the local computer's own device(s). This means that under Plan 9, multiple file servers provide access to devices, classing them as file systems. Servers for "synthetic" file systems can also run in user space bringing many of the advantages of micro kernel systems while maintaining the simplicity of the system.

Everything on a Plan 9 system has an abstraction as a file; networking, graphics, debugging, authentication, capabilities, encryption, and other services are accessed via I-O operations on file descriptors. For example, this allows the use of the IP stack of a gateway machine without need of NAT, or provides a network-transparent window system without the need of any extra code. Another example: a Plan-9 application receives FTP service by opening an FTP site. The `ftpf`s server handles the open by essentially mounting the remote FTP site as part of the local file system. With `ftpf`s as an intermediary, the application can now use the usual file-system operations to access the FTP site as if it were part of the local file system. A further example is the mail system which uses file servers that synthesize virtual files and directories to represent a user mailbox as `/mail/fs/mbox`. The `wikifs` provides a file system interface to a wiki. These file systems are organized with the help of private, per-process namespaces, allowing each process to have a different view of the many file systems that provide resources in a distributed system. The Inferno operating system shares these concepts with Plan 9.

File Systems under Microsoft Windows

Windows makes use of the FAT and NTFS file systems. Unlike many other operating systems, Windows uses a *drive letter* abstraction at the user level to distinguish one disk or partition from another. For example, the path C:\WINDOWS represents a directory WINDOWS on the partition represented by the letter C. The C drive is most commonly used for the primary hard disk partition, on which Windows is usually installed and from which it boots. This “tradition” has become so firmly ingrained that bugs came about in older applications which made assumptions that the drive that the operating system was installed on was C. The tradition of using “C” for the drive letter can be traced to MS-DOS, where the letters A and B were reserved for up to two floppy disk drives. This in turn derived from CP/M in the 1970s, and ultimately from IBM’s CP/CMS of 1967. Network drives may also be mapped to drive letters.

FAT

The File Allocation Table (FAT) filing system, supported by all versions of Microsoft Windows, was an evolution of that used in Microsoft’s earlier operating system (MS-DOS which in turn was based on 86-DOS). FAT ultimately traces its roots back to the short-lived M-DOS project and Standalone disk BASIC before it. Over the years various features have been added to it, inspired by similar features found on file systems used by operating systems such as Unix. Older versions of the FAT file system (FAT12 and FAT16) had file name length limits, a limit on the number

of entries in the root directory of the file system and had restrictions on the maximum size of FAT-formatted disks or partitions. Specifically, FAT12 and FAT16 had a limit of 8 characters for the file name, and 3 characters for the extension (such as .exe). This is commonly referred to as the 8.3 filename limit. VFAT, which was an extension to FAT12 and FAT16 introduced in Windows NT 3.5 and subsequently included in Windows 95, allowed long file names (LFN). FAT32 also addressed many of the limits in FAT12 and FAT16, but remains limited compared to NTFS. exFAT (also known as FAT64) is the newest iteration of FAT, with certain advantages over NTFS with regards to file system overhead. exFAT is only compatible with newer Windows systems, such as Windows 2003, Windows Vista, Windows 2008, Windows 7 and more recently, support has been added for WinXP.

NTFS

NTFS, introduced with the Windows NT operating system, allowed ACL-based permission control. Hard links, multiple file streams, attribute indexing, quota tracking, sparse files, encryption, compression, reparse points (directories working as mount-points for other file systems, symlinks, junctions, remote storage links) are also supported, though not all these features are well-documented.

Other File Systems

- The Prospero File System is a file system based on the Virtual System Model. The system was created by Dr. B. Clifford Neuman of the Information Sciences Institute at the University of Southern California.

Computer Memory System

- RSRE FLEX file system - written in ALGOL 68
- The file system of the Michigan Terminal System (MTS) is interesting because: (i) it provides “line files” where record lengths and line numbers are associated as metadata with each record in the file, lines can be added, replaced, updated with the same or different length records, and deleted anywhere in the file without the need to read and rewrite the entire file; (ii) using programme keys files may be shared or permitted to commands and programmes in addition to users and groups; and (iii) there is a comprehensive file locking mechanism that protects both the file’s data and its metadata.

6

Computer External Memory System

MAGNETIC DISKS

A hard drive is a special disk that is usually mounted permanently inside your computer's cabinet. You rarely see the hard drive, and almost never take it out. Hard drives are made of different material than floppies, and they are physically hard (although if you touched the actual hard surface, you would destroy it!) They spin much more quickly than floppies, and require much more precision. They are sealed inside a special case, and that is sealed inside the computer case. A hard drive has a much larger capacity than a floppy, and is much faster at saving and retrieving information. Modern computers frequently have hard drives with 500 MB or more of capacity. As this capacity grows, people are beginning to measure it in terms of gigabytes.

Software programmes are always becoming larger and taking more room on hard drives. It never takes long to completely fill up the capacity of a drive.

READ ONLY MEMORY (ROM)

There is another memory in computer, which is called Read Only Memory (ROM). Again it is the ICs inside the PC that form the ROM. The storage of programme and data in the ROM is permanent. The ROM stores some standard processing programmes supplied by the manufacturers to operate the personal computer.

The ROM can only be read by the CPU but it cannot be changed. The basic input/output programme is stored in the ROM that examines and initializes various equipment attached to the PC when the switch is made ON. The memories, which do not lose their content on failure of power supply, are known as non-volatile memories. ROM is non-volatile memory.

Memory is like the banks of switches we have been thinking about that simply contain 1/0 binary patterns. Read Only Memory (abbreviated ROM) is a special kind of memory that cannot be changed. The most basic instructions for the CPU are built in to the ROM at the factory. To the end user, there really is very little to worry about regarding ROM. The amount of ROM in your computer doesn't really matter to you. All your programmes and information will go into another kind of memory that we will learn about soon.

Definition

ROM is an integrated-circuit memory chip that contains configuration data. ROM is commonly called firmware

because its programming is fully embedded into the ROM chip. As such, ROM is a hardware and software in one.

Because data is fully incorporated at the ROM chip's manufacture, data stored can neither be erased nor replaced. This means permanent and secure data storage. However, if a mistake is made in manufacture, a ROM chip becomes unusable. The most expensive stage of ROM manufacture, therefore, is creating the template. If a template is readily available, duplicating the ROM chip is very easy and affordable.

A ROM chip is also non-volatile so data stored in it is not lost when power is turned off. There is a type of memory that stores data without electrical current; it is the ROM (Read Only Memory) or is sometimes called non-volatile memory as it is not erased when the system is switched off.

This type of memory lets you stored the data needed to start up the computer. Indeed, this information cannot be stored on the hard disk since the disk parameters (vital for its initialisation) are part of these data which are essential for booting.

Different ROM-type memories contain these essential start-up data, that is:

- The BIOS is a programme for controlling the system's main input-output interfaces, hence the name BIOS ROM which is sometimes given to the read-only memory chip of the mother board which hosts it.
- *The bootstrap loader:* a programme for loading (random access) memory into the operating system and launching it. This generally seeks the operating system on the floppy drive then on the hard disk, which

allows the operating system to be launched from a system floppy disk in the event of malfunction of the system installed on the hard disk.

- The CMOS Setup is the screen displayed when the computer starts up and which is used to amend the system parameters (often wrongly referred to as BIOS).
- The Power-On Self Test (POST), a programme that runs automatically when the system is booted, thus allowing the system to be tested (this is why the system "counts" the RAM at start-up).

Given that ROM are much slower than RAM memories (access time for a ROM is around 150 ns whereas for SDRAM it is around 10 ns), the instructions given in the ROM are sometimes copied to the RAM at start-up; this is known as shadowing, though is usually referred to as shadow memory). One major type of memory that is used in PCs is called read-only memory, or ROM for short. ROM is a type of memory that normally can only be read, as opposed to RAM which can be both read and written.

There are two main reasons that read-only memory is used for certain functions within the PC:

- *Permanence:* The values stored in ROM are always there, whether the power is on or not. A ROM can be removed from the PC, stored for an indefinite period of time, and then replaced, and the data it contains will still be there. For this reason, it is called non-volatile storage. A hard disk is also non-volatile, for the same reason, but regular RAM is not.
- *Security:* The fact that ROM cannot easily be modified provides a measure of security against accidental (or

malicious) changes to its contents. You are not going to find viruses infecting true ROMs, for example; it's just not possible. (It's technically possible with erasable EPROMs, though in practice never seen.)

Read-only memory is most commonly used to store system-level programmes that we want to have available to the PC at all times. The most common example is the system BIOS programme, which is stored in a ROM called (amazingly enough) the system BIOS ROM.

Having this in a permanent ROM means it is available when the power is turned on so that the PC can use it to boot up the system. Remember that when you first turn on the PC the system memory is empty, so there has to be something for the PC to use when it starts up.

While the whole point of a ROM is supposed to be that the contents cannot be changed, there are times when being able to change the contents of a ROM can be very useful.

There are several ROM variants that can be changed under certain circumstances; these can be thought of as "mostly read-only memory":^)

Types of ROMs

There are types of ROMs with a description of their relative modifiability.

ROM

Read-only memory (ROM) is a class of storage medium used in computers and other electronic devices. Data stored in ROM cannot be modified, or can be modified only slowly or with difficulty, so it is mainly used to distribute firmware (software that is very closely tied to specific hardware, and

unlikely to need frequent updates). A regular ROM is constructed from hard-wired logic, encoded in the silicon itself, much the way that a processor is. It is designed to perform a specific function and cannot be changed. This is inflexible and so regular ROMs are only used generally for programmes that are static (not changing often) and mass-produced. This product is analagous to a commercial software CD-ROM that you purchase in a store.

Programmable ROM (PROM)

A programmable read-only memory (PROM) or field programmable read-only memory (FEPROM) or one-time programmable non-volatile memory (OTP NVM) is a form of digital memory where the setting of each bit is locked by a fuse or antifuse.

Such PROMs are used to store programmes permanently. The key difference from a strict ROM is that the programming is applied after the device is constructed.

There is another type of primary memory in computer, which is called Programmable Read Only Memory (PROM). You know that it is not possible to modify or erase programmes stored in ROM, but it is possible for you to store your programme in PROM chip. Once the programmes are written it cannot be changed and remain intact even if power is switched off. Therefore programmes or instructions written in PROM or ROM cannot be erased or changed.

This is a type of ROM that can be programmed using special equipment; it can be written to, but only once. This is useful for companies that make their own ROMs from software they write, because when they change their code they can create new PROMs without requiring expensive equipment. This is

similar to the way a CD-ROM recorder works by letting you "burn" programmes onto blanks once and then letting you read from them many times. In fact, programming a PROM is also called burning, just like burning a CD-R, and it is comparable in terms of its flexibility.

Erasable Programmable ROM (EPROM)

This stands for Erasable Programmable Read Only Memory, which over come the problem of PROM and ROM. EPROM chip can be programmed time and again by erasing the information stored earlier in it. Information stored in EPROM exposing the chip for some time ultraviolet light and it erases chip is reprogrammed using a special programming facility. When the EPROM is in use information can only be read.

An EPROM (rarely EROM), or erasable programmable read only memory, is a type of memory chip that retains its data when its power supply is switched off. In other words, it is non-volatile. It is an array of floating-gate transistors individually programmed by an electronic device that supplies higher voltages than those normally used in digital circuits. Once programmed, an EPROM can be erased by exposing it to strong ultraviolet light source (such as from a mercury-vapor light). EPROMs are easily recognizable by the transparent fused quartz window in the top of the package, through which the silicon chip is visible, and which permits exposure to UV light during erasing.

An EPROM is a ROM that can be erased and reprogrammed. A little glass window is installed in the top of the ROM package, through which you can actually see the chip that holds the memory. Ultraviolet light of a specific frequency can be shined through this window for a specified

period of time, which will erase the EPROM and allow it to be reprogrammed again. Obviously this is much more useful than a regular PROM, but it does require the erasing light. Continuing the "CD" analogy, this technology is analogous to a reusable CD-RW.

Electrically Erasable Programmable ROM (EEPROM)

Short for Electrically Erasable Programmable Read-Only Memory, EEPROM is a PROM that can be erased and reprogrammed using an electrical charge that was developed by George Perlegos while at Intel in 1978. Unlike most memory inside a computer, this memory remembers data when the power is turned off.

EEPROM was a replacement for PROM and EPROM chips and is used for later computer's BIOS that were built after 1994. Having a computer with an EEPROM allows a computer user to update the BIOS in their computer without having to open the computer or remove any chips.

The next level of erasability is the EEPROM, which can be erased under software control. This is the most flexible type of ROM, and is now commonly used for holding BIOS programmes.

When you hear reference to a "flash BIOS" or doing a BIOS upgrade by "flashing", this refers to reprogramming the BIOS EEPROM with a special software programme. Here we are blurring the line a bit between what "read-only" really means, but remember that this rewriting is done maybe once a year or so, compared to real read-write memory (RAM) where rewriting is done often many times per second!

Cache Memory

The speed of CPU is extremely high compared to the access time of main memory. Therefore the performance of CPU decreases due to the slow speed of main memory. To decrease the mismatch in operating speed, a small memory chip is attached between CPU and Main memory whose access time is very close to the processing speed of CPU.

It is called CACHE memory. CACHE memories are accessed much faster than conventional RAM. It is used to store programmes or data currently being executed or temporary data frequently used by the CPU. So each memory makes main memory to be faster and larger than it really is. It is also very expensive to have bigger size of cache memory and its size is normally kept small.

SECONDARY STORAGE DEVICES

Alternatively referred to as external memory and auxiliary storage, secondary storage is a storage medium that holds information until it is deleted or overwritten regardless if the computer has power. For example, a floppy disk drive and hard drive are both good examples of secondary storage devices. As can be seen by the below picture there are three different storage on a computer, although primary storage is accessed much faster than secondary storage because of the price and size limitations secondary storage is used with today's computers to store all your programmes and your personal data.

Definition

Storage Devices are the data storage devices that are used in the computers to store the data. The computer has many

types of data storage devices. Some of them can be classified as the removable data Storage Devices and the others as the non removable data Storage Devices. The data Storage Devices come in many sizes and shapes. And more over the technology used for the storage of the data over them is also altogether different.

The storage devices are one of the most important components of the computer system. The memory is of two types; one is the primary memory and the other one is the secondary memory. The primary memory is the volatile memory and the secondary memory is the non- volatile memory. The volatile memory is the kind of the memory that is erasable and the non- volatile memory is the one where in the contents cannot be erased.

Basically when we talk about the data storage devices it is generally assumed to be the secondary memory. The secondary memory is used to store the data permanently in the computer. The secondary storage devices are usually as follows: hard disk drives - this is the most common type of storage device that is used in almost all the computer systems. The other ones include the floppy disk drives, the CD ROM, and the DVD ROM. The flash memory, the USB data card etc.

The storage devices are used to record the data over any storage surface. The memories may also be of different types depending upon the architecture and the design like the optical data storage memory, magnetic media storage and the mechanical storage media etc and also the flash memory devices etc. The storage devices are actually defined as the peripheral unit which holds the data like the tape, disk, or

flash memory card etc. The most of the drives that are used for the purpose of data storage are fragile and the data can be easily corrupted in them. The data storage devices are the ones that are also used for the backup and the archiving of the data. The data storage devices were at a time in the past used to be too costly and expensive. But these days the data storage devices are becoming cheap day by day. Hence the data storage devices price is falling. So, we are in a position to get a storage device for a comparatively cheaper price than the earlier drive. The technology is improving a lot and now the memory storage capacity has gone up TB.

The data in the storage devices can be in the form of the files, data bases, digital video and the audio etc. The storage devices that are called as the non- volatile can store the data permanently until otherwise erased purposely. This is in the case of the hard disk drives or the floppy disk drives. The other kinds of the storage media like for example the CD and the DVD can even have again two types of the storage; the first one is that in which the data once written cannot be erased.

It is stored permanently over it. While the second type of the CD's or the DVD's are called as the rewritable; where in the data that is once written can be erased completely and the same storage device can be used again for storing the different data. The storage devices are used to record the data over any storage surface.

The memories may also be of different types depending upon the architecture and the design like the optical data storage memory, magnetic media storage and the mechanical storage media etc and also the flash memory devices etc.

The storage devices are actually defined as the peripheral unit which holds the data like the tape, disk, or flash memory card etc. The most of the drives that are used for the purpose of data storage are fragile and the data can be easily corrupted in them. The data storage devices are the ones that are also used for the backup and the archiving of the data. The data storage devices were at a time in the past used to be too costly and expensive.

But these days the data storage devices are becoming cheap day by day. Hence the data storage devices price is falling. So, we are in a position to get a storage device for a comparatively cheaper price than the earlier drive. The technology is improving a lot and now the memory storage capacity has gone up TB. The data in the storage devices can be in the form of the files, databases, digital video and the audio etc. The storage devices that are called as the non-volatile can store the data permanently until otherwise erased purposely. This is in the case of the hard disk drives or the floppy disk drives. The other kinds of the storage media like for example the CD and the DVD can even have again two types of the storage; the first one is that in which the data once written cannot be erased. It is stored permanently over it. While the second type of the CD's or the DVD's are called as the rewritable; where in the data that is once written can be erased completely and the same storage device can be used again for storing the different data.

SECONDARY STORAGE DEVICES

A secondary storage device is another memory device for a computer that will hold and store more information.

Examples of secondary storage devices are floppy disks and external hard drives. These devices will hold information regardless of if the computer itself has power.

Secondary Storage refers to non-volatile data storage which is not directly accessible by the CPU and only accessible via primary storage devices using I/O (Input/Output) channels or device drivers. Typical examples of Secondary Storage are hard disks, floppy disks, optical storage devices (such as CD, DVD drives), magnetic tape data storage devices, RAM drives, and flash memory (such as USB sticks/keys.)

The secondary storage helps in securing the data on media types and storage categories, as data is vulnerable to network attacks, administrative access and media theft. There are many organizations who are working with third parties and disaster recovery efforts. Often data goes offsite and is in the hands of employees that are not authorized to see critical company data. And storage consolidation opens the door to greater administrative access. All of these trends drive the need to ensure the data at rest is secure. Secondary storage devices hold files that are not currently being used.

For a file to be used it must first be copied to main memory first. After any modifications files must be saved to secondary storage. It is advisable to save your data files at the regular intervals as you work on them as data can be lost unexpectedly because of various reasons like interruption in power supply, memory management problems, freezing keyboard, etc. Compared to main memory, secondary memory is slower, bigger, and cheaper (per unit storage.)

Any method for secondary storage must involve two physical parts: a peripheral device (the component of the

computer which 'reads' in or 'writes' out the information to/from the system unit,) and an input/output medium, on which the information is actually stored. Diskettes and cassettes are types of media.

The medium has to be 'in' a corresponding peripheral device for information transfer to take place. In most methods of secondary storage, this transfer is realized by passing the medium by a read/write head, which is capable of sensing/writing information from/to that type of medium. Secondary storage media may be removable (easily separable from the computer) like your diskettes, or fixed, like hard disks, which, for example, can be found in relatively expensive PCs.

A fixed medium generally has much greater capacity (and is faster) compared to a removable one of the same type. On the other hand, by replacing your removable medium by another one when it is full, you can attain a virtually unlimited storage capacity. Another nice feature of removable media is that they allow backup copies of important material to be taken and kept away from the computer, so that they do not get corrupted if some disaster strikes the computer. The act of retrieving pieces of stored information is called access.

There may be two kinds of access to data stored as a sequence of items: Sequential access, where the items are traversed one by one from the beginning of the sequence to the desired one, and direct access (or random access), where any item can be accessed relatively independently of its location in the sequence. Different types of secondary memory media support one or more of these different methods, depending on their nature. Cassette tapes, for example, are sequential access media: You have to 'go through' the first

10 minutes if you want to listen to a piece which begins at the 11th minute of the tape. The most common type of secondary storage medium is the magnetic disk.

Diskettes (small disks made of flexible plastic) and hard disks (made of rigid metal) are the two different kinds of magnetic disks. Magnetic disks are coated with a magnetizable substance, the spots on which mean '0' or '1', depending on whether they are magnetized or not. Each surface of the disk is subdivided into concentric rings called tracks. Disks with bigger capacity have more tracks than others. In larger computers, one stores the same amount of data in each track and keeps several disks mounted on a shaft on top of each other as a disk pack. The group of tracks which are at the same position in their respective disks on a disk pack is called the disk cylinder for that position. The associated peripheral device is called the disk unit. At least one read-write head is assigned for each surface.

The read-write heads are mounted on a device called the access mechanism, which positions them on the cylinder in which the appropriate data item is to be located. The time for the read-write head to move from its initial location to the appropriate track is called seek time. In larger computers, the shaft on which the disk is mounted continuously rotates at a great speed, and the time required for the relevant portion of the track to come near the head is called rotational delay.

The time needed for the actual information transfer between head and disk is called data movement time. Disk access time is the sum of these three components. In larger computers, the read/write head never touches the disk. Diskettes, on the other hand, rotate only when a read/write

is taking place. Furthermore, the head actually touches the diskette. Diskettes usually come in one of two diameters: 3.5 inches or 5.25 inches. 3.5-inch diskettes are products of a newer technology and can store more than 5.25-inch diskettes. Different types of computers employ different methods of data organization on a diskette. For this reason, a new diskette has to be formatted before it is used on a particular computer. Formatting a disk involves designating parts of each track as individually addressable (directly accessible) sectors. Each diskette contains a file directory, in which the names, lengths and starting addresses of the files (collections of related data) stored on it are listed.

People with many files to store in the same disk are encouraged to use a tree-like directory structure in which related files are kept in separate subdirectories. If a computer system has enough main memory and a programme designed to read/write data from/to a disk is to be executed, one can transfer the contents of the disk to main memory and then use a method called disk emulation (or RAM disk) to 'trick' the computer into thinking that it is accessing the disk, while it actually deals with the much faster main memory. An older type of secondary storage medium is magnetic tape.

Characters of data are represented in byte form across the tracks which go along the length of the tape. Other types of secondary storage media include mass storage units, which are combinations of spools of magnetic tape that can be automatically retrieved from a huge 'library' and read to disk. These store enormous amounts of information and do not require the intervention of a human operator, note that some

disk packs and tapes need this kind of intervention. Optical disks employ a relatively new technology.

They use laser beams for the read/write operations. Most optical disks are of the CD-ROM (Compact Disk/Read-Only Memory) type; whose vendor-determined contents cannot be changed by the user. Because of their great capacity, these are used for storing multimedia (*i.e.* mixed text, graphics, sound, etc.) applications of great size. Their speed is low, compared to magnetic hard disks. As secondary storage media can be damaged and files on them become corrupted, it is suggested to make back-up copies of valuable files on a regular basis.

Lots of people skip the last but very important step in the backup procedure - check that the backup copy of files is not damaged.

Types of Secondary Storage Devices

It is important to know the difference between secondary storage and a computer's main memory. Secondary storage is also called auxiliary storage and is used to store data and programmes when they are not being processed. Secondary storage is more permanent than main memory, as data and programmes are retained when the power is turned off. The needs of secondary storage can vary greatly between users. A personal computer might only require 20,000 bytes of secondary storage but large companies, such as banks, may require secondary storage devices that can store billions of characters. Because of such a variety of needs, a variety of storage devices are available. The two most common types of secondary storage are magnetic tapes and magnetic disks.

Floppy Disk

They are plastic square disks, usually with a silver or black sliding piece going across the top. These disks come in a variety of colours and they hold about 144 million bytes. (Bytes are characters, symbols and letters).

Magnetic Tape Storage

Magnetic tape data storage uses digital recording on magnetic tape to store digital information. Modern magnetic tape is most commonly packaged in cartridges and cassettes. The device that performs actual writing or reading of data is a tape drive. Autoloaders and tape libraries automate cartridge handling.

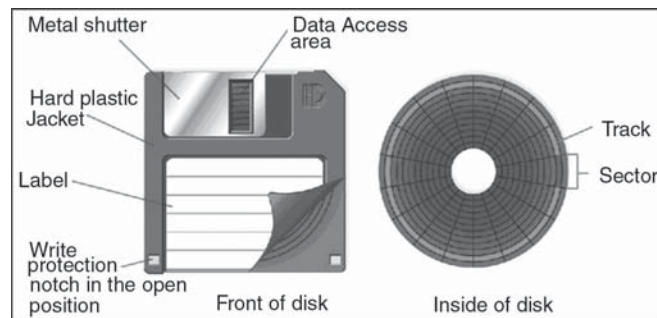
Magnetic tape is a one-half inch or one-quarter inch ribbon of plastic material on which data is recorded. The tape drive is an input/output device that reads, writes and erases data on tapes. Magnetic tapes are erasable, reusable and durable. They are made to store large quantities of data inexpensively and therefore are often used for backup. Magnetic tape is not suitable for data files that are revised or updated often because it stores data sequentially.

Diskettes

A diskette is a random access, removable data storage medium that can be used with personal computers. The term usually refers to the magnetic medium housed in a rigid plastic cartridge measuring 3.5 inches square and about 2 millimeters thick. Also called a "3.5-inch diskette," it can store up to 1.44 megabytes (MB) of data.

Although many personal computers today come with a 3.5-inch diskette drive pre-installed, some notebook computers and centrally-administered desktop computers

omit them. Made of flexible Mylar, a diskette can record data as magnetized spots on tracks on its surface. Diskettes became popular along with the personal computer. The older diskette, 5-1/4 inches in diameter, is still in use, but newer computers use the 3-1/2 inch diskette.



The 3-1/2 inch diskette has the protection of a hard plastic jacket, a size to fit conveniently in a shirt pocket or purse, and the capacity to hold significantly more data than a 5-1/4 inch diskette. Diskettes offer particular advantages, which, as you will see, are not readily available with hard disk:

Portability.

Diskettes easily transport data from one computer to another. Workers, for example, carry their files from office computer to home computer and back on a diskette instead of in a briefcase. Students use the campus computers but keep their files on their own diskettes.

Backup

It is convenient to place an extra copy of a hard disk file on a diskette.

New Software

Although, for convenience, software packages are kept on hard disk, new software out of the box may come on diskettes

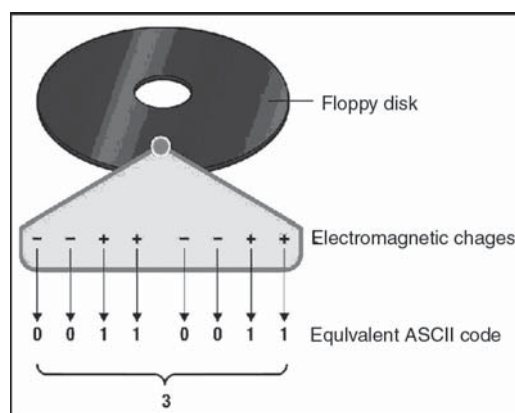
(new software also may come on CD-ROM disks, which we will discuss shortly). The end of the diskettes useful life-time may be upon us. In 1998 Macintosh introduced its new computer, the IMAC, without a floppy disk drive.

Alternatives such as Zip disks (discussed later), or transferring data via networks are making the low-capacity diskette become obsolete. The diskette was introduced in the early 1970s by IBM as a new type of secondary storage. Originally they were eight inches in diameter and were thin and flexible which gave them the name floppy disks, or floppies. Diskettes are used as the principle medium of secondary storage for personal computers. They are available in two different sizes: 3 1/2 inch and 5 1/4 inch.

Storage Capacity

Before you can store data on your diskette, it must be formatted (G).

The amount of data you can store on a diskette depends on the recording density and the number of tracks on the diskette. The recording density is the number of bits (G) that can be recorded on one inch of track on the diskette, or bits per inch (bpi).



The second factor that influences the amount of data stored on a diskette is the number of tracks on which the data can be stored or tracks per inch (tpi). Commonly used diskettes are referred to as either double-density or high-density (single-density diskettes are no longer used). Double-density diskettes (DD) can store 360K for a 5 1/4 inch diskette and 720K for a 3 1/2 inch diskette. High-density diskettes (HD) can store 1.2 megabytes (G) on a 5 1/4 inch diskette and 1.44 megabytes on a 3 1/2 inch diskette.

Care Of Diskettes

You should keep diskettes away from heat, cold, magnetic fields (including telephones) and contaminated environments such as dust, smoke, or salt air. Also keep them away from food and do not touch the disk surface.

Magnetic Disk Storage

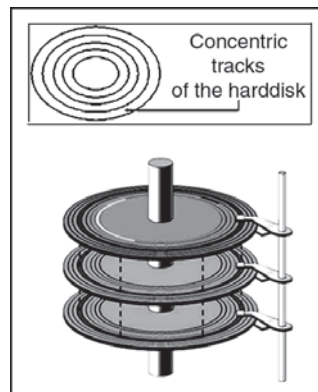
The most common physical device for storing files is the magnetic disk. Actually, a disk typically contains several rotating disks, or platters. The surfaces of the platters are covered in metal oxide, and read/written by electromagnetic recording heads, rather like those on an audio cassette recorder.

There is one head for each surface, and all the heads move together. The disk rotates at around 3600 rpm (or approx 90mph), with the heads floating microscopic distances above the surfaces. Modern disks for workstations typically hold 500MB - 9GB, and cost of the order of £200 - 3000; prices are currently dropping rapidly.

Magnetic disks are the most widely used storage medium for computers. A magnetic disk offers high storage capacity, reliability, and the capacity to directly access stored data.

Magnetic disks hold more data in a small place and attain faster data access speeds. Types of magnetic disks include diskettes, hard disks, and removable disk cartridges.

Diskettes and hard disks are magnetic media; that is, they are based on a technology of representing data as magnetized spots on the disk with a magnetized spot representing a 1 bit and the absence of such a spot representing a 0 bit.



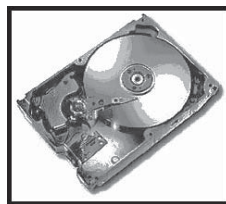
Reading data from the disk means converting the magnetized data to electrical impulses that can be sent to the processor. Writing data to disk is the opposite: sending electrical impulses from the processor to be converted to magnetized spots on the disk. The surface of each disk has concentric tracks on it. The number of tracks per surface varies with the particular type of disk.

Hard Disks

A hard disk drive (HDD) is a data storage device used for storing and retrieving digital information using rapidly rotating disks (platters) coated with magnetic material. An HDD retains its data even when powered off. Data is read in a random-access manner, meaning individual blocks of data can be stored or retrieved in any order rather than sequentially. An HDD consists of one or more rigid ("hard")

rapidly rotating disks (platters) with magnetic heads arranged on a moving actuator arm to read and write data to the surfaces.

A hard disk is a metal platter coated with magnetic oxide that can be magnetized to represent data. Hard disks come in a variety of sizes.



Hard disk for mainframes and minicomputers may be as large as 14 inches in diameter. Several disks can be assembled into a disk pack. There are different types of disk packs, with the number of platters varying by model.

Each disk in the pack has top and bottom surfaces on which to record data. Many disk devices, however, do not record data on the top of the top platter or on the bottom of the bottom platter.

A disk drive is a machine that allows data to be read from a disk or written on a disk. A disk pack is mounted on a disk drive that is a separate unit connected to the computer. Large computers have dozens or even hundreds of disk drives. In a disk pack all disks rotate at the same time although only one disk is being read or written on at any one time.

The mechanism for reading or writing data on a disk is an access arm; it moves a read/write head into position over a particular track. The read/write head on the end of the access arm hovers just above the track but does not actually touch the surface. When a read/write head does accidentally touch the disk surface, this is called a head crash and all data is

destroyed. Data can also be destroyed if a read/write head encounters even minuscule foreign matter on the disk surface. A disk pack has a series of access arms that slip in between the disks in the pack.

Two read/write heads are on each arm, one facing up for the surface above it and one facing down for the surface below it. However, only one read/write head can operate at any one time. In some disk drives the access arms can be retracted; then the disk pack can be removed from the drive. Most disk packs, however, combine the disks, access arms, and read/write heads in a sealed module called a Winchester disk. Winchester disk assemblies are put together in clean rooms so even microscopic dust particles do not get on the disk surface.

Hard disks for personal computers are 5-1/4 inch or 3-1/2 inch disks in sealed modules and even gigabytes are not unusual. Hard disk capacity for personal computers has soared in recent years; capacities of hundreds of megabytes are common and gigabytes are not unusual. Although an individual probably cannot imagine generating enough output-letters, budgets, reports, and so forth-to fill a hard disk, software packages take up a lot of space and can make a dent rather quickly.

Furthermore, graphics images and audio and video files require large file capacities. Perhaps more important than capacity, however, is the convenience of speed. Personal computer users find accessing files on a hard disk is significantly faster and thus more convenient than accessing files on a diskette. Hard disks provide larger and faster secondary storage capabilities than diskettes. Usually hard

disks are permanently mounted inside the computer and are not removable like diskettes.

On minicomputers (G) and mainframes (G), hard disks are often called fixed disks. They are also called direct-access storage devices (DASD). Most personal computers have two to four disk drives. The input/output device that transfers data to and from the hard disk is the hard disk drive.

Hard Disks in Groups

A concept of using several small disks that work together as a unit is called a redundant array of inexpensive disks, or simply RAID.

The group of connected disks operates as if it were just one large disk, but it speeds up reading and writing by having multiple access paths. The data file for, say, aircraft factory tools, may be spread across several disks; thus, if the computer is used to look up tools for several workers, the computer need not read the data in turn but instead read them at the same time in parallel.

Furthermore, data security is improved because if a disk fails, the disk system can reconstruct data on an extra disk; thus, computer operations can continue uninterrupted. This is significant data insurance.

How Data Is Organized on a Disk

There is more than one way of physically organizing data on a disk. The methods we will consider here are the sector method and the cylinder method.

The Sector Method

In the sector method each track is divided into sectors that hold a specific number of characters. Data on the track

is accessed by referring to the surface number, track number, and sector number where the data is stored. The sector method is used for diskettes as well as disk packs.

Zone Recording

The fact that a disk is circular presents a problem: The distances around the tracks on the outside of the disk are greater than that of the tracks on the inside.

A given amount of data that takes up 1 inch of a track on the inside of a disk might be spread over several inches on a track near the outside of a disk. This means that the tracks on the outside are not storing data as efficiently. Zone recording involves dividing a disk into zones to take advantage of the storage available on all tracks, by assigning more sectors to tracks in outer zones than to those in inner zones. Since each sector on the disk holds the same amount of data, more sectors mean more data storage than if all tracks had the same number of sectors.

The Cylinder Method

A way to organize data on a disk pack is the cylinder method. The organization in this case is vertical. The purpose is to reduce the time it takes to move the access arms of a disk pack into position. Once the access arms are in position, they are in the same vertical position on all disk surfaces. To appreciate this, suppose you had an empty disk pack on which you wished to record data.

You might be tempted to record the data horizontally-to start with the first surface, fill track 000, then fill track 001, track 002, and so on, and then move to the second surface and again fill tracks 000, 001, 002, and so forth.

Each new track and new surface, however, would require movement of the access arms, a relatively slow mechanical process. Recording the data vertically, on the other hand, substantially reduces access arm movement. The data is recorded on the tracks that can be accessed by one positioning of the access arms—that is, on one cylinder.

To visualize cylinder organization, pretend a cylindrically shaped item, such as a tin can, were figuratively dropped straight down through all the disks in the disk pack.

All the tracks thus encountered, in the same position on each disk surface, comprise a cylinder. The cylinder method, then, means all tracks of a certain cylinder on a disk pack are lined up one beneath the other, and all the vertical tracks of one cylinder are accessible by the read/write heads with one positioning of the access arms mechanism.

Tracks within a cylinder are numbered according to this vertical perspective: A 20-surface disk pack contains cylinder tracks numbered 0 through 19, top to bottom.

Hard Disk Storage Capacity

Like diskettes, hard disks must be formatted before they can store information. The storage capacity for hard drives is measured in megabytes. Common sizes for personal computers range from 100MB to 500MB of storage. Each 10MB of storage is equivalent to approximately 5,000 printed pages (with approximately 2,000 characters per page).

Removable Storage: Zip Disks

The Zip drive is a medium-capacity removable disk storage system that was introduced by Iomega in late 1994. Originally, Zip disks launched with capacities of 100 MB,

but later versions increased this to first 250 MB and then 750 MB.

Personal computer users, who never seem to have enough hard disk storage space, may turn to a removable hard disk cartridge. Once full, a removable hard disk cartridge can be replaced with a fresh one.

In effect, a removable cartridge is as portable as a diskette, but the disk cartridge holds much more data. Removable units also are important to businesses concerned with security, because the units can be used during business hours but hidden away during off hours. A disadvantage of a removable hard disk is that it takes longer to access data than a built-in hard drive.

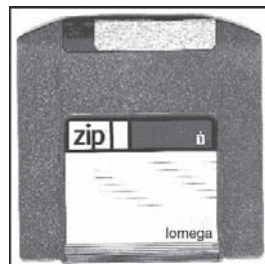


Fig. Zip Disk

The most popular removable disk media is the Zip drive as shown in Figure. Over 100's of millions have been sold, making it the de facto standard. The disk cartridges look like a floppy disk, but are slightly bigger in all dimensions. Older Zip disks hold 100MB, newer ones hold 250MB and cost \$8-\$10 a piece (Floppies hold 1.4MB and cost around \$2). The drive sells for around \$80- \$125. Many new PCs come with Zip drives built in addition to floppy drives. Zip disks are a great way to store large files and software programmes.

Backup

Backup means creating a copy of important programmes and data. To backup diskettes, copy the data from one to the other. Diskettes are frequently used to backup important files on hard drives of personal computers.

Disk Cartridges

Removable disk cartridges are another form of disk storage for personal computers. They offer the storage and fast access of hard disks and the portability of diskettes. They are often used when security is an issue since, when you are done using the computer, the disk cartridge can be removed and locked up leaving no data on the computer.

Data Recovery

Data Recovery is the process of salvaging data from damaged, failed, corrupted, or inaccessible secondary storage media when it cannot be accessed normally. Often the data are being salvaged from storage media such as internal or external hard disk drives, solid-state drives (SSD), USB flash drive, storage tapes, CDs, DVDs, RAID, and other electronics. Recovery may be required due to physical damage to the storage device or logical damage to the file system that prevents it from being mounted by the host operating system.

There does exist certain software programmes that, under special circumstances, can recover data and programmes that have been "lost." Often your operating system (G) will be able to do this, or applications such as Norton Utilities.

Optical Disk Storage

In computing and optical disc recording technologies, an optical disc (OD) is a flat, usually circular disc which encodes

binary data (bits) in the form of pits (binary value of 0 or off, due to lack of reflection when read) and lands (binary value of 1 or on, due to a reflection when read) on a special material (often aluminium) on one of its flat surfaces.

The explosive growth in storage needs has driven the computer industry to provide cheaper, more compact, and more versatile storage devices with greater capacity. This demanding shopping list is a description of the optical disk, like a CD. The technology works like this: A laser hits a layer of metallic material spread over the surface of a disk.

When data is being entered, heat from the laser produces tiny spots on the disk surface. To read the data, the laser scans the disk, and a lens picks up different light reflections from the various spots. Optical storage technology is categorized according to its read/write capability. Read-only media are recorded on by the manufacturer and can be read from but not written to by the user. Such a disk cannot, obviously, be used for your files, but manufacturers can use it to supply software. Applications software packages sometimes include a dozen diskettes or more; all these could fit on one optical disk with plenty of room to spare. The most prominent optical technology is the CD-ROM, for compact disk read-only memory. CD-ROM has a major advantage over other optical disk designs: The disk format is identical to that of audio compact disks, so the same dust-free manufacturing plants that are now stamping out digital versions of Mozart or Mary Chapin Carpenter can easily convert to producing anything from software to an encyclopaedia. Furthermore, CD-ROM storage is large -up to 660 megabytes per disk, the equivalent of over 400 3-1/2 inch diskettes.



When buying a computer the speed of the CD-ROM drive is advertised using an "X" factor, like 12X, or 24X. This indicates the speed at which the CD can transfer data to the CPU - the higher the X factor, the faster the CD. Modern computers now offer a write CD drive or, CD-RW as an option. CD-RW is a write-once, read-many media. With a CD-RW drive, you can create your own CDs. This offers an inexpensive, convenient, safe way to store large volumes of data such as favourite songs, photographs, etc.

DVDs

DVD (sometimes explained as "digital video disc" or "digital versatile disc") is a digital optical disc storage format, invented and developed by Philips, Sony, Toshiba, and Panasonic in 1995. DVDs offer higher storage capacity than compact discs while having the same dimensions.

Digital Versatile Disk (DVD) drives are now widely available in computers as well as home entertainment centres. DVD-ROM drives can read data, such as stored commercial videos for playing. DVD-RW allow DVDs to be created on a computer.

The DVD is a flat disk, the size of a CD - 4.7 inches diameter and .05 inches thick. Data are stored in a small indentation in a spiral track, just like in the CD. DVD disks are read by a laser beam of shorter wave-length than used by the CD ROM drives. This allows for smaller indentations and

increased storage capacity. The data layer is only half as thick as in the CD-ROM. This opens the possibility to write data in two layers. The outer gold layer is semi transparent, to allow reading of the underlying silver layer. The laser beam is set to two different intensities, strongest for reading the underlying silver layer. A 4.7 GB side of a DVD can hold 135 minutes top quality video with 6 track stereo. This requires a transmission rate of 4692 bits per second. The 17 GB disk holds 200 hours top quality music recording.



Fig. DVD Disk and Drive

DVD movies are made in two "codes." Region one is USA and Canada, while Europe and Asia is region two. When you play movies, your hardware (MPEG decoder. MPEG is the data coding for movies similar to JPEG for pictures.) must match the DVD region. The movies are made in two formats, each with their own coding. The DVD drives come in 2X, 4X, etc. versions, like the CD-ROM's.

The DVD drives will not replace the magnetic hard disks. The hard disks are being improved as rapidly as DVD, and they definitely offer the fastest seek time and transmission rate (currently 5-10 MB/second). No optic media can keep up with this. But the DVD will undoubtedly gain a place as the successor to the CD ROM and is playing an important role in the blending of computers and entertainment centres.

Magnetic Tape Storage

We saved magnetic tape storage for last because it has taken a subordinate role in storage technology. Magnetic tape looks like the tape used in music cassettes plastic tape with a magnetic coating. As in other magnetic media, data is stored as extremely small magnetic spots. Tapes come in a number of forms, including 1/2-inch-wide tape wound on a reel, 1/4-inch-wide tape in data cartridges and cassettes, and tapes that look like ordinary music cassettes but are designed to store data instead of music. The amount of data on a tape is expressed in terms of density, which is the number of characters per inch (cpi) or bytes per inch (bpi) that can be stored on the tape.

The highest-capacity tape is the digital audio tape, or DAT, which uses a different method of recording data. Using a method called helical scan recording, DAT wraps around a rotating read/write head that spins vertically as it moves. This places the data in diagonal bands that run across the tape rather than down its length.

This method produces high density and faster access to data. Two reels are used, a supply reel and a take-up reel. The supply reel, which has the tape with data on it or on which data will be recorded, is the reel that is changed. The take-up reel always stays with the magnetic tape unit. Many cartridges and cassettes have the supply and take-up reels built into the same case. Tape now has a limited role because disk has proved the superior storage medium.

Disk data is quite reliable, especially within a sealed module. Furthermore, as we will see, disk data can be accessed directly, as opposed to data on tape, which can be

accessed only by passing by all the data ahead of it on the tape. Consequently, the primary role of tape today is as an inexpensive backup medium.

Backup Systems

Although a hard disk is an extremely reliable device, a hard disk drive is subject to electromechanical failures that cause loss of data. Furthermore, data files, particularly those accessed by several users, are subject to errors introduced by users.

There is also the possibility of errors introduced by software. With any method of data storage, a backup system a way of storing data in more than one place to protect it from damage and errors is vital. As we have already noted, magnetic tape is used primarily for backup purposes.

For personal computer users, an easy and inexpensive way to back up a hard disk file is to simply copy it to a diskette whenever it is updated. But this is not practical for a system with many files or many users. Personal computer users have the option of purchasing their own tape backup system, to be used on a regular basis for copying all data from hard disk to a high-capacity tape. Data thus saved can be restored to the hard disk later if needed. A key advantage of a tape backup system is that it can copy the entire hard disk in minutes, saving you the trouble of swapping diskettes in and out of the machine.

A rule of thumb among computer professionals is to estimate disk needs generously and then double that amount. But estimating future needs is rarely easy. Many users, therefore, make later adjustments like adding a removable hard disk cartridge to accommodate expanding storage needs.

Benefits of Secondary Storage Devices

Secondary storage, sometimes called auxiliary storage, is storage separate from the computer itself, where you can store software and data on a semi permanent basis. Secondary storage is necessary because memory, or primary storage, can be used only temporarily.

If you are sharing your computer, you must yield memory to someone else after your programme runs; if you are not sharing your computer, your programmes and data will disappear from memory when you turn off the computer. However, you probably want to store the data you have used or the information you have derived from processing; that is why secondary storage is needed.

Furthermore, memory is limited in size, whereas secondary storage media can store as much data as necessary.

- *Capacity:* Organizations may store the equivalent of a roomful of data on sets of disks that take up less space than a breadbox. A simple diskette for a personal computer holds the equivalent of 500 printed pages, or one book. An optical disk can hold the equivalent of approximately 400 books.
- *Reliability:* Data in secondary storage is basically safe, since secondary storage is physically reliable. Also, it is more difficult for unscrupulous people to tamper with data on disk than data stored on paper in a file cabinet. Convenience. With the help of a computer, authorized people can locate and access data quickly. Cost. Together the three previous benefits indicate significant savings in storage costs. It is less expensive to store data on tape or disk (the principal means of

Computer Memory System

secondary storage) than to buy and house filing cabinets. Data that is reliable and safe is less expensive to maintain than data subject to errors. But the greatest savings can be found in the speed and convenience of filing and retrieving data.

These benefits apply to all the various secondary storage devices but, as you will see, some devices are better than others. We begin with a look at the various storage media, including those used for personal computers, and then consider what it takes to get data organized and processed.