

COMP20081

Systems Software

Revision Guide



Contents



1	Introduction	1
1.1	History of performing computations	1
1.2	Hardware	3
1.3	What is an operating system (OS)?	4
1.4	Functions of an operating system	6
1.5	Current operating system (OS) trends	8
1.6	Operating system (OS) layers	10
1.7	Kernel mode	10
2	Process Management	14
2.1	Programs and processes	14
2.2	Process life cycle	15
2.3	Program execution	21
2.4	Concurrency	23
2.5	Process scheduling	24
3	Threads and Concurrency	33
3.1	What are threads?	33
3.2	Sequential and concurrent programming	36
3.3	Sequential execution	38
3.4	Introduction to concurrent execution	40
3.5	Interleaving	41
3.6	User and kernel threads	45
3.7	Multi-threading models	47

3.8	Evaluation of concurrent programming	50
4	Synchronisation and Mutual Exclusion	52
4.1	Race condition	52
4.2	Inter-process synchronisation	54
4.3	The producer-consumer problem	57
4.4	The dining philosophers problem	64
5	Memory Management	70
5.1	Memory hierarchy	70
5.2	Primary memory	71
5.3	Secondary memory	74
5.4	Read-only memory (ROM)	80
5.5	Concept and aim of memory management	81
5.6	Overview of memory abstraction	82
5.7	Memory abstraction - physical address	83
5.8	Memory abstraction - logical address	85
5.9	Memory Management Unit (MMU)	86
5.10	Memory partitioning	88
5.11	Swapping	91
5.12	Managing free space	92
6	Virtual memory	99
6.1	Introduction	99
6.2	Paging	99
6.3	Segmentation	109
6.4	Comparison of paging and segmentation	112
6.5	Paging segments	113
7	File Management	114
7.1	Non-volatile storage	114
7.2	Filing system	115
7.3	Files	115
7.4	Directories	123
8	File System Implementation	129
8.1	File system layout	129
8.2	File block allocation methods	130

8.3	Implementation of directories/folders	137
8.4	Locating files using i-nodes	139
8.5	Determining block size	141
8.6	Tracking free space	142
8.7	Consistency	144
8.8	Journaling	145
8.9	Backups	146
8.10	Drives and mount points	146
8.11	RAID	147
9	Input/Output (I/O)	150
9.1	Input/Output (I/O) devices	150
9.2	Input/Output (I/O) system	152

1. Introduction

1.1	History of performing computations	1
1.2	Hardware	3
1.3	What is an operating system (OS)?	4
1.4	Functions of an operating system	6
1.5	Current operating system (OS) trends	8
1.6	Operating system (OS) layers	10
1.7	Kernel mode	10

1.1 History of performing computations

1950s

A user may want to perform computations that may have a complexity that is not feasible by a standard calculator.

A description of the computations to be performed could be expressed in a physical form using a perforated card.



Figure 1.1: Perforated card

There are multiple lines on the card, each with varying holes. Each line describes a particular instruction that is necessary for the computations to be completed.



Figure 1.2: Process for performing computations in the 1950s

>1960s

The process of performing computations after 1960 changed dramatically and required far less human input.

Automation was brought about by the “Operating System Paradigm” in which users could interface with a computer system directly through an operating system (OS).

User → Operating System (OS) → Results
--

Figure 1.3: Process for performing computations after 1960

History details

Late 1950s	<ul style="list-style-type: none"> Standard subroutines were produced that were loaded at start-up. These contained features similar to those found on an operating system. Magnetic tapes were used for storage and were later replaced by disks. Assemblers started to be used. These are programs that takes basic computer instructions and convert them in to machine code; this is a pattern of binary bits (0's and 1's) that the computer system's processor can use to perform its basic operations. High-level languages, which consisted of more natural and human-readable language, started to be used. For example, FORTRAN is a general-purpose, compiled imperative programming language that is especially suited to numeric computation and scientific computing and was introduced in 1957.
1960s	<p>Automated batch system.</p> <ul style="list-style-type: none"> This replaced the computer operator. Several programs could be loaded in to memory and automatically processed in a “first in, first out” (FIFO) fashion.
1970s	<p>Multiprogramming.</p> <ul style="list-style-type: none"> The computer could switch between jobs, which allows processing and input/output (I/O) interaction simultaneously.

1980s	<p>Graphical user interfaces (GUIs).</p> <ul style="list-style-type: none"> The interaction between a computer system and a user through the medium of a mouse and keyboard.
-------	---

1.2 Hardware

External hardware

A **peripheral** is any external hardware device that provides input/output (I/O) for the computer.

For example, a keyboard and mouse are input peripherals, while a monitor and printer are output peripherals. Some peripherals, such as external hard drives, provide both input and output for the computer.

A computer system generally has many internal hardware components and hardware peripherals.



Internal hardware

A **processor** or **central processing unit (CPU)** is the hardware within a computer that carries out the instructions of a computer program by performing the basic arithmetic, logical, and input/output operations of the system.

A **motherboard** is the main printed circuit board (PCB) in a computer. The motherboard is a computer's central communications backbone connectivity point, through which all components and external peripherals connect.

Inside of a computer system, there are many components connected to the processor via the motherboard.

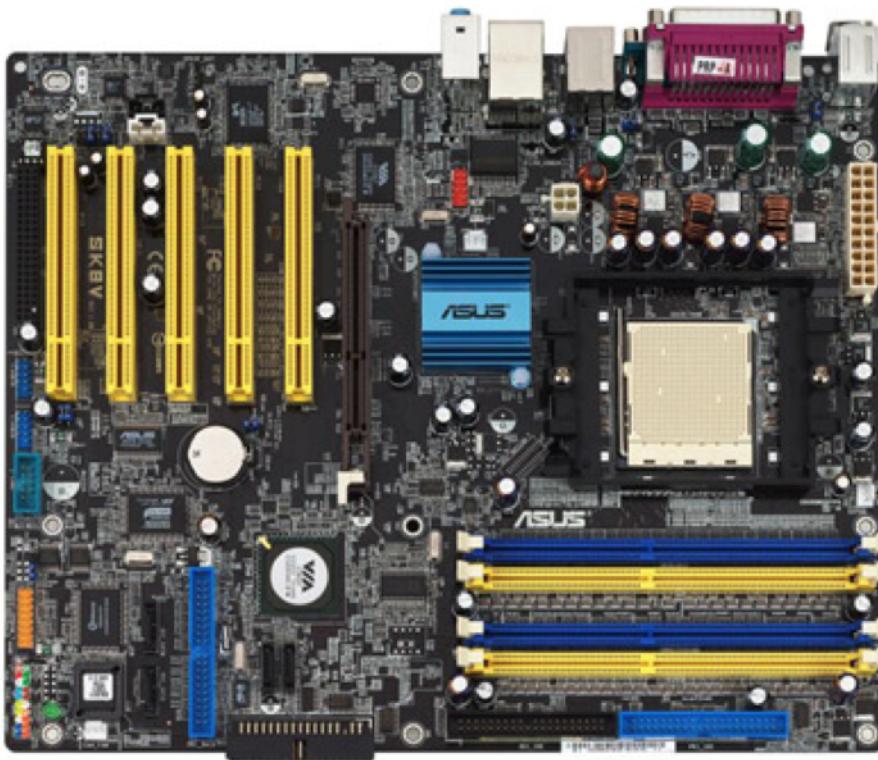


Figure 1.5: Typical computer motherboard

The motherboard connects:

- all of the internal components via the data bus; and
- the peripherals.

1.3 What is an operating system (OS)?

Definition

An **operating system (OS)** is a collection of system programs that manage the hardware resources and peripherals connected to a computer system. It is also responsible for the graphical user interface or command line interface and all other software running on the computer system.

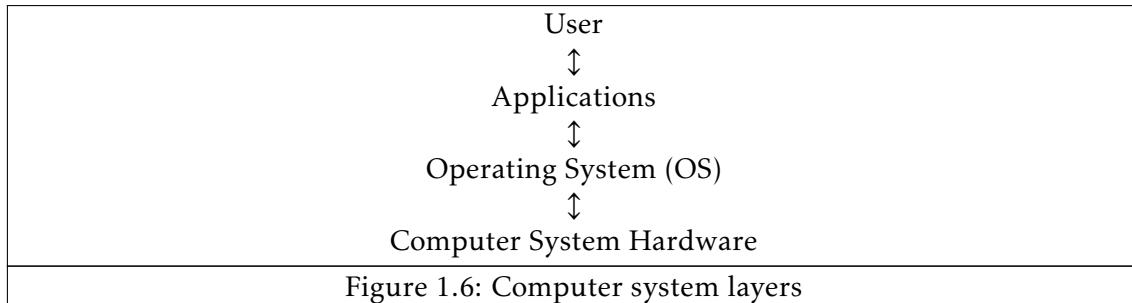
Purpose of an operating system (OS)

An operating system (OS) is designed to:

- eliminate the need to have hardware knowledge to operate a computer system;

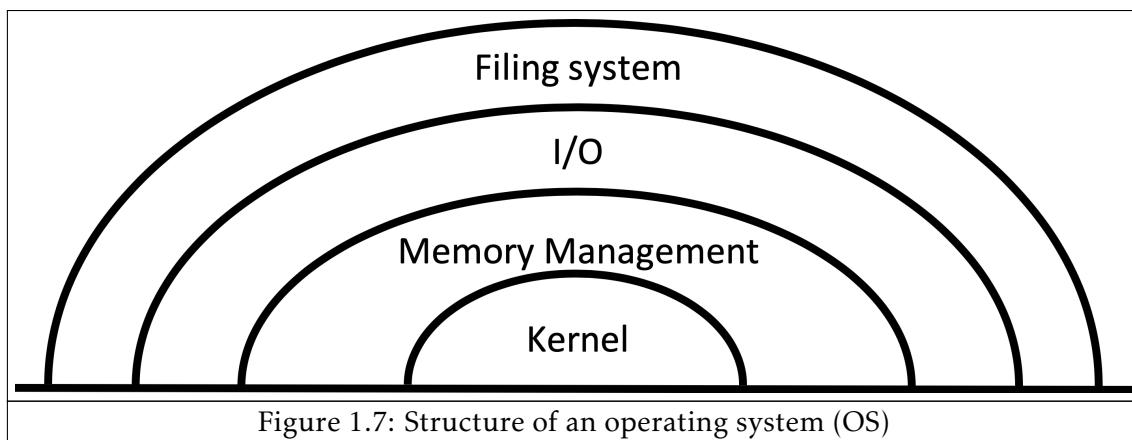
- make the boundary between hardware and software transparent, allowing the user to not be concerned with the technical details; and
- provide a user-friendly environment to execute and develop programs.

These attributes are achieved by layering the computer system such that the user can interface with applications, rather than the operating system (OS) or the hardware directly.



Structure of an operating system (OS)

The structure of an operating system (OS) can be said to resemble an onion.



An operating system (OS) has four main components.

The **kernel** hides the complexity of how a computer system works from users. It is responsible for:

- process management;
- CPU scheduling; and
- handling interrupts.

Memory management is responsible for allocating and deallocating memory to processes.

Input/output (I/O) includes any interaction between the internal computer system components and peripherals.

The **filing system** is comprised of file management subsystems.

Each layer in the operating system (OS) structure provides functions to the above layers. Each layer uses facilities provided by layers within and below that layer.

Practical features of current operating systems (OSs)

Concurrency	Allows overlapping input/output (I/O) operations with computations and several programs to be stored in memory at a single time.
Sharing of resources	Sharing hardware and peripherals, such as hard disks and printers.
Access to long term storage	Important for saving important files on mediums such as hard disk drives (HDDs) and solid-state drives (SSDs).
Non-determinacy	The ability to cope with unpredictable events without crashing.

1.4 Functions of an operating system

An operating system (OS) has two main complementary functions:

- resource managing; and
- machine extending.

Resource managing

It manages resources shared among users and user programs and maximises their utilisation of the CPU, RAM and other resources. This is done simultaneously in order to increase the availability.

This is similar to the role of computer operators in the 1950s.

Machine extending

It presents a virtual machine (or extended machine) to users that is much easier to access than the underlying physical machine.



The virtual machine presented to the user provides an abstraction of the computer system. This hides the complexity of the hardware from the user; this means that the user need only be concerned with the details of the hardware if they desire.

This is a way of translating the functions needed by a user from the hardware to a presentable and user-friendly medium. As a result, the operating system (OS) acts as an intermediary layer between the user and machine language.

The benefit of this abstraction can be demonstrated when comparing how computations may be processed with and without an operating system (OS).

Without Operating System (OS)	With Operating System (OS)
<p>The instructions written in machine code or assembly language much interface directly with memory hardware. As such, the memory locations to load the two numbers from must be explicitly defined and the memory location to which the result is stored must also be defined.</p> <p>The example below shows a possible assembly code implementation of a computation that is capable of adding <u>two numbers</u>.</p>	<p>The instructions can be written in a high-level language, such as C++.</p>
<p>LDAA \$80 (load number at memory location 80)</p>	
<p>LDAB \$81 (load number at memory location 81)</p>	
<p>ADDB (add these two numbers)</p>	<pre>int a, b, c; a = 1; b = 2; c = a + b;</pre>
<p>STAA &55 (store the sum to memory location 55)</p>	

Figure 1.9: Adding two numbers

This demonstrates that program development is much more user-friendly with an operating system (OS). This is because without an operating system (OS) the user must have knowledge of the system hardware; in this case, the necessary memory locations.

In addition, it is possible that the machine code or assembly language written may not work on another computer system as that computer system may have a different architecture or the memory locations may be different. For example, in another computer system:

- memory location 81 may not exist as the memory is smaller; or
- memory location 55 contains important data or instructions that should not be overwritten and therefore the computer system may crash.

By contrast, with an operating system (OS), it is possible to perform computations without interfacing directly with hardware. In the example above, variables (a, b and c) can

be used to access data in memory rather than addressing memory locations. The only concern here is that the variable *c* is able to store the result of *a + b*.

The operating system (OS) provides a unified environment to users to run their computations in different systems. The operating system (OS) is capable of taking high-level code and translating it into the machine code that can be executed on a particular computer system.

1.5 Current operating system (OS) trends

Hardware evolution

Due to fast rate of hardware evolution, operating systems (OSs) are more wide-spread than just traditional desktop computers. They can be found on hardware such as:

- mobile devices, such as smartphones and tablets; and
- embedded systems.

Desktop	Mobile	Embedded
<ul style="list-style-type: none"> • Windows • macOS • Linux 	<ul style="list-style-type: none"> • iOS • Android • Symbian OS 	<ul style="list-style-type: none"> • Windows Embedded

Figure 1.10: Popular operating systems (OSs)

Multiprocessor systems

Definition

A **multiple processor computer system** makes use of two or more processors and has the ability to allocate tasks between them.

Workstations

A single machine may contain multiple processors and therefore have large computing power.

Distributed and network systems

These computer systems share computing power and peripherals.



The diagram shows an abstraction of a distributed or network system. P1 and P2 are the connected computer systems that both share memory, access to the disk control and, if the system is a network system, the network interface.

However, there is a distinguishable difference between distributed and network systems.

Similarities	Differences
Consist of multiple systems that are interconnected to exchange information.	In distributed systems, users are not aware of the multiplicity of computer systems available.
	In network systems, users explicitly move/share files, submit jobs for processing and other perform other similar tasks.
	In distributed systems, tasks such as moving/sharing files, submitting jobs for processing and other similar tasks are handled automatically by the operating system (OS).

Evaluation

Advantages	Disadvantages
Increase processor throughput due to the use of parallel processing.	A more complex operating system is required in order to be able to interface and manage multiple processor units.
Lower cost than using multiple processors across multiple computer systems because the processors share resources such as the power supply and motherboard.	

Increased reliability because failure of one processor does not affect the other processors and will only slow down the computer system.	
---	--

1.6 Operating system (OS) layers

User interfaces

In an operating system (OS), the top layer is the user interface. This is the only layer explicitly visible to the user.

The user interface may be a:

- terminal
- graphical user interface (GUI)
- text-based command prompt; and/or
- a visual way of interacting with a computer using items such as windows, icons and menus.

History of the graphical user interface (GUI)

The first company to develop a graphical user interface (GUI) was Xerox PARC. They developed the “Alto” personal computer. It had a bitmapped screen and was the first computer to have a “desktop” screen with a graphical user interface (GUI).

The “Alto” personal computer was not a commercial product. However, several thousand units were manufactured and used at Xerox’s offices and several universities.

This development was a large influence on the design of personal computers during the late 1970s and early 1980s. Notable examples include:

- Three Rivers PERQ;
- Apple Lisa;
- Apple Macintosh; and
- the first Sun Workstations.

1.7 Kernel mode

Kernel mode vs user mode

Kernel Mode	User Mode
-------------	-----------

<p>Operating systems (OSs) run in kernel mode. This allows:</p> <ul style="list-style-type: none"> • execution of privileged machine instructions; and • complete access and control of all the hardware. 	<p>Other software runs in user mode. In this mode, instructions that affect control of the machine are forbidden. For example:</p> <ul style="list-style-type: none"> • web browsers; • e-mail software; and • music players.
---	--

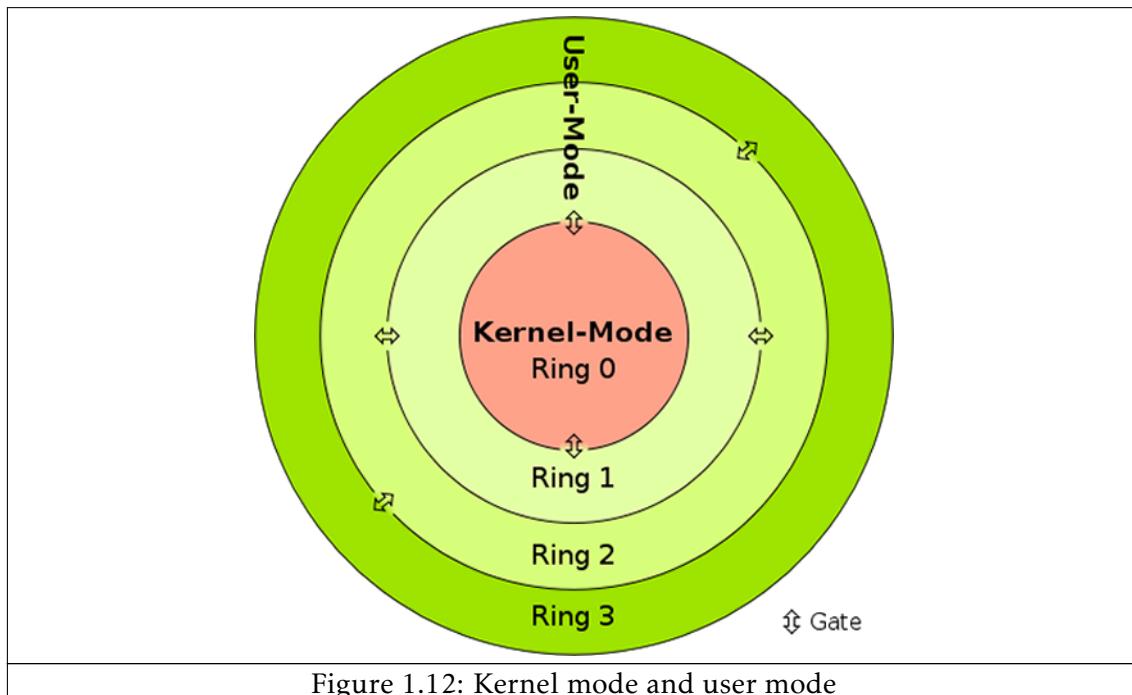


Figure 1.12: Kernel mode and user mode

Ring 0 represents the kernel mode.

Rings 1-3 represent the user mode.

Kernel mode protection

These rings allow separation between the operating system (OS) and user programs for security and protection purposes.

If user mode had unrestricted access to all of the machine instructions:

- a user could inadvertently obtain a virus or write code that is capable of causing damage to the system, and therefore it is necessary to prevent any instructions from directly controlling the machine;
- a program may use resources unfairly, such as holding the CPU or memory, and therefore harm the execution of other programs; and/or
- sensitive machine instructions could be used improperly which may lead to kernel mode errors.

Kernel mode errors are catastrophic.

Kernel Mode	User Mode
<p>A kernel panic represents the operating system (OS) attempting to prevent software causing any harm to the computer system and to recover from a kernel mode error on reboot.</p> <p>An example of a kernel panic is the “Blue Screen of Death” (BSOD) in Microsoft’s Windows.</p>	<p>Application errors where an exception was thrown due to an attempt to execute a privileged instruction, this is one that should only be accessed and executed by the kernel mode, represents the operating system (OS) preventing an application from having unrestricted access to all of the machine instructions.</p>

Figure 1.13: Kernel mode errors

System calls

A **system call** is the programmatic way in which a computer program requests a service from the kernel of the operating system on which it is executed. A system call is a way for programs to interact with the operating system.

Some privileged instructions can be called by a programmer via system calls. Operating systems (OSs) contain system calls for which provide a layer of services for user programs to implement some activities/request services. These are usually sensitive or privileged from the kernel.

All interactions with the hardware are implemented via system calls. For example, a system call may occur if an application requires interaction with a peripheral such as a printer.

Invoking a system call is similar to calling a general function. However, the difference is that a general function’s code is part of program itself, while the system call code is part of the operating system (OS). Different operating systems (OSs) offer different (limited) sets of system calls.

Call	Description
Process Mangament	
pid = fork()	Create a child process identical to the parent.
exit(status)	Terminate process execution and return status.
File Mangament	
n = read(fd, buffer, nbytes)	Read data from a file in to a buffer.
n = write(fd, buffer, nbytes)	Write data to a file
Process Mangament	
seconds = time(&seconds)	Get the elapsed time since Jan 1. 1970

Figure 1.14: Unix system calls

2. Process Management

2.1	Programs and processes	14
2.2	Process life cycle	15
2.3	Program execution	21
2.4	Concurrency	23
2.5	Process scheduling	24

2.1 Programs and processes

Definitions

A **program** is the code written by a programmer.

A **process** (or job/task) shows a program in execution and is a particular instance of a program. These may be shown in a monitoring software, such as Windows Task Manager.

Data are stored values used for the computations by the process.

How it works

A single program may have multiple processes that are currently running.

Each process can share the same code for the program. This is possible as each process uses its own address space, a list of memory locations which the process can read and write. These memory locations contain the program's code instructions and data.

A program is only code however, once it is run, a process is started, and it becomes instructions and data.

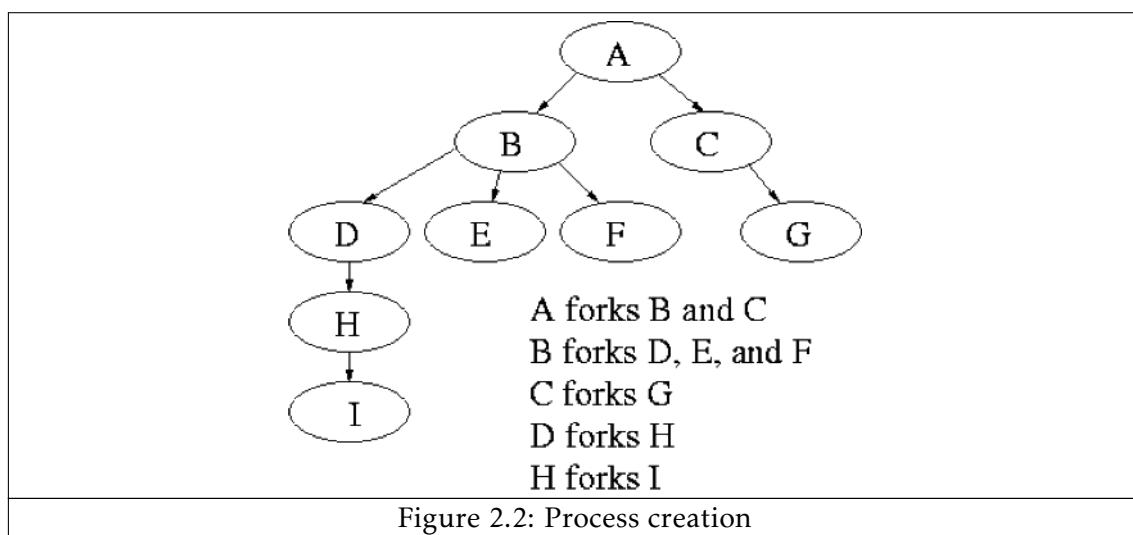


2.2 Process life cycle

Process creation

A process can be created by:

- a user
 - a program may be executed by a user, such as via a double-click using a graphical user interface (GUI) or by typing in a command, and “trigger” the processor to load the program’s executable file containing the program code; or
- another process
 - an existing process may create another process by spawning/forking – the process that creates a new process is called the parent process while the new process is called the child process. The child process may also spawn a new process forming a tree of processes, as demonstrated in the diagram below.



Process table and process control block

A **process identification number (PID)** is a unique identifier given to a new process when it is created.

A **process control block (PCB)** holds all of the information about a process. It is created when new process is created.

A **process table** stores the process identification numbers (PIDs) for a process and a pointer to the respective process control block (PCB) for that process. This is managed by the operating system (OS).



Figure 2.3: Process table with respective process control blocks

The process descriptor fields in the process control block (PCB) may differ between operating system (OS). An example of possible process descriptor fields is shown by those used by the Minix operating system (OS) are shown below.

Process Management	Memory Management	File Management
Registers	Pointer to text segment	UMASK mask
Program counter	Pointer to data segment	Root directory
Program status word	Pointer to bss segment	Working directory
Stack pointer	Exit status	File descriptors
Process state	Signal status	Effective uid
Time when process started	Process ID	Effective gid
CPU time used	Parent process	System call parameters
Children's CPU time	Process group	Various flag bits
Time of next alarm	Real uid	
Message queue pointers	Effective uid	
Pending signal bits	Real gid	
Process ID	Effective gid	
Various flag bits	Various flag bits	

Figure 2.4: Process descriptor fields in Minix

Although Minix is a fully-featured operating system (OS), it is a small operating system (OS) and therefore the processor descriptor fields are less complex than other operating systems (OSs).

Three-state model

The **three-state model** shows how a process, once initiated, can be in one of three states. The current state of a process is stored in its respective process control block (PCB).

A process, once initiated, can be in one of three main states:

- running – actually using the CPU to perform a task;
- ready – ready to run but waiting for the CPU as it has not yet had time on the CPU has been temporarily stopped to let another process run; or

- ready
 - unable to run until some external event occurs, such as:
 - waiting for an interrupt, this is a message from the hardware saying that a resource is now available to read from – such as waiting for an input/output (I/O) operation to complete; and
 - waiting for another process to finish accessing a shared resources – for example: a file; memory; or an external peripheral, such as a printer.



In reference to the diagram above, transitions can occur between process states.

Transition	Diagram Number	Explanation
running → blocked	1	Process blocks for input. For example, if the process is waiting for some input from I/O.
running → ready	2	Scheduler picks another process. The process has had opportunity to run and is flagged as no longer currently running, so that another process can run.
ready → running	3	Scheduler picks this process. The next process that is ready to run is set to running to allow access to the CPU.
blocked → ready	4	Input becomes available. The process has received input from I/O or the process sends an interrupt and the interrupt service routine (ISR) is executed, the scheduler is called to transition the process from blocked to ready.

Figure 2.6: Transitions between processor states

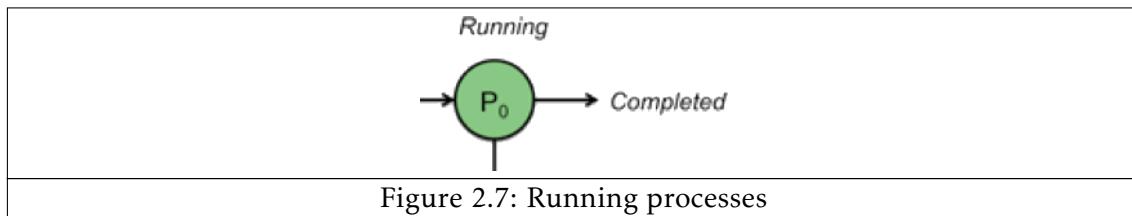
The transition between process states is dependent on the scheduling algorithm used. A clock is used to send a signal to stop the current process, move it from running to ready and then run the scheduler to find out what process should be processed next and the next process will be made to running.

State queues

At any time, a process is in only one state.

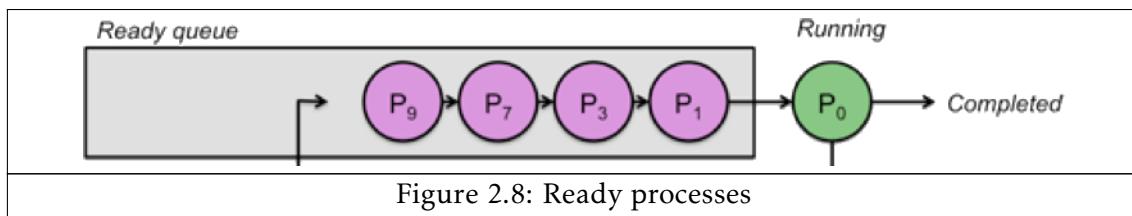
Running processes

At any time, at most one process is in the running state. This is because a single-core processor is only able to process one instruction coming from a single core at a time. If a computer system has a multi-core processor, this rule applies to each individual core on the processor, rather than the processor as a whole, as they are able to complete parallel execution.



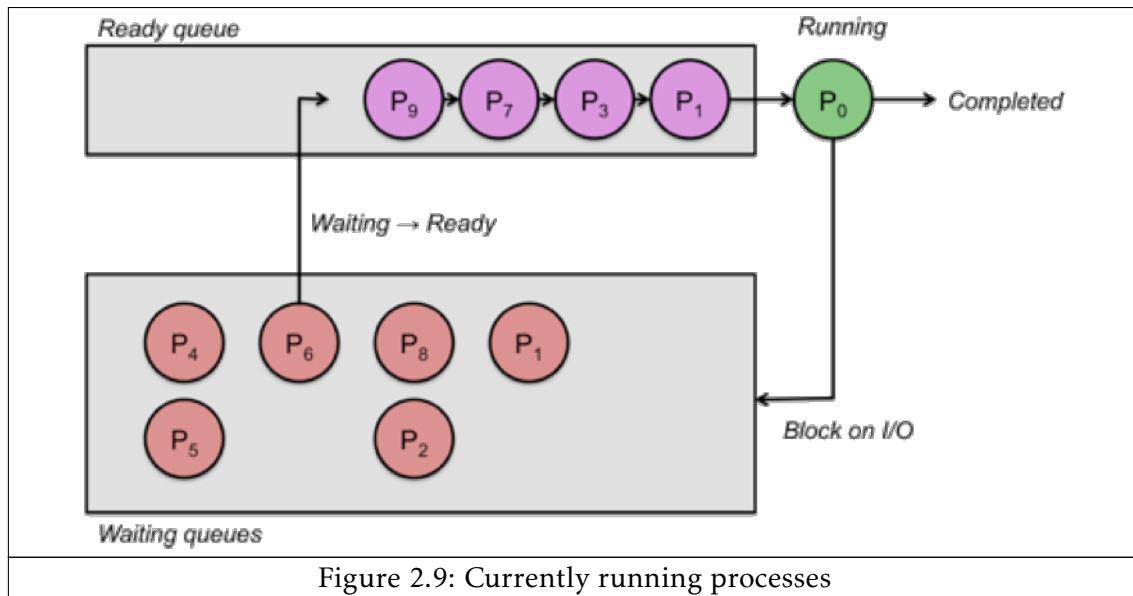
Ready processes

There may be a queue of processes in the ready state.



Blocked processes

There may be several queues of processes in the blocked state, where each queue represents one resource for which processes in that queue are waiting.



Process termination

Process termination is the end of life for a process, this can occur in two general ways.

Voluntary termination

Voluntary termination of a process represents the end of life for a process where its termination was intended by the user or the programmer.

This can be a:

- **normal exit** – the process has done its work; or
- **error exit** – the process itself handles and “catches” an error – for example, some try-catch code is implemented to check if a condition is met, such as if the definition of a variable is present.



Involuntary termination

Involuntary termination of a process represents the end of life for a process where its termination was not intended by the user or the programmer.

This can be due to:

- **a fatal error** – an error is detected by the operating system's (OS's) protected mode – for example, an exception has been thrown due to reference to a non-existent memory location or division by zero; or
- **being killed by another process** – a process may execute a system call that causes the operating system (OS) to kill another process – this process may have control over the killed process and this may be due to that process being the parent process of the killed process.

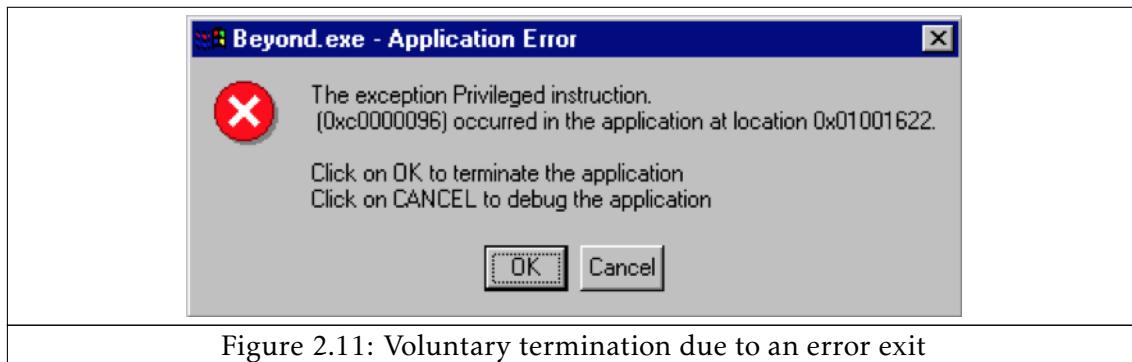


Figure 2.11: Voluntary termination due to an error exit

2.3 Program execution

The simple fetch-execute cycle

Definition

The **fetch-execute cycle** is an operational process in which a computer system retrieves a program instruction from its memory, determines what actions the instruction dictates, and carries out those actions.

How it works



Once a process is started, instructions are fetched from memory and executed using the CPU. This process will continue until the process is terminated.

This predictable cycle is only feasible to an extent as some processes may be slow or blocked and other may require immediate attention and cannot wait for the current process to terminate. For example, if a process blocks the processor (CPU) because it is waiting for an event to occur, such as a printer to finish its job, or if a high-priority process is supposed to execute as soon as possible. In these cases, interrupts are required.

Interrupts

Definition

An **interrupt** is a signal sent to the processor indicating that an event caused by hardware or software requires immediate attention.

Types of interrupts

- **Input/output – interrupt** – Caused by an input/output (I/O) device to signal completion or an error.
- **Timer interrupt** – Caused by a processor timer and is used to alert the operating system (OS) at specific instants.
- **Program interrupt** – Caused by error conditions within user programs or fatal errors.

When do interrupts occur?



Interrupts enable operating systems (OSs) to oversee several programs and input/output (I/O) events simultaneously.

This also means that single-core processors can effectively emulate the way in which multi-core processors deal with multiple instructions at a given time by switching between instructions intelligently. Due to the high clock speeds of modern processors, it

is easy to give the illusion that true multi-tasking, where two instructions are being processed at once, is taking place on a single-core processor.

Updated fetch-execute cycle

Given the introduction of interrupts, it is now necessary to update the fetch-execute cycle in order to include the possibility of an interrupt.



2.4 Concurrency

Definition

Concurrency describes the ability for a program to be decomposed into parts that can run independently of each other. This means that tasks can be executed out of order and the result would still be the same as if they are executed in order.

Why is concurrency required

Concurrency allows the processor (CPU) to run several processes. An example of this is shown by an interrupt occurring due to an input/output (I/O) event occurring (page 24).

Interleaving

Concurrency is able to achieve multitasking, that does not require parallel execution, by performing interleaved execution.



2.5 Process scheduling

The scheduler

Definition

A **scheduler** uses a scheduling algorithm to determine how to share processor time.

How it works

After an input/output (I/O) system call or interrupt handling, control is passed to the scheduler to decide which process to execute next.



The scheduler checks if the current process is still the most suitable to run at this moment in time. If it is control is returned to the process, otherwise:

- the state of the current process is saved in the process control block (PCB);
- the state of the most suitable process is retrieved from the process control block (PCB); and
- control is transferred to the newly selected process at the point indicated by the restored program counter (PC).

The action of storing the state of the current process and activating another process is called a context switch.

Context switching must be minimised to reduce overheads created by copying the state of processes and the time taken to perform the switch. However, it is still regarded as important to context switch when appropriate to avoid longer waiting times in the event of a blocked process.

Scheduling

Definition

Scheduling is the act of determining the optimal sequence and timing of assigning execution to processes.

Scheduling criteria

Different scheduling criteria may be selected depending on the use case for a given computer system.

- **CPU utilisation** – Aims to keep the CPU as busy as possible.
- **Efficiency** – Aims to maximise system throughput.
- **Fairness** – Aims to be fair to all running processes or to all users on a multi-user operating system (OS).

This means that different policies and algorithms for scheduling will exist to match these criteria.

Scheduling policies

A **non-preemptive scheduling policy** is one that allows processes to run until complete or incurring an input/output (I/O) wait. These scheduling policies can be described as passive.

A **preemptive scheduling policy** is one that allows processes to be interrupted and replaced by other processes, generally through timer interrupts.

Scheduling algorithms

Definitions

Arrival time is the instant at which a process is first created.

Service time is the time that it takes for a process to complete if it is in continuous execution.

The **waiting time** for a process is the sum of time spent in the ready queue during the life of the process. This does not include time that the process is blocked or waiting for input/output (I/O).

In order for a scheduling algorithm to be deemed as fair to all processes, the ratio between waiting time and run time should be about the same for each process.

Case study

Process	Arrival Time	Service Time
A	0	3
B	2	6
C	4	4
D	6	5
E	8	2

Figure 2.17: Case study

The case study shows different processes (A-E) that all have different arrival times and different service times.

The following examples of scheduling algorithms will refer to this case study to demonstrate their function.

First come, first served (FCFS) / First in, first out (FIFO) – Non-preemptive

In this algorithm, the first process to arrive is assigned to the processor (CPU) until it is finished. Meanwhile, any other processes that come along are queued up waiting to be processed.



Advantages	Disadvantages
Simple to implement.	Does not consider the priority of a process and therefore the important processes may not be completed quickly.
	Prevents other processes from starting while another is in progress and therefore, if processes are of varying sizes, there may be inefficiencies since a single process may take a long time to complete therefore leaving the user waiting before they can perform any other actions.
Favours long processes as all processes are given the opportunity to run until completion and therefore, some shorter processes may not be able to start processing until the longer processes are finished.	

Shortest job first (SJF) – Non-preemptive

In this algorithm, the process with the shortest estimated run time is assigned to the processor (CPU) until it is finished. Meanwhile, any other processes that come along are queued up waiting to be processed.

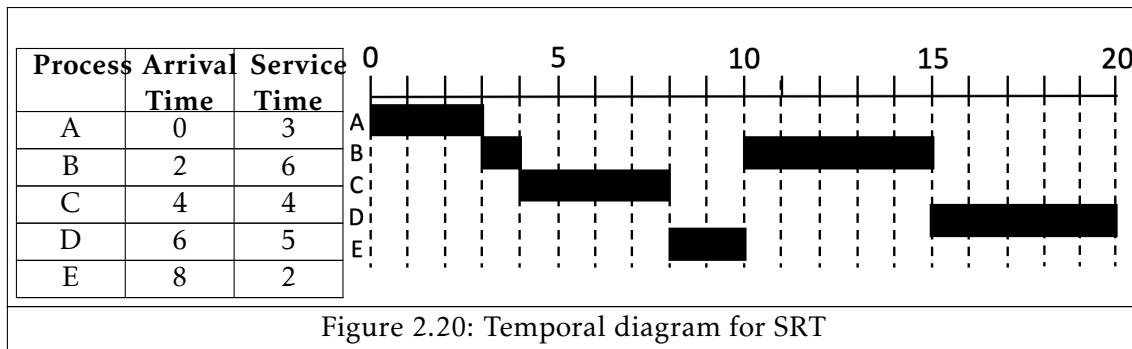


Figure 2.19: Temporal diagram for SJF

Advantages	Disadvantages
Simple to implement.	Does not consider the priority of a process and therefore the important processes may not be completed quickly.
Shorter processes are processed quickly because they take precedence.	Relies on an estimation of how long a process will take which could be incorrect.
Minimises the average time taken to complete a process because the shortest processes take precedence.	Prevents other processes from starting while another is in progress and therefore, if processes are of varying sizes, there may be inefficiencies since a single process may take a long time to complete therefore leaving the user waiting before they can perform any other actions.
Favours long processes as the processes with the shortest estimated time are given the opportunity to run until completion and therefore, some longer processes may not be able to start processing until there are no more shorter processes currently running or ready to be processed.	

Shortest remaining time (SRT) – Preemptive

In this algorithm, the process with the shortest estimated remaining time is assigned to the processor (CPU). If a process becomes ready that has a shorter remaining time, it will be pre-empted to allow the new process to start and the scheduler is switched to that new process.



Advantages	Disadvantages
High throughput because the number of processes completed is high due to the shortest processes taking precedence.	Does not consider the priority of a process and therefore the important processes may not be completed quickly.
Shorter processes are processed quickly because they take precedence.	Relies on an estimation of how long a process will take which could be incorrect.
Can be inefficient if a large process is in progress and shorter processes are being added to the queue because they will take precedence.	
Favours long processes as the processes with the shortest estimated time are given the opportunity to run until a process with a shorter estimated time is ready and therefore, some longer processes may not be able to start processing until there are no more shorter processes currently running or ready to be processed.	

Highest response ratio next (HRRN) – Preemptive

In this algorithm, the process with the highest priority is assigned to the processor (CPU). If a process becomes ready that has a higher priority, it will be pre-empted to allow the new process to start and the scheduler is switched to that new process.

The priority of a process may be based on memory, time and/or any other resource requirement such as:

$$\frac{\text{waitingtime} + \text{runtime}}{\text{runtime}}$$



Advantages	Disadvantages
High throughput because the number of processes completed is high due to the shortest processes taking precedence.	Does not consider the priority of a process and therefore the important processes may not be completed quickly.
Shorter processes are processed quickly because they take precedence.	Relies on an estimation of how long a process will take which could be incorrect.
	Can be inefficient if a large process is in progress and shorter processes are being added to the queue because they will take precedence.

Round robin (RR) – Preemptive

In this algorithm, each process is dispatched to the processor (CPU) on a “first in, first out” (FIFO) basis with a fixed time quantum.

Each time quantum is typically 10-20ms. Modern processor (CPU) clock frequencies are typically greater than 2GHz, which implies clock periods of 5×10^{-7} ms. This shows that the given time quantum is relatively large compared to a typical processor's (CPU's) clock period.

If a process experiences a timeout, this will mean that it has run over its fixed time quantum. In which case, the process will be interrupted and returned to the back of the queue.

A system designer may wish to choose a time quantum that is most appropriate for a given system. This may be done by measuring the average service time and waiting time for the processes that will be running on the system and design the fixed time quantum around these figures. It may be that a system designer wishes to minimise the number of processes that are interrupted by another process.



Figure 2.22: RR



Advantages	Disadvantages
Simple to implement.	Heavy overhead due to continuous context switches.
Suitable for some types of computer systems , such as those which will be running processes of similar priority and size.	Does not consider the priority of a process and therefore the important processes may not be completed quickly.
	Does not consider the size of a process and therefore, if processes are of varying sizes, there may be inefficiencies since a single process may take a long time to complete thus leaving the user waiting before they can perform any other actions.

Multi-level queueing

Definition

Multi-level queueing is a queue with a predefined number of levels which consist of several independent queues.

How it works

Multi-level queueing makes use of other existing scheduling algorithms.



Each queue has a different priority; the top queue has the highest priority and the bottom queue has the lowest priority. Processes are able to move between queues if their priority changes.

There is some form of inter-queue scheduling policy that governs the assignment of processes from each queue to the processor (CPU).

The queues in a multi-level queueing system may differ between operating system (OS). For example, the scheduler used by the Minix operating system (OS) uses multi-level queueing and implements 16 queues.

Advantages	Disadvantages
Allows scheduling optimisation as a system designer may be able to leverage the advantages of a range of different scheduling algorithms.	
Maintains common processes as it is possible to split processes into different queues depending on their nature. For example, input/output processes could be in one queue while processor (CPU) processes are in another queue.	
Helps to prevent bottlenecks because input/output (I/O) devices are slower than the processor speed and therefore maximising processes involving input/output (I/O) devices ensures that these devices are continuously busy.	

3. Threads and Concurrency

3.1	What are threads?	33
3.2	Sequential and concurrent programming	36
3.3	Sequential execution	38
3.4	Introduction to concurrent execution	40
3.5	Interleaving	41
3.6	User and kernel threads	45
3.7	Multi-threading models	47
3.8	Evaluation of concurrent programming	50

3.1 What are threads?

Definition

A **thread** is an independent path/sequence of execution within a process and can be managed independently by a scheduler. A process may contain many threads.

Multi-threading

How it works

As mentioned in the previous section, a process (or job/task) shows a program in execution and is a particular instance of a program. By default, these processes are run by means of “single execution thread”.

However, different parts of the same process could be parallelised in order to allow multi-threading.

Multi-threading provides a way of improving application performance and therefore improving the efficiency and/or usability of a computer system.



Figure 3.1: Single-threaded vs multi-threaded

Example



Threads vs processes

Comparison

Threads and processes share some similarities however, there are also some distinct differences.

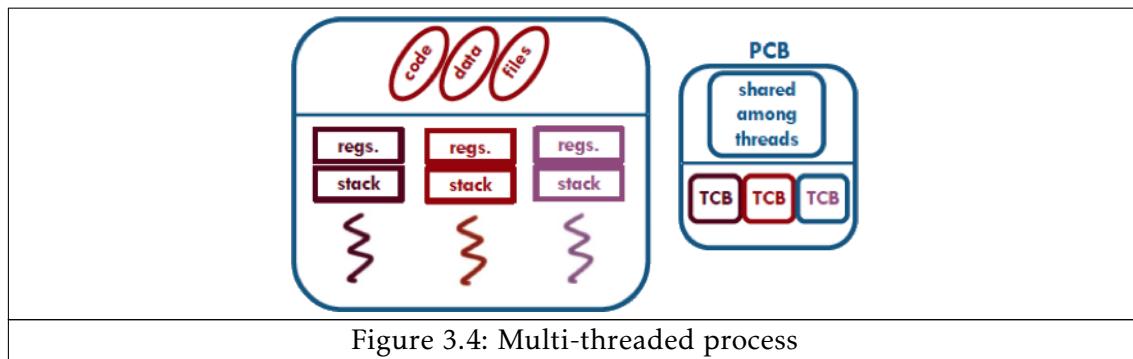
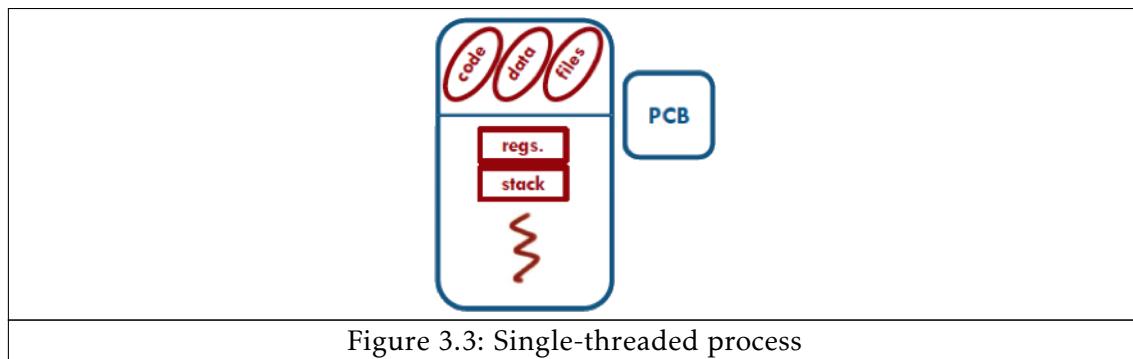




Figure 3.5: Process table and process control blocks

	Threads	Processes
Similarities	Sequential flow of control with <u>start</u> and <u>end</u> .	Sequential flow of control with <u>start</u> and <u>end</u> .
	At any time, a thread has a single point of execution.	At any time, a process has a single point of execution.
	Has its own execution context, stack (history) and program counter stored in a thread control block (TCB).	Has its own execution context, stack (history) and program counter stored in a process control block (PCB).
	Follows the three-state model in which the thread can be running, blocked or ready.	Follows the three-state model in which the process can be running, blocked or ready.
	Context switching can happen for threads.	Context switching can happen for processes.
	A thread can spawn another thread.	A process can spawn another process.
	A thread is often called a lightweight process.	
Differences	A thread cannot exist on its own, instead it exists within a process.	A process does not require a parent entity.
	Usually created and/or controlled by a process.	A process is not typically created and/or controlled by another process.
	Threads can share process properties, including memory and open files.	Processes cannot share process properties with other processes.
	Inexpensive creation and context switching as does not require separate address space.	Expensive creation and context switching as requires separate address space.
	When running multiple threads concurrently, they share an address space.	When running multiple processes concurrently, they are resources, such as memory, disk and printers.

Figure 3.6: Threads vs processes

Properties

Per process items	Per thread items
Address space	Program counter
Global variables	Registers
Open files	Stack
Child processes	State
Pending alarms	
Signals and signal handlers	
Accounting information	

Figure 3.7: Threads vs processes

Process properties are shared between threads. Thread properties are local and private to each thread.

3.2 Sequential and concurrent programming

Sequential programming

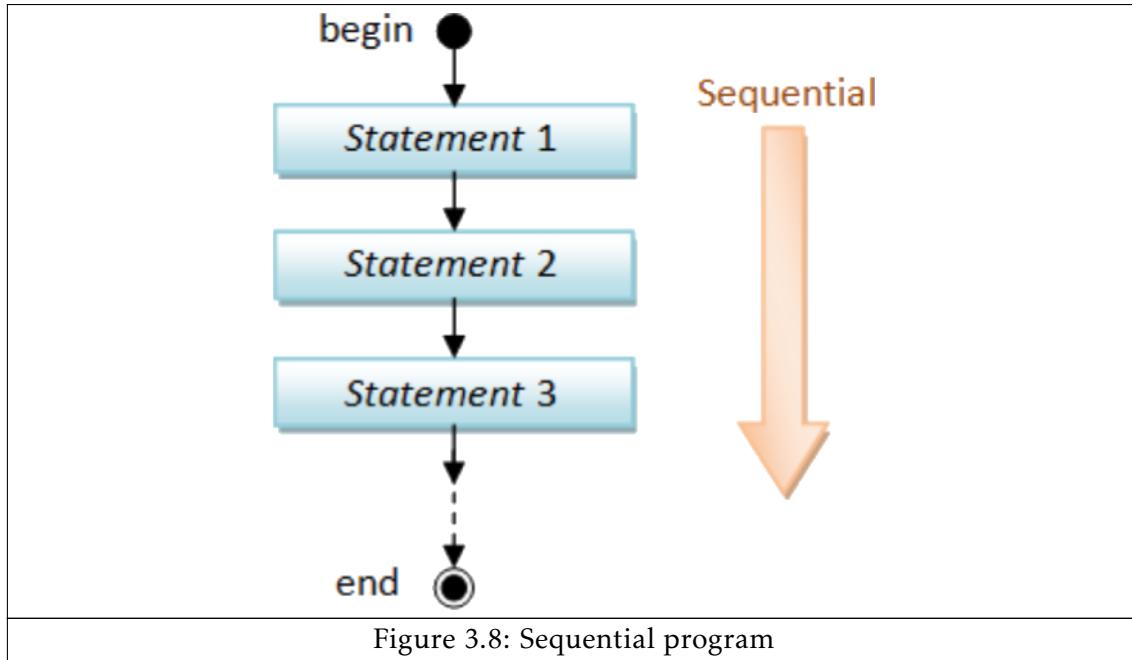
Definition

Sequential programming is the traditional activity of constructing a computer program using a sequential programming language.

How it works

This involves a programming methodology that assumes statements are executed in order/sequence.

Programs written using sequential programming are assumed to execute on a single-CPU system and have a single thread of control.



Evaluation

Advantages	Disadvantages
No additional support required from the programming language.	Lower processor throughput than concurrent programming as it cannot benefit from multitasking or concurrent processing.
No additional support required from the operating system as most old-school operating systems were generally single-threaded and therefore later generations of operating systems typically inherit this functionality.	Multiple computer systems that each have their own CPU may yield a higher cost than multi-CPU systems as more will be spent on resources, such as the power supply and motherboard, that could be otherwise shared in a multi-CPU system.
	Lower reliability than a multi-CPU system because, in the case of the failure of the processor, there is no redundancy.

Concurrent programming

Definition

Concurrent programming is the activity of constructing a computer program that takes advantage of concurrency allowed by the use of multiple threads of control.

How it works

Multiple threads of control allow a given process to perform multiple computations in parallel and to control simultaneous external activities.

The program may be run on both:

- a single-CPU system
 - where the computer program will take advantage of multitasking; and
- a multi-CPU system
 - where the computer program will take advantage of true parallelism.

Evaluation

Advantages	Disadvantages
Increase processor throughput due to the use of multitasking in a single-CPU system or parallel processing in a multi-CPU system.	Requires support from the programming language as it must implement techniques to deal with multitasking on a single-CPU system and/or parallel processing on a multi-CPU system.
A multi-CPU system generally yields a lower cost than using multiple CPUs across multiple computer systems because the processors share resources such as the power supply and motherboard.	Requires support from the operating system as it must support multi-threading in order to allow multitasking on a single-CPU system and/or interface and manage multiple CPU units on a multi-CPU system.
Increased reliability in a multi-CPU system because failure of one processor does not affect the other processors, instead the computer system may experience lower performance until fixed.	

3.3 Sequential execution

Definition

Sequential execution is where the execution of threads in a sequential program is executed in sequence/order with no overlapping.

Order and precedence

Explanation

In sequential execution, there is only one possible sequence of execution. This is because a sequential program gives the system strict instructions on the order of executing the statements in the program.

Importance

For example, a simple hypothetical program could be:

P;

Q;

R;

This tells the computer system that the statements must be executed in the order they are written, such that:

- P must precede Q; and
- Q must precede R.

High level

The importance of the order of precedence can be highlighted by demonstrating this idea in a high-level programming language.

Given the following program written in a high-level language:

```
x = 1;  
y = x + 1;  
x = y + 2;
```

it is possible to see that the final values of x and y depend on the order of execution of the statements.

System level

Given the following program written in a high-level language:

```
x = 1; P  
y = x + 1; Q  
x = y + 2; R
```

where each statement is assigned a letter respectively, each statement may be compiled in to several machine instructions.

Statement **P** is treated as a single machine instruction:

- **P1**: store 1 at the memory address of x.

Statement **Q** is broken in to three machine instructions:

- **Q1**: load the value of x in to a CPU register;
- **Q2**: increment the value in this register by 1; and
- **Q3**: store the value in this register at the memory address of y.

Statement **R** is broken in to three machine instructions:

- **R1**: load the value add of y in to a CPU register;

- R2: increment the value in this register by 2; and
- R3: store the result at the memory address of x.

The nature of sequential execution

The execution of statements **P**, **Q** and **R** at the program level (or high-level) as

$$P \rightarrow Q \rightarrow R$$

implies that the execution at the system level is as follows

$$P_1 \rightarrow Q_1 \rightarrow Q_2 \rightarrow Q_3 \rightarrow R_1 \rightarrow R_2 \rightarrow R_3,$$

given that **P** is compiled to a single machine instruction, whilst **Q** and **R** are compiled to three machine instructions – as seen on page 40.

Sequential execution has the following assumptions:

- total ordering – there is single-threaded computation, and therefore no overlap in the execution of the statements;
- deterministic – the same input will always result in the same output; and
- sequence – users will specify a strict sequence of steps required in order to achieve a desired goal.

However, this does not apply in many practical applications, for which a sequence of steps is not required.

3.4 Introduction to concurrent execution

Definition

Concurrent execution is where the execution of threads in a concurrent program is occurring asynchronously, meaning that the order in which tasks are executed is not pre-determined.

The squares example

In this hypothetical example, a person desires to have a list with the results of all of the squares (2^n) from 1 to 100000.

A group of 100000 people are split into heavily uneven teams and assigned the same task to complete all of the calculations in order to achieve the desired result.

It is given that each calculation takes n amount of time.

Team 1		Team 2	
Number of members	1	Number of members	99999

Strategy	One person should complete all of the calculations.	Strategy	Each member is assigned a number between 1 to 100000. Each member should calculate the respective square for the number they are assigned.
Time taken	$100000n$	Time taken	n

This shows that Team 2 was $100000x$ faster than Team 1. This was because it was possible to decompose the larger task in to smaller sub-tasks and assign each of those tasks to a separate resource, which in this case is one person.

This example forms that basic concept of concurrent execution.

The nature of concurrent execution

Concurrent execution dismisses many of the assumptions required for sequential execution (page 41).

Calculations may be carried out without total ordering. As a result, calculations may be carried out in parallel and overlapping is therefore allowed.

In the example above, each individual person in team 2 carried out their operations in sequence.

In the example above, the operations in the whole computation can be viewed as being in partial order. However, the order does not matter here because there is no dependency between the calculations. This is because the output from any given calculation is not required as an input to any other given calculation.

However, in general, concurrent execution is non-deterministic, and therefore the same input generally means different output due to ordering. This is because there are many cases where the order of operation does matter.

3.5 Interleaving

Why is interleaving required?

Concurrent execution on a computer system with a multi-core CPU or multiple CPUs can make use of parallel processing in order to run threads asynchronously. However, this is not possible on a computer system with a single-CPU that consists of only one core.

As a result, interleaving is used in order to switch execution between threads.

It is important to note that the operations within each thread are strictly ordered, but the interleaving of the operations are not ordered and are interleaved in an unpredictable order.

Calculating interleavings

Formula

It is possible to calculate the number of interleavings given the formula:

$$\text{number of interleavings} = \frac{t_1 + t_2 + \dots + t_n}{t_1!t_2!\dots t_n}$$

where

- t_n represents the number of statements/operations in each thread.

For example, in a concurrent program that has two threads, the formula may be adjusted to:

$$\text{number of interleavings} = \frac{(n+m)!}{n!m!}$$

where

- n represents the number of statements/operations in the first thread (t_1); and
- m represents the number of statements/operations in the second thread (t_2).

Example

A high-level concurrent program may spawn new threads.



In this example, the concurrent program spawns the threads t_1 and t_2 where:

- t_1 has three statements – A, B and C; and
- t_2 has two statements – S and T.

There are two threads and therefore it is possible to use the adjusted formula to calculate the number of interleavings in this concurrent program:

$$\text{number of interleavings} = \frac{(n+m)!}{n!m!}$$

$$\text{number of interleavings} = \frac{(3+2)!}{3! \times 2!}$$

$$\text{number of interleavings} = \frac{6!}{3! \times 2!}$$

$$\text{number of interleavings} = \frac{6 \times 5 \times 4 \times 3 \times 2 \times 1}{(3 \times 2 \times 1) \times (2 \times 1)}$$

$$\text{number of interleavings} = \frac{120}{6 \times 2}$$

$$\text{number of interleavings} = \frac{120}{12}$$

$$\text{number of interleavings} = 10$$

There are 10 possible interleavings, thus yielding 10 possible different execution sequences.

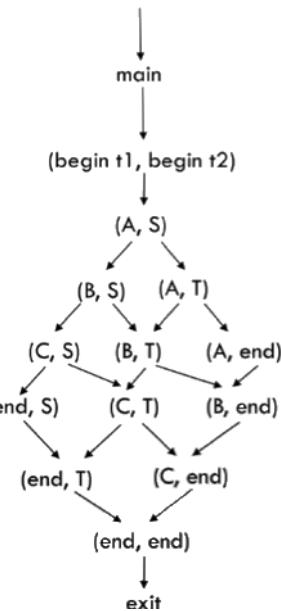


Figure 3.10: Visual representation of the possible execution sequences

A run of the program corresponds to an interleaving sequence. Each interleaving sequence determines a unique sequence of executing the statements. Repeated runs with the same input will likely trace different interleavings.

Growth of interleavings

The number of interleavings grows extremely quickly given an increase in:

- the number of threads in the concurrent program; or
- the number of statements/operations in one or more of the concurrent program's threads.

This can be demonstrated by increasing the number of operations in the previous example:

- t_1 now has four statements – A, B, C and D; and
- t_2 now has five statements – S, T, U, V and W.

and therefore:

$$\begin{aligned} \text{number of interleavings} &= \frac{(n+m)!}{n!m!} \\ \text{number of interleavings} &= \frac{(4+5)!}{4! \times 5!} \\ \text{number of interleavings} &= \frac{9!}{4! \times 5!} \\ \text{number of interleavings} &= \frac{9 \times 8 \times 7 \times 6 \times 5 \times 4 \times 3 \times 2 \times 1}{(4 \times 3 \times 2 \times 1) \times (5 \times 4 \times 3 \times 2 \times 1)} \\ \text{number of interleavings} &= \frac{362880}{24 \times 120} \\ \text{number of interleavings} &= \frac{362880}{2880} \\ \text{number of interleavings} &= 126 \end{aligned}$$

3.6 User and kernel threads

User threads

User threads are created and managed by a user level library, typically without the knowledge of the kernel.



The diagram shows that:

- all of the threads for a given process is present within the user space; and
- the thread table is present within the process.

User threads are:

- fast to create and manage; and
- portable to any operating system (OS).

If one user thread is blocked, all other threads in the same process are also blocked. For example, in a word processor application, a thread that handles a printing event would block all other threads and therefore prevent the user from interacting with other aspects of the application.

Multi-threaded applications cannot take advantage of parallel execution on computer systems with a multi-core CPU or computer systems with multiple CPUs.

Kernel threads

Kernel threads are directly managed and supported by the operating system's (OS's) kernel.



The diagram shows that:

- all of the threads for a given process is present within the user space; and
- the thread table is present within the kernel space, rather than the process itself or the user space.

Kernel threads are:

- slower to create and manage than user threads; and
- specific to the operating system (OS).

If one user thread is blocked, all other threads in the same process are scheduled and not blocked. For example, in a word processor application, a thread that handles a printing event would no longer block all other threads and therefore would allow the user to interact with other aspects of the application whilst the printing event occurs.

Can take advantage of parallel execution on computer systems with a multi-core CPU or computer systems with multiple CPUs.

3.7 Multi-threading models

Why is multi-threading mapping required?

The kernel is generally not aware of the user threads present in a process. Therefore, a thread library must map user threads to kernel threads.



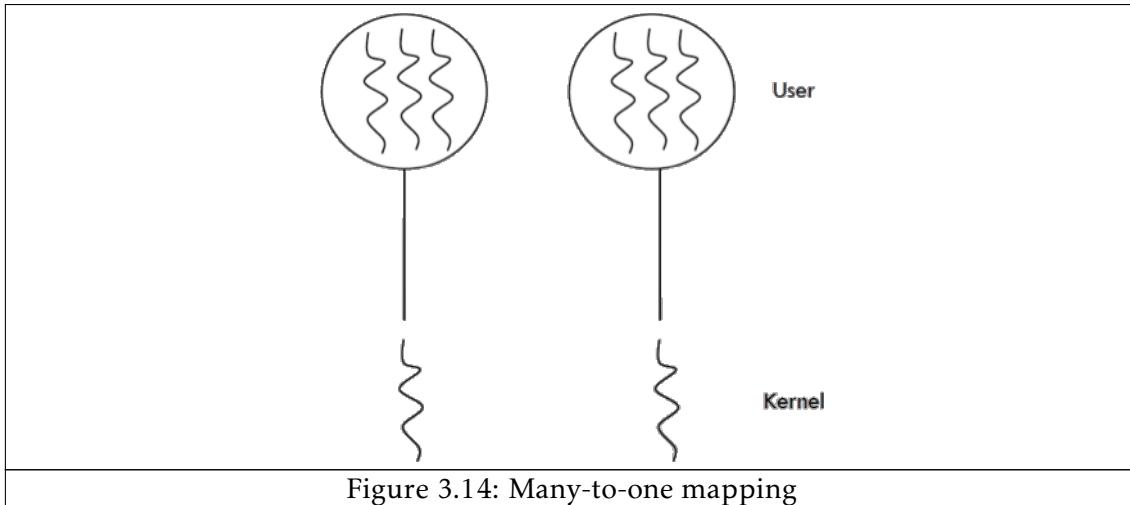
The diagram shows that there must be some relationship between the user threads and the kernel threads. This relationship may be defined using different mappings, including:

- many-to-one;
- one-to-one; and
- many-to-many.

Many-to-one mapping

How it works

All user threads from each process map to one kernel thread.



Evaluation

Advantages	Disadvantages
Portable as there are few system dependencies.	No parallel execution of threads.
	No concurrency as all threads in a process are blocked if another thread is blocked, for example if the thread is waiting for an input/output (I/O) interrupt.

One-to-one mapping

How it works

Each user thread maps to a single kernel thread.



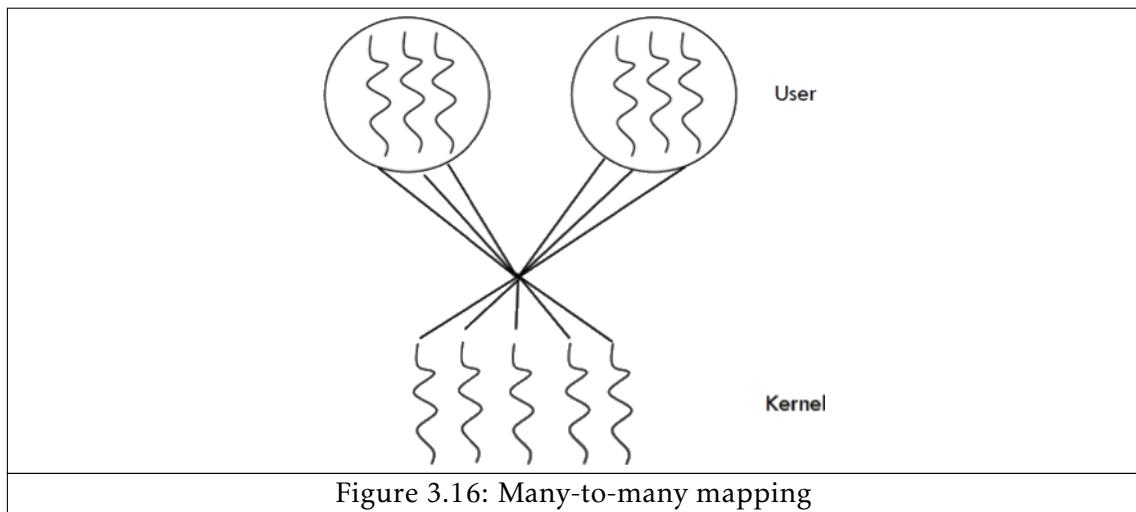
Evaluation

Advantages	Disadvantages
Concurrency as all threads in a process are not blocked if any given thread becomes blocked.	Slow as there is management overhead because the kernel is involved for every user thread.
Performance as it can take advantage of multiple CPUs.	Restricted as there is typically a limit on the number of threads.
	Creating user threads requires the corresponding kernel support.

Many-to-many mapping

How it works

Many user threads multiplex to an equal or smaller number of kernel threads.



Advantages	Disadvantages
Performance as it can take advantage of multiple CPUs.	Complexity and therefore implementation difficulties.
Flexible as there is no limit on the number of threads.	

3.8 Evaluation of concurrent programming

Advantages	Disadvantages
 <p>Parallelism. It improves the efficiency of program execution in computer systems with multiple CPUs by allows tasks/operations to be split up and executed independently on each CPU.</p>	<p>Debugging complexity as concurrent programs are non-deterministic and therefore it can be difficult to trace a problem/bug in the code as the same input will generally not result in the same output.</p>
 <p>Multi-tasking. It improves the utilisation of the CPU in a computer system that only has a single CPU. This allows multiple tasks/operations to run alongside each other and appear to be processed simultaneously.</p>	<p>No protection between threads.</p>
<p>Increases application responsiveness, for example, in a word processor application one thread could be responsible for responding to user input/output (I/O) while other threads perform tasks in the background.</p>	<p>Concurrent processes must interact with each other in order to share resources or exchange data.</p>
<p>Suited to some applications as there are some practical applications that are non-deterministic and concurrent as the order of program operations is determined by other external events. This is useful for applications that need to handle multiple events.</p>	<p>Synchronisation must be promoted in order to determine when, how and with what language abstractions computation events can be synchronised in order to eliminate unacceptable outputs.</p>
	<p>Distribution must be taken care of in order to consider how threads can be distributed among a number of CPUs and how one thread is able to interact with another thread on a different CPU.</p>

	<p>Error-prone. Examples of major concurrent programming errors include:</p> <ul style="list-style-type: none">• Therac-25 - A computerised radiation therapy machine whose errors contributed to accidents causing deaths and serious injuries.• Mars Rover Pathfinder – Problems with interaction between concurrent tasks caused periodic software resets, thus reducing availability for exploration.
--	---

4. Synchronisation and Mutual Exclusion

4.1	Race condition	52
4.2	Inter-process synchronisation	54
4.3	The producer-consumer problem	57
4.4	The dining philosophers problem	64

4.1 Race condition

Definition

A **race condition** describes the competition for resources in a critical section caused by interleaving/thread interference.

Why do race conditions happen?

Race conditions occur due to interleaving/thread interference.

Interleaving/thread interference describes an undesired outcome resulting from non-deterministic, concurrent usage of shared resources.

This happens because, in general, concurrent execution is non-deterministic, and therefore the same input generally means different output due to ordering. This is because there are many cases where the order of operation does matter.

Examples

Racing for memory access

A race condition may occur when two threads attempt to access the same location memory, such as registers or RAM, at the same time.



In this example, the two threads t_1 and t_2 manipulate the same variable where:

- t_1 increments the variable c ; and
- t_2 decrements the variable c .

As seen before (page 40), each statement may be compiled in to several machine instructions.

The increment ($c++$) instruction is broken in to three machine instructions:

- retrieve c ;
- increment retrieved value; and
- store result in c .

The decrement ($c--$) instruction is also broken in to three machine instructions:

- retrieve c ;
- decrement retrieved value; and
- store result in c .

As a result, one interleaving possibility is as follows:

- t_1 : retrieve c ;
- t_2 : retrieve c ;
- t_1 : increment retrieved value; (result is 1)
- t_2 : decrement retrieved value; (result is -1)
- t_1 : store result in c ; (c is now 1)
- t_2 : store result in c ; (c is now -1)

This example shows that the race condition has caused the result from thread t_1 to be lost as it has been overwritten by the result from thread t_2 .

Racing for peripheral access

A race condition may also occur when two threads attempt to access the same peripheral, such as a printer spooler directory, at the same time.



Figure 4.2: Printer spooler directory

In this example, the two processes *A* and *B* attempt to access the printer spooler directory at the same time:

- the next available printer job slot is 7;
- process *A* and *B* access printer job slot 7 simultaneously;
- process *A* reads printer job slot 7 and a timer interrupt occurs that causes a context switch to process *B* before process *A* has opportunity to store any data;
- process *B* reads printer job slot 7 and stores its job data and increments the values; and
- another timer interrupt occurs that causes a context switch to process *A* that then stores its job at printer job slot 7.

This example shows that the race condition has caused the job data stored in the printer spooler directory by process *B* to be lost as it has been overwritten by the job data from process *A*.

4.2 Inter-process synchronisation

Definition

Inter-process synchronisation involves techniques that are designed to prevent race conditions and allows threads/processes to share resources.

Behaviour of threads

Threads in a computer system may behave in two possible ways:

- competing – two or more processes compete for the same computing resource, for example access to a particular memory cell; or
- cooperating – two or more processes may need to communicate with one another, thus causing information to be passed from one to the other.

Inter-process synchronisation is required to manage both threads that are competing and threads that are cooperating.

An operating system (OS) itself contains both threads that are competing and threads that are cooperating.

How mutual exclusion works

The solution to preventing race conditions is by implementing mutual exclusion on any given critical section/region.

The **critical section/region** is code in a process that involves sensitive operations on a shared resource.

Mutual exclusion is the requirement that one thread of execution never enters its critical section/region at the same time that another concurrent thread of execution enters its own critical section/region.



When a thread/process enters its critical section/region, no other thread/process may also enter its critical section/region.

This is demonstrated in the diagram above:

- at time interval T_1 , process A enters its critical section/region;
- at time interval T_2 , process B attempts to enter its critical section/region;
- process B is blocked from entering its critical section/region until process A leaves its critical section/region;
- at time interval T_3 , process A leaves its critical section/region and therefore process B is allowed to enter its critical section/region; and
- at time interval T_4 , process B leaves its critical section/region.

This shows that mutual exclusion is enforced as no two threads/processes are simultaneously inside their critical sections/regions.

For mutual exclusion to be effective:

- no assumptions may be made about the speeds or the number of threads/processes;
- no threads/processes running outside its critical section/region may block other threads/processes, such that a thread/process that is not in its critical section/region cannot prevent other threads from entering their critical section/region; and
- no thread/process should have to wait forever to enter its critical section/region.

Mutual exclusion is a major design issue in operating systems (OSs) as consideration must be taken in order to prevent race conditions while maintaining parallelism and efficiency.

How mutual exclusion is implemented

Mutual exclusion is implemented using semaphores.

A **semaphore** is a system tool used for the design of correct synchronisation protocols. This was introduced by Edsger Dijkstra in the 1962/1963. Semaphores are implemented using a variable or abstract data type and are used to control thread/process access to a resource. They are typically integer values that accept only non-negative values.

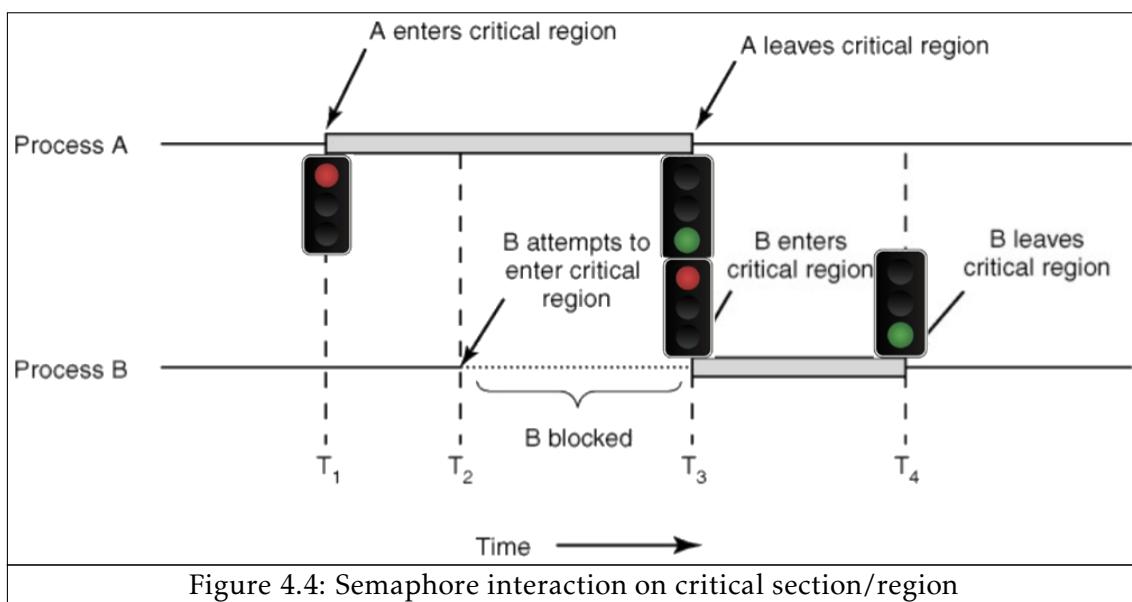


Figure 4.4: Semaphore interaction on critical section/region

The diagram shows that semaphores allow the CPU to context switch between threads/processes when one becomes blocked.

It is convenient to write entry and exit protocols using a single atomic statement. This statement is atomic and therefore is indivisible, meaning that the statement cannot be interrupted.

As mentioned before, a semaphore, denoted by S , is an integer that takes only non-negative values. Only two atomic (indivisible) statements are permitted, as shown below.

Statement	Statement Implementation	Usage
-----------	--------------------------	-------

<code>wait(s)</code>	<code>wait(s)</code> { if (S > 0) { S--; } }	If a thread/process is in its non-critical section/region and wishes to enter its critical section/region, this statement will be performed. This means that the thread/process will be blocked until $S = 0$ evaluates to <i>True</i> .
<code>signal(s)</code>	<code>signal(s)</code> { S++; }	If a thread/process is in its critical section/region, this statement will be performed. This helps to achieve mutual exclusion as it prevents $S = 0$ from evaluating to <i>True</i> until the thread/process has left its critical section/region.

This is a good solution as there is no possibility for a race condition as these statements will always be enforced due to the fact that they are atomic (indivisible) statements and cannot be interrupted.

4.3 The producer-consumer problem

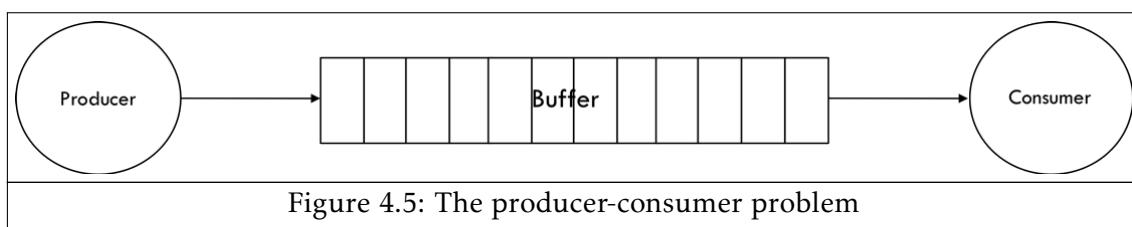
Problem description

The producer-consumer problem is a classical inter-process communication problem in which:

- a producer repeatedly produces items and places them in to a buffer; and
- a consumer consumes the items one-by-one by taking them from the buffer.

This problem has the following requirements:

- the buffer must be assumed to be first in, first out (FIFO);
- the producer may produce a new item only at a time when the buffer is not full;
- the consumer may consume an item only at a time when the buffer is not empty; and
- the process terminates when all items produced are eventually consumed.



The problem arises when attempting to devise a method that is able to:

- put the producer to “sleep” when the buffer is full to prevent further items being produced when there is no space in the buffer; and
- “wake” the consumer when the buffer is not empty as there is possibility to consumer when the buffer is not empty.

Possible solution

This problem could be solved by keeping track of the number of items in the buffer.

This could be achieved by implementing loops in the producer class and consumer class.

```

LOOP
{
    Produce item i          //produce item
    if ( itemCount == N )   //end of buffer
    {
        sleep(producer);
    }

    Put item i;            // place item in to buffer
    itemCount++;           // increment buffer count
    if ( itemCount == 1 )   // buffer nearly empty
    {
        wakeup(consumer);
    }
}

```

Producer class

```

LOOP
{
    if ( itemCount == 0 )   // buffer empty
    {
        sleep(consumer);
    }

    Remove item j;         // remove item from buffer
    itemCount--;           // decrement buffer count
    if ( itemCount == N-1 ) // buffer has space
    {
        wakeup(producer);
    }
    Consume item j;        // consume item
}

```

Consumer class

The loop in the producer class would be running as one thread and the loop in the consumer thread would be running as another thread. These two threads would be running in parallel. As a result, if the threads in the solution are interleaved, a race condition may occur, which in turn, may cause a deadlock.

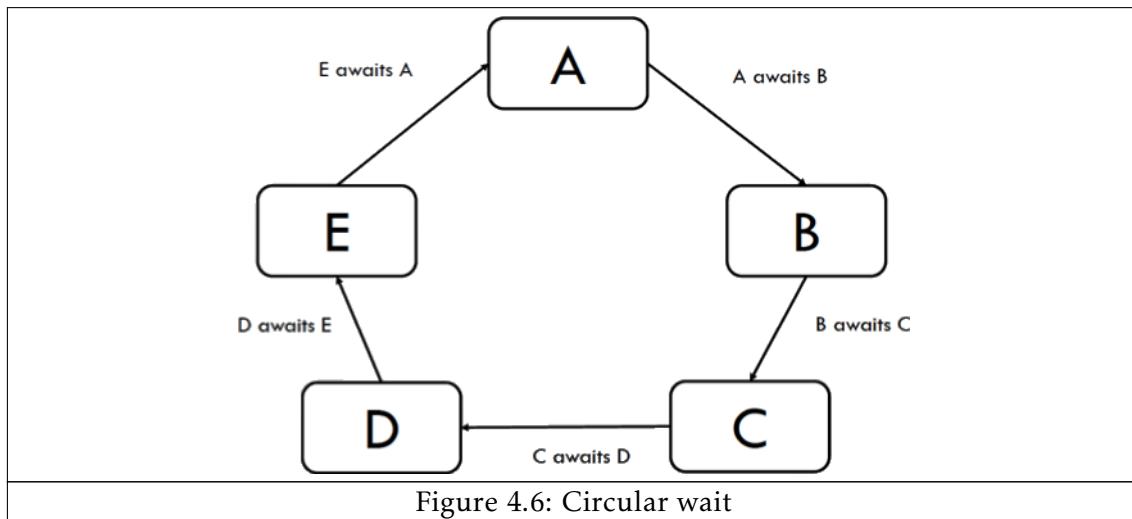
Deadlocks

A **deadlock** occurs when two or more threads wait for each other to finish.

Four conditions must be hold simultaneously in order for a deadlock to occur:

- mutual exclusion – a resource can be assigned to, at most, one process at a time;
- hold and wait – processes holding resources are permitted to request and wait for additional resources;

- no pre-emption – resources previously locked cannot be forcefully unlocked by another process, instead they must be released by the holding process; and
- circular wait – there must be a chain of processes, such that each member of the chain is waiting for a resource held by the next member of the chain, as shown in the diagram below.



A deadlock may occur in the possible solution described previously.

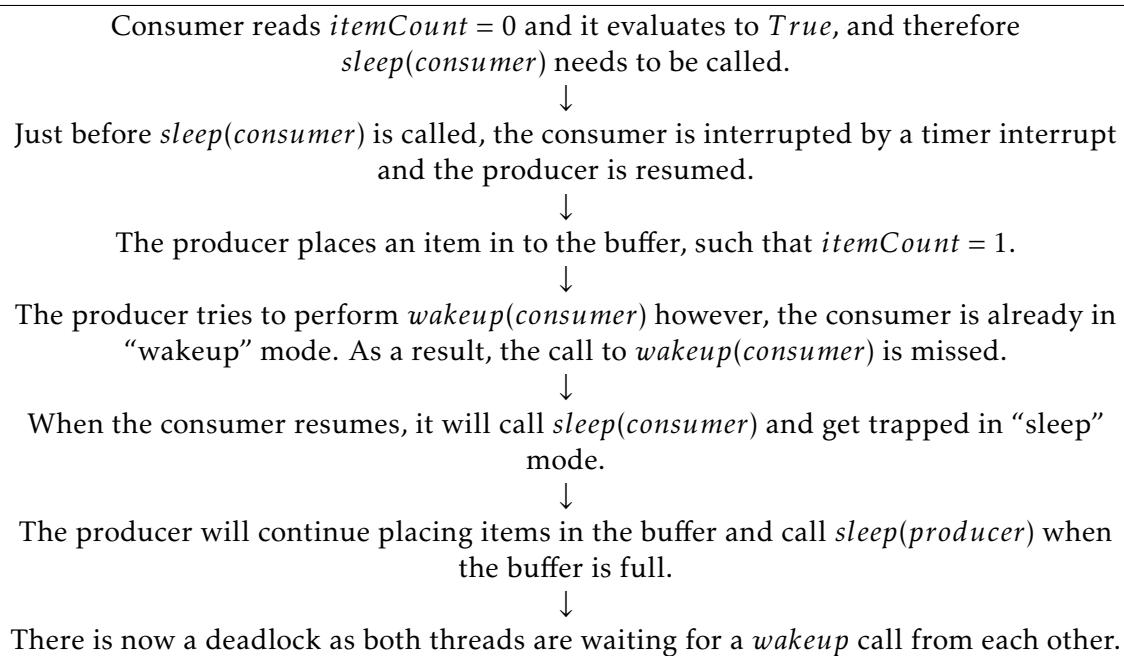


Figure 4.7: Possible deadlock

This possible deadlock shows that another solution is required to effectively solve the producer-consumer problem.

Solving the problem using semaphores

It is assumed that

```
ItemsReady = 0
SpacesLeft = N //size of buffer
```

```
LOOP
{
    Produce item i      // produce item
    Wait(SpacesLeft)   // decrement semaphore

    Put item i;        // place item in to buffer
    Signal(ItemsReady) // increment
}
```

Producer class

```
LOOP
{
    Wait(ItemsReady)   // decrement semaphore

    Get item j;        // remove item from buffer
    Signal(SpacesLeft) // increment semaphore
    Consume item i;   // consume item
}
```

Consumer class

If this solution uses semaphores correctly, then

$$N = \text{SpacesLeft} + \text{ItemsReady}$$

as the producer will always be placing items in to the buffer when there are spaces available in the buffer.

However, this solution does not consider situations in which there are multiple producers and/or multiple consumers.

The multiple producer-consumer problem

Problem description

The multiple producer-consumer problem is a classical inter-process communication problem in which:

- multiple producer repeatedly produces items and places them in to a buffer; and
- multiple consumer consumes the items one-by-one by taking them from the buffer.

As with the previous producer-consumer problem, this problem has the following requirements:

- the buffer must be assumed to be first in, first out (FIFO);
- the producers may produce a new item only at a time when the buffer is not full;
- the consumers may consume an item only at a time when the buffer is not empty; and
- the process terminates when all items produced are eventually consumed.



The problem arises when attempting to devise a method that is able to manage:

- two producers placing items in to the same slot in the buffer; and
- two consumers removing items from the same slot in the buffer.

This is similar to the problem discussed in the printer spooler example (page n).

A race condition may also occur when producers attempt to access a variable at the same time.

To demonstrate the race condition, it is necessary to consider the following possible interleaving of the threads/processes:

- two producers access the *SpacesLeft* variable at the same time, which corresponds to decrementing the semaphore;
- both producers get the same next empty slot in the buffer at the same time; and
- both producers write in to the same slot.

This example shows that the race condition has caused the data stored in the buffer slot by the first producer to be lost as it has been overwritten by the data stored in the buffer slot by the second producer.

In order to ensure mutual exclusion when multiple users are involved, an additional semaphore must be introduced.

Mutex

A **mutex** (or **binary semaphore**) is a semaphore with ownership that can only be released by its owner and is initially set to 1.

Problem solution

It is now possible to construct a solution, using a mutex (or binary semaphore), that will ensure mutual exclusion even when there are multiple producers and/or multiple consumers.

It is assumed that

```
ItemsReady = 0
SpacesLeft = N //size of buffer
```

```
LOOP
{
    Produce item i      // produce item
    Wait(SpacesLeft)   // decrement semaphore
    Wait(BusyBuffer)   // mutex

    Put item i;        // place item in to buffer
    Signal(BusyBuffer) // release mutex
    Signal(ItemsReady) // increment
}
```

Producer class

```
LOOP
{
    Wait(ItemsReady)   // decrement semaphore
    Wait(BusyBuffer)   // mutex

    Get item j;        // remove item from buffer
    Signal(SpacesLeft) // increment semaphore
    Signal(BusyBuffer) // release mutex
    Consume item i;   // consume item
}
```

Consumer class

The mutex *BusyBuffer* has ownership and therefore can only be incremented/decremented by the same thread/process.

The order in which semaphores are incremented and decremented is essential. This can be demonstrated by inspecting the effect of switching around two statements in the Consumer class:

```
Wait(ItemsReady) //decrement semaphore
Wait(BusyBuffer) //mutex
...
Wait(BusyBuffer) //mutex
Wait(ItemsReady) //decrement semaphore
```

This switching would cause ...

4.4 The dining philosophers problem

Problem description

The dining philosophers problem is a classical inter-process communication problem in which:

- five philosophers are seated around a circular table eating and thinking; and
- each philosopher has a plate of spaghetti that they can eat with forks.

This problem has the following requirements:

- the life of a philosopher consists of only alternating periods of eating and thinking;
- between each pair of plates is one fork;
- each philosopher needs two forks to eat the spaghetti; and
- no two philosophers may hold the same fork simultaneously.



Figure 4.9: Layout of table

The problem arises when attempting to devise a method that is able to:

- allow each philosopher to have alternating periods of eating and thinking; and
- not result in a deadlock.

It could be said that the problem requirement for each philosopher to need two forks to eat the spaghetti is somewhat artificial. As a result, we can substitute the spaghetti for rice and substitute chopsticks for forks.

The problem now has the following updated requirements:

- the life of a philosopher consists of only alternating periods of eating and thinking;
- between each pair of plates is one chopstick;
- each philosopher needs two chopsticks to eat the rice; and
- no two philosophers may hold the same chopstick simultaneously.



Figure 4.10: Updated layout of table

Problem solutions

Problem solution 1

This problem can be solved using semaphores, using the following assumptions:

- each philosopher is a different thread with a unique ID;
- one semaphore is implemented per chopstick; and
- chopsticks are identified by using the unique ID of a philosopher.

As chopsticks are identified by using the unique ID of a philosopher, it could be that:

- the chopstick to the left of a given philosopher is $\text{chopstick}[i]$; and
- the chopstick to the right of a given philosopher is $\text{chopstick}[(i - 1) + N] \% N$.

where i is the ID of the philosopher and N is the number of philosophers.

The identification of the chopsticks works as demonstrated by using Philosopher 1 as an example.



The diagram shows that:

- the chopstick to the left of Philosopher 1 is *chopstick[1]*; and
- the chopstick to the right of Philosopher 1 is *chopstick[0]*.

Using this example, it is possible to validate the chopstick identification formulas discussed before.

Given that $i = 1$ for Philosopher 1 and that $N = 5$ as there are five philosophers in total,

```

left = chopstick[i]
left = chopstick[1]
right = chopstick[(i-1) + N] % N
right = chopstick[(1-1) + 5] % 5
right = chopstick[(0 + 5) % 5]
right = chopstick[5 % 5]
right = chopstick[0]

```

This shows that the chopstick identification formulas work for Philosopher 1.

```

LOOP
{
    think();
    wait(chopstick[left]);      // take left chopstick
    wait(chopstick[right]);    // take right chopstick
    eat();
    signal(chopstick[left]);   // release left chopstick
    signal(chopstick[right]); // release right chopstick
}

```

Producer class

However, this solution has the possibility to cause a race condition.

This can be demonstrated in the situation where all five philosophers wish to eat at the same time and therefore all take their left chopsticks:

- Philosopher 0 takes *chopstick[left]*;
- Philosopher 1 takes *chopstick[left]*;
- Philosopher 2 takes *chopstick[left]*;
- Philosopher 3 takes *chopstick[left]*; and

- Philosopher 4 takes *chopstick[left]*.

This causes a race condition as:

- each philosopher will now be waiting to take *chopstick[right]*;
- no philosopher can take *chopstick[right]* as this chopstick is the subsequent philosopher's *chopstick[left]* and this has already been taken; and
- no philosopher can perform their *eat* function and therefore no chopsticks will be released as the *signal* functions are only performed after the *eat* function has completed.

This shows that, if the threads in the solution are interleaved, a race condition may occur. In this case, a circular wait is caused and therefore there is a deadlock.

Possible solution 2

As shown in the multiple producer-consumer problem, a mutex (or binary semaphore) can be introduced to prevent this deadlock.

```
LOOP
{
    think();
    wait(busy);           // mutex
    wait(chopstick[left]); // take left chopstick
    wait(chopstick[right]); // take right chopstick
    eat();
    signal(chopstick[left]); // release left chopstick
    signal(chopstick[right]); // release right chopstick
    signal(busy);          // release mutex
}
```

Producer class

This solves the deadlock as the mutex (or binary semaphore) signals the critical section/region and prevents other philosophers from attempting to take a chopstick that is currently in use by another philosopher.

Although this solution is deadlock-free, it has a performance bug. The mutex means that only one philosopher can be eating at any instant and, with five chopsticks available, it should be possible to allow two philosophers to eat at the same time. As a result, there are more efficient solutions to this problem that achieve maximum parallelism.

Final revision solution

This solution is deadlock-free and allows the maximum parallelism for an arbitrary number of philosophers.

```
#define N      5           // number of philosophers
#define LEFT    (i + N - 1) % N // ID of i's left neighbour
#define RIGHT   (i + 1) % N   // ID of i's right neighbour
#define THINKING 0          // philosopher is thinking
#define HUNGRY   1          // philosopher is trying to acquire
                           // chopsticks
```

```

#define EATING    2           // philosopher is eating

typedef int semaphore        // semaphores are a special kind of
                             // integer
int state[N];               // array to keep track of philosopher'
                             // state
semaphore mutex = 1;        // mutual exclusion for critical
                             // regions
semaphore s[N];             // one semaphore per philosopher

// i: philosopher unique ID, from 0 to N-1
void philosopher (int i)
{
    while (TRUE)           // repeat indefinitely
    {
        think();            // philosopher is thinking
        take_chopsticks(i); // acquire two chopsticks or be blocked
        eat();                // philosopher is eating
        put_chopsticks(i);   // place both chopsticks back on the
                             // table
    }
}

// i: philosopher unique ID, from 0 to N-1
void take_chopsticks(int i)
{
    wait(&mutex);          // enter critical section/region
    state[i] = HUNGRY;       // record fact that philosopher i is
                             // hungry
    test(i);                // attempt to acquire two chopsticks
    signal(&mutex);         // exit critical section/region
    wait(&s[i]);            // block if chopsticks were not
                             // acquired
}

// i: philosopher unique ID, from 0 to N-1
void put_chopsticks(int i)
{
    wait(&mutex);          // enter critical section/region
    state[i] = THINKING;    // record fact that philosopher i is
                             // thinking
    test(LEFT);              // check if left neighbour can now eat
    test(RIGHT);             // check if right neighbour can now eat
    signal(&mutex);          // exit critical section/region
}

// i: philosopher unique ID, from 0 to N-1
void test(int i)
{
    if (state[i] == HUNGRY && state[LEFT] != EATING

```

```
&& state[RIGHT] != EATING)
{
    state[i] = EATING;           // record fact that philosopher i is
    hungry
    signal(&s[i]);            // exit critical section/region
}
}
```

Producer class

This solution uses the array *state* to keep track of whether a philosopher is:

- eating;
- thinking; or
- hungry (trying to acquire chopsticks).

This represents an array of semaphores and each philosopher has its own *state* array. This enables hungry philosophers to be blocked if the required chopsticks are currently busy.

A philosopher may move in to eating state only if neither neighbouring philosopher is eating. For example, Philosopher 1 cannot enter eating state if:

- Philosopher 0 is currently in eating state; or
- Philosopher 2 is currently in eating state.

5. Memory Management



5.1	Memory hierarchy	70
5.2	Primary memory	71
5.3	Secondary memory	74
5.4	Read-only memory (ROM)	80
5.5	Concept and aim of memory management	81
5.6	Overview of memory abstraction	82
5.7	Memory abstraction - physical address	83
5.8	Memory abstraction - logical address	85
5.9	Memory Management Unit (MMU)	86
5.10	Memory partitioning	88
5.11	Swapping	91
5.12	Managing free space	92

5.1 Memory hierarchy

Definition

The **memory hierarchy** separates computer storage into a hierarchy based on response time. A computer system is usually composed by a layered memory system

Diagram



Higher layers correspond to faster devices that:

- have lower capacity;
- require more power and generate more heat; and
- are more expensive.

Lower layers correspond to slower devices that:

- have higher capacity;
- require less power and generate less heat; and
- are cheaper.

The temporary storage areas can be described as volatile; the data does not persist after power-down.

The permanent storage areas can be described as non-volatile; the data persists after power-down.

5.2 Primary memory

Definition

Primary memory is volatile, meaning that the data contained within does not persist after power-down, and stores data and instructions for use by the CPU in currently running processes.

Registers

Size	Very small.
Speed	Extremely fast.
Purpose	Individual containers for single values. Almost all computers load data from a memory lower in the hierarchy into registers, where it is used for arithmetic operations and is manipulated or tested by machine instructions.
Location	Inside the CPU.

Registers are typically referred to by “name” (i.e. individual identifying code), not by address, as happens with other types of memory.



There are two main types of registers:

- generic registers – allow the CPU to temporarily store data on which it will perform operations, and consequently data that results from those operations; and
- specialised registers – manipulated in mostly the same fashion as generic registers but are used for specialised purposes.

Specialised registers include:

- Instruction Register (IR) – holds the instruction currently being executed;
- Memory Data Register (MDR) – holds the piece of data that has been fetched from memory;
- Memory Address Register (MAR) – holds the address of the next piece of memory to be fetched;

- Program Counter (PC) – holds the location of the next instruction to be fetched from memory and is automatically incremented between supplying the address of the next instruction and the instruction being executed; and
- Accumulator (ACCU) – used as the default location to store any calculations performed by the arithmetic and logic unit.

Generic registers are available to store any transient data required by the program.

For example, when a program is interrupted its state, the values stored in the specialised registers may be saved into the generic registers, ready for recall when the program is ready to start again.

In general, the more registers a CPU has available, the faster it can work.

Cache memory

Size	Small (KB/MB).
Speed	Fast.
Purpose	Serves as an intermediary between main memory and the registers. When data or instructions are fetched from main memory, they are copied to the cache for further use, and therefore, reduce the average cost (time or energy) to access data from the main memory.
Location	Inside the CPU.

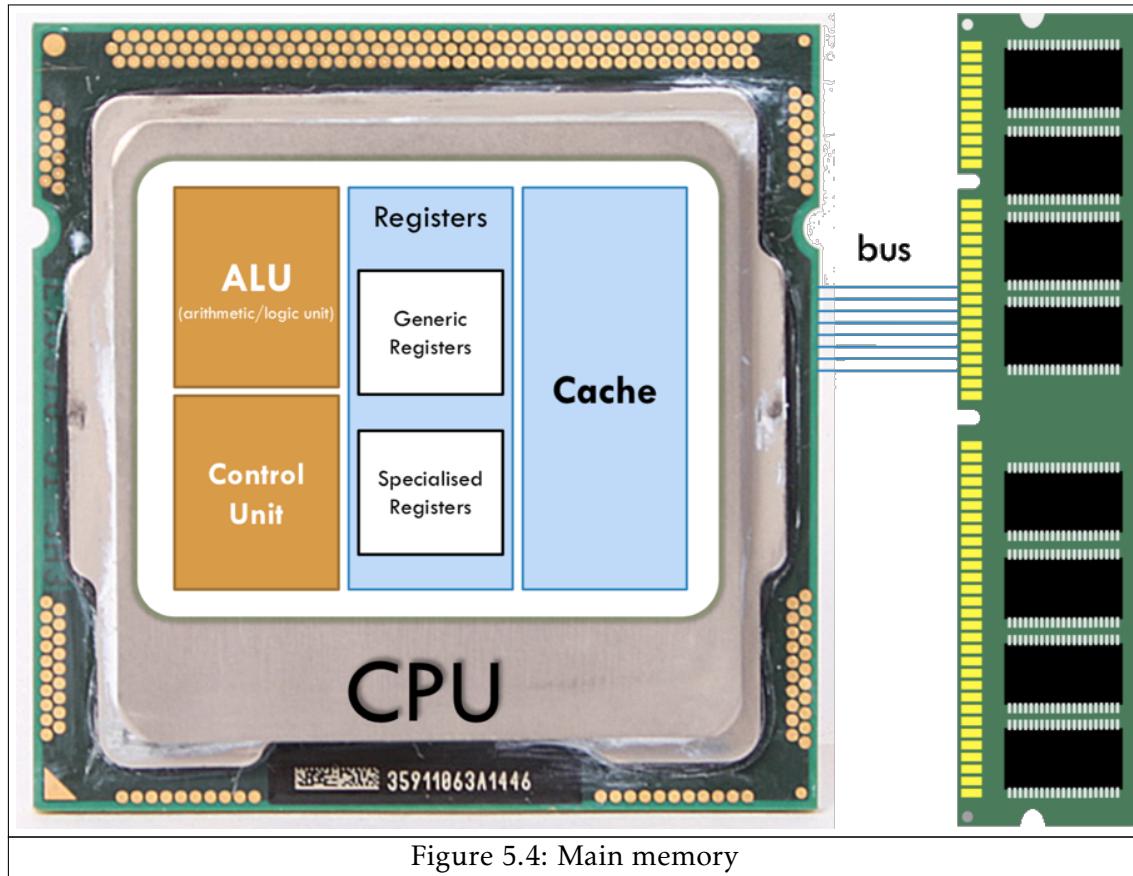


Modern computers typically have a sub-hierarchy of cache memories (L1, L2, L3, L4, etc.), for example:

- level 1 cache (L1) – very fast and small size (2KB - 64KB); and
- level 2 cache (L2) – fast and medium size (256KB – 2MB).

Main memory

Size	Medium (MB/GB).
Speed	Fast.
Purpose	Stores data and machine code currently being used.
Location	Motherboard, connected to the CPU by a bus.



Main memory is typically a random-access memory (RAM) device. These devices allow data items to be read from or written to in almost the same amount of time, irrespective of the physical location of the data.

5.3 Secondary memory

Definition

Secondary memory is a non-volatile and persistent store of data, which is used for items such as user files, programs and the operating system.

Magnetic – Hard-Disk Drive (HDD)

Description

A **Hard-Disk Drive (HDD)** contains rotating platters which are coated with magnetic material. This contains ferrous particles, containing iron, which are polarised to become either a north or south state; these states can correspond to either a 0 or a 1.

The platters spin at speeds of up to 10,000 revolutions per minute (RPM) and a drive head moves across the platters to access different tracks and sectors. Data is read from or written to the platters as it passes under the drive head. When no operations are being performed, the drive head is parked to the side of the platters to prevent damage from movement.

There may be several platters and several drive heads.

Size	xxx'GB and x'TB.
Speed	Medium.
Purpose	Consumer desktop computer system secondary storage.
Reliability	Prone to drop damage as it has moving parts and magnetic fields.

Evaluation

Advantages	Disadvantages
Large storage capacity.	Relatively slow read and write speeds because the disk head must physically move to the location of the data on the platter.
Affordable large capacities, HDDs are available in large storage capacities at relatively low prices.	Can be loud because there are moving parts, such as the platters and disk head, which generate noise.
	Low reliability because they are prone to drop damage, due to their moving parts, magnetic fields.
	Relatively large form factor since there must be room to enclose the platters and disk heads.

Magnetic – Magnetic tape

Description

A **magnetic tape** is a reel of tape which is coated with a magnetic material. This contains ferrous particles, containing iron, which are polarised to become either a north or

south state; these states can represent either a 1 or a 0.

The data is accessed using serial access; this is where a tape reader starts at the beginning of the tape and “fast forwards” to the portion of the tape containing the required data.

Size	xxx'GB and x'TB.
Speed	Slow.
Purpose	Server backup storage.
Reliability	Prone to magnetic fields.

Evaluation

Advantages	Disadvantages
Large storage capacity.	Relatively slow read and write speeds because it uses serial data access.
Affordable large capacities, magnetic tapes are available in high capacities at relatively low prices	Low reliability because they are prone to magnetic fields.
	Requires specialist equipment to read the data on the tape, known as a tape reader.

Flash – Solid-State Drive (SSD)

Description

A **Solid-State Drive (SSD)** contains millions of NAND flash memory cells and a controller which manages pages and blocks of memory.

Each NAND flash memory cell delivers a current along the bit and word lines to activate the flow of electrons from the source towards the drain. The current on the word line is strong enough to push a few electrons across an insulated oxide layer into a floating gate. The charge in the flowing gate is measured: if there is some charge, where there are some electrons flowing, a 1 is represented; and if there is no charge, where there are no electrons flowing, a 0 is represented.

Size	xx'GB and xxx'GB.
Speed	Very fast.
Purpose	Portable device secondary storage, such as smartphones and laptops.
Reliability	Less prone to drop damage as there are no moving parts, limited write cycles before damage.

Flash drives and SD cards also use solid-state technologies. They are often used for portable storage; for example, flash drives may be used for transferring documents and SD cards may be used to storage images in a camera.

Evaluation

Advantages	Disadvantages
Fast read and write speeds because the data can be accessed directly on the drive using electrical currents.	Expensive large capacities, SSDs are more expensive than other storage options.
Available in different form factors, such as USB memory sticks and SD cards.	Low durability because they have a finite number of read and write operations which can be performed before the performance of the drive begins to downgrade due to the discharge of its electrons.
High reliability because there are no moving parts.	

Silent operation because there are no moving parts.	
--	--

Optical – Optical disk

Description – Optical disk

An **optical disk** works by using:

- a high-powered laser to change the properties of its surface to make them less reflective, by a process called “burning”; and
- a low-powered laser to read the disk by shining light onto the surface and a sensor is used to measure the amount of light that is reflected back.

There are multiple types of optical disks, including:

- CD-ROM;
- CD-RW;
- DVD-RW; and
- Bluray.

Description – CD-ROM

A **CD-ROM (read only Compact Disk)** is pressed during its manufacture and has pits, where the surface has been depressed, and lands, where the surface has not been depressed.

When the low-powered laser is shone on the surface:

- the start or end of a pit, the light is scattered and not reflected very well — this is used to represent a 0; and
- a land, the light is not scattered and reflected normally -- this is used to represent a 1.

Size	xxx'MB.
Speed	Slow.
Purpose	Commercial music distribution.
Reliability	Prone to scratches and cracks.

Description – CD-RW

A **CD-RW (re-writable Compact Disk)** uses a laser and a magnet to heat a spot on the disk. When the spot is heated the magnetic orientation is changed to represent a 0 or a 1 before the magnet is cooled.

Size	xxx'MB.
Speed	Slow.
Purpose	Small file storage such as documents.

Reliability	Prone to scratches and cracks.
--------------------	--------------------------------

Description – DVD-RW

A **DVD-RW** (**re-writable Digital Versatile Disc**) uses a phase changed alloy that can change between amorphous and crystalline states by changing the power of the laser beam.

Size	x'GB.
Speed	Slow.
Purpose	Standard definition videos.
Reliability	Prone to scratches and cracks.

Description – Blu-Ray

A **Blu-Ray** disk has a higher capacity than a CD-ROM because the laser used has a shorter wavelength which creates smaller pits and lands and therefore, a greater number of pits and lands can be present on the same surface area.

Size	xx'GB.
Speed	Slow.
Purpose	High definition videos.
Reliability	Prone to scratches and cracks.

Evaluation

Advantages	Disadvantages
Relatively portable because they can be carried in a disk case.	Relatively slow read and write speeds because the disk must spin in the optical drive to read the data.
Immune to magnetic fields meaning that the close proximity of a magnet is not going to corrupt data.	Low reliability because they are prone to scratches and cracks.

5.4 Read-only memory (ROM)

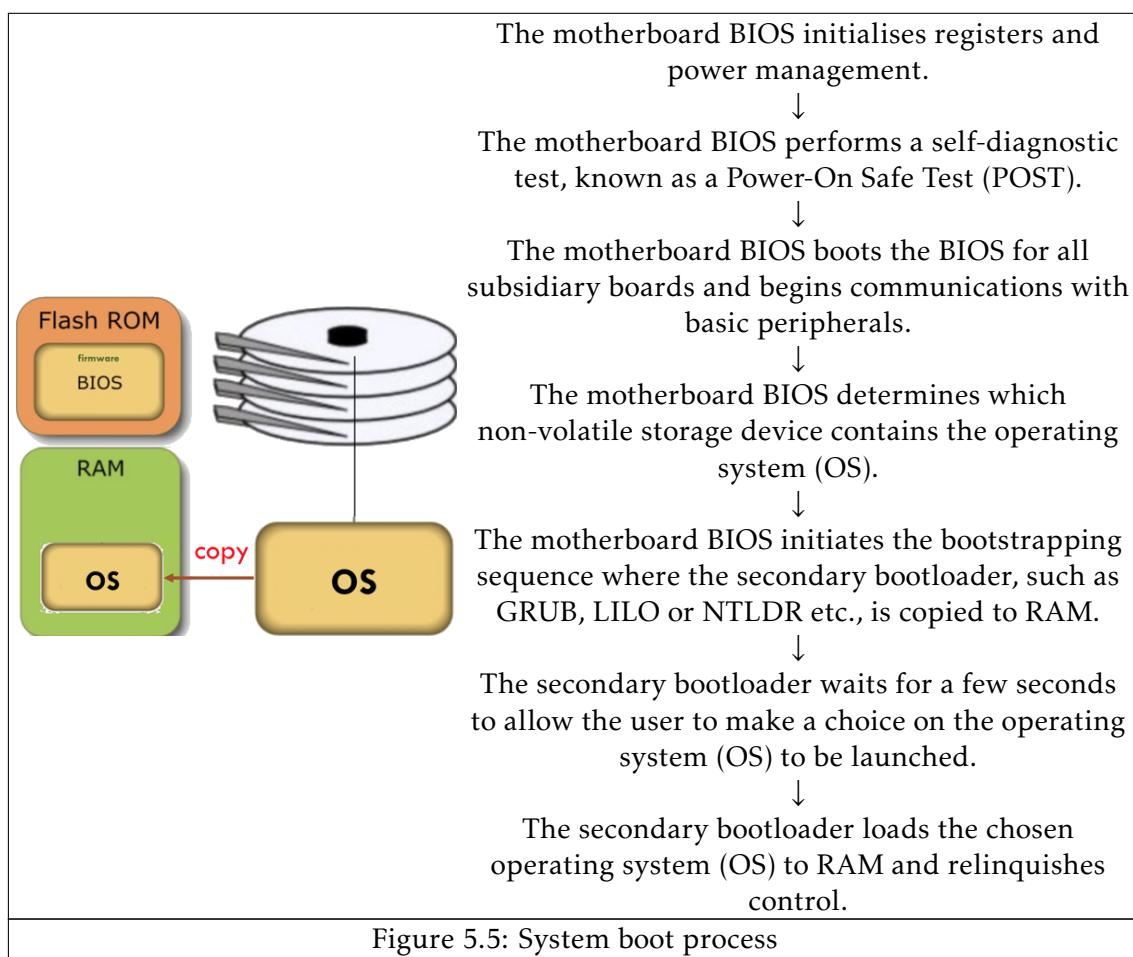
Properties

Size	Medium (MB/GB).
Speed	Fast.
Purpose	Used in the system boot process and, in some cases, stores the operating system in computer systems where the software is not updated.
Reliability	Motherboard.

Use in the system boot process

A **Basic Input/Output System (BIOS)** is a set of instructions which control the boot process of a computer system and the communication between the operating system and peripherals.

A **peripheral** is a piece of hardware outside a computer system that is not required for the computer system to operate.

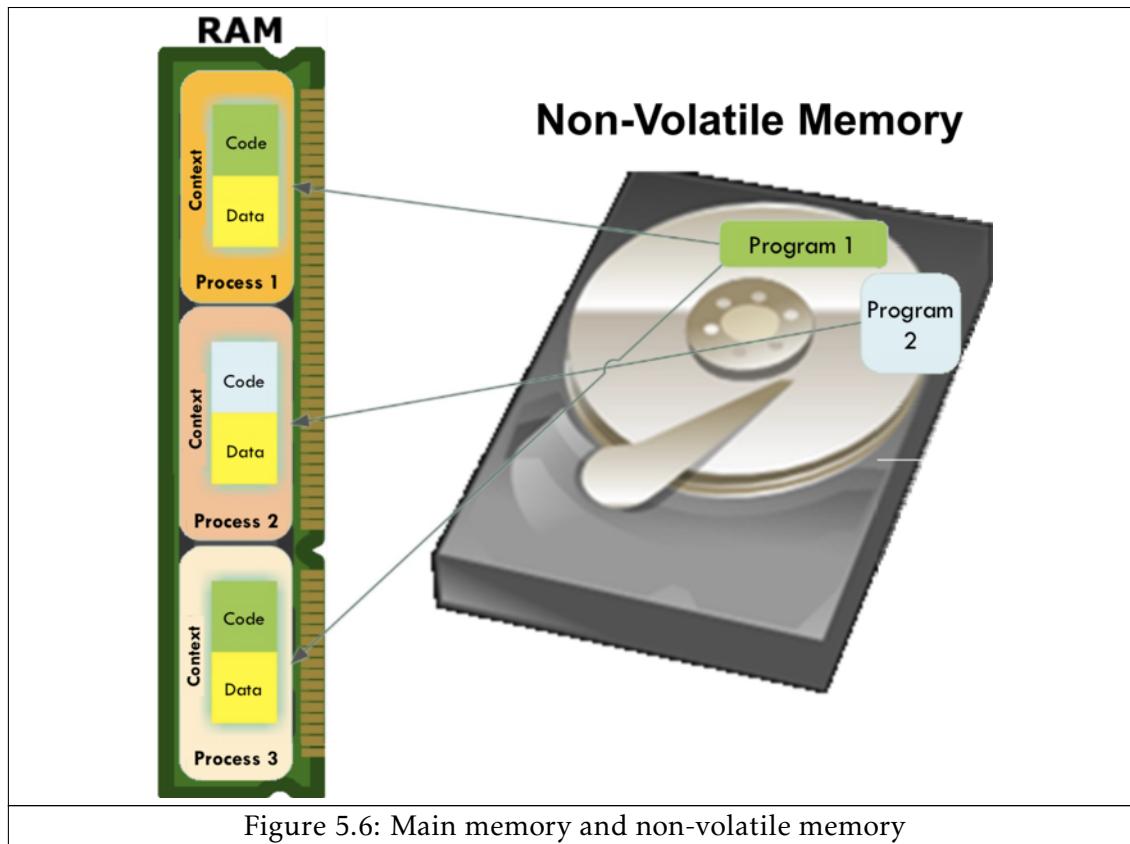


5.5 Concept and aim of memory management

Concept

Processes share physical memory, as well as the CPU. Main memory (RAM) is one of the most important resources in a computer system as:

- data and instructions are fetched from RAM in to the CPU for execution; and
- memory can store several running processes at once.



As seen in the diagram above, RAM can contain the data and instructions for multiple processes. These are loaded in to RAM from non-volatile memory, such as a hard-disk drive (HDD).

Aim

Aims of memory management include:

- providing the memory space for concurrent processes;
- taking advantage of the locality of reference, which is the tendency of a CPU to access the same set of memory locations repetitively over a short period of time, by estimating the data and/or instructions to be used next and transfer them to the most effective place in the memory hierarchy pyramid;
- protect processes from one another;

- relocated memory to new processes; and
- make the addressing of memory space transparent by providing memory abstraction.

5.6 Overview of memory abstraction

Definition

Memory abstraction provides an abstraction layer between the program execution and the memory that provides a different "view" of a memory location depending on the execution context in which the memory access is made. This is achieved by creating an "abstract memory" to allow for co-existence in physical memory.

Computer system with no memory abstraction

Early mainframe computers had no memory abstraction and every program had access to physical memory.

The memory model presented to the programmer was a set of addresses from zero (0) to max, in which each address corresponds to a cell containing some number of bits.



The diagram above shows three different ways of organising memory with an operating system and one user process.

Model	Organisation	Uses
(a)	The operating system may be at the bottom of memory in RAM.	This model was formerly used on mainframes and minicomputers but is rarely used any more.
(b)	The operating system may be in ROM at the top of memory.	This model is used on some handheld computers and embedded systems.

(c)	The device drivers may be at the top of memory in ROM and the operating system may be in RAM at the bottom of memory.	This model was formerly used by early personal computers, such as those running MS-DOS, where part of the operating system (OS) was stored in BIOS.
-----	---	---

These models refer to fixed memory addresses and therefore prevent more than one process running at a time.

5.7 Memory abstraction - physical address

Definition

An **address space** is a set of addresses that a process can use to reference memory.

A **physical address** (or **memory address**) identifies a physical location of required data in a memory.

Logical address space is the set of all physical addresses corresponding to the logical addresses in a logical address space.

Referencing physical memory

It is possible to reference memory by using the physical address locations.



Drawbacks

However, this is not ideal as:

- user programs can address physical memory directly and trash the OS intentionally/accidentally, i.e. models (a) and (c); and
- it is difficult to have multiple programs running at once (via context switch).

For these reasons, users should not be able to reference physical memory directly.

5.8 Memory abstraction - logical address

Definition

An **address space** is a set of addresses that a process can use to reference memory.

A **logical address** (or **program address**) is generated by the CPU while a program is running. The logical address is virtual address as it does not exist physically, therefore, it is also known as **virtual address**. This address is used as a reference to access the physical memory location by CPU.

Logical address space is the set of all logical addresses generated by a program's perspective.

Using logical address spaces

Programmers can refer to their own address space, the logical address space for a given program.

It is possible for the same logical address, as seen by two different processes, to correspond to different physical addresses. It is important that a distinction is made between the two physical addresses.



When a process is loading, the **loader** is responsible for transferring the program from mass storage to main memory for execution.



It is not possible to determine where the program is loaded in physical memory as its location depends on the current contents of memory. As a result, the program's code cannot refer to fixed memory addresses.

As discovered before, machine instructions should not directly access the physical address space and therefore logical addresses are used. These must be converted to physical addresses by the Memory Management Unit (MMU).

5.9 Memory Management Unit (MMU)

Definition

The **Memory Management Unit (MMU)** is responsible for the logical-physical conversion. This is achieved by using registers to record the location of the partition in order to map the logical addresses to the correct physical addresses.

The ultimate aim for the MMU is to perform efficiently such that the CPU does not have to access memory past the cache level. This does not always happen but, in an ideal world, the CPU would only be accessing the faster memory, its registers and cache, as to prevent bottlenecks.

Base and limit registers

The **base register** stores the start address of the partition.

The **limit register** stores the length of the partition.

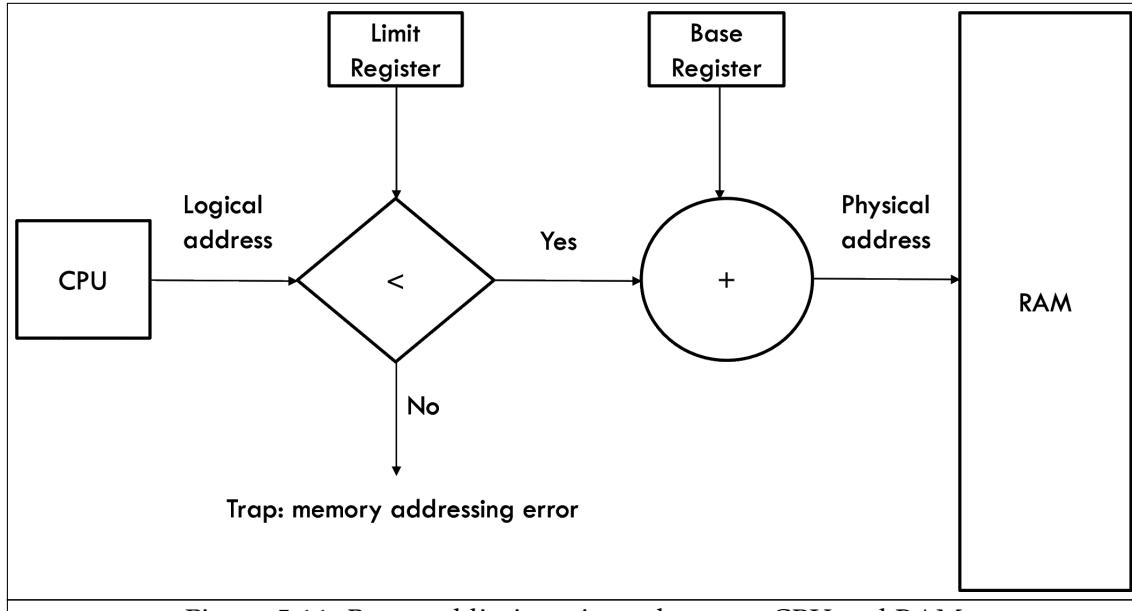


Figure 5.11: Base and limit registers between CPU and RAM

Memory protection and relocation

Definitions

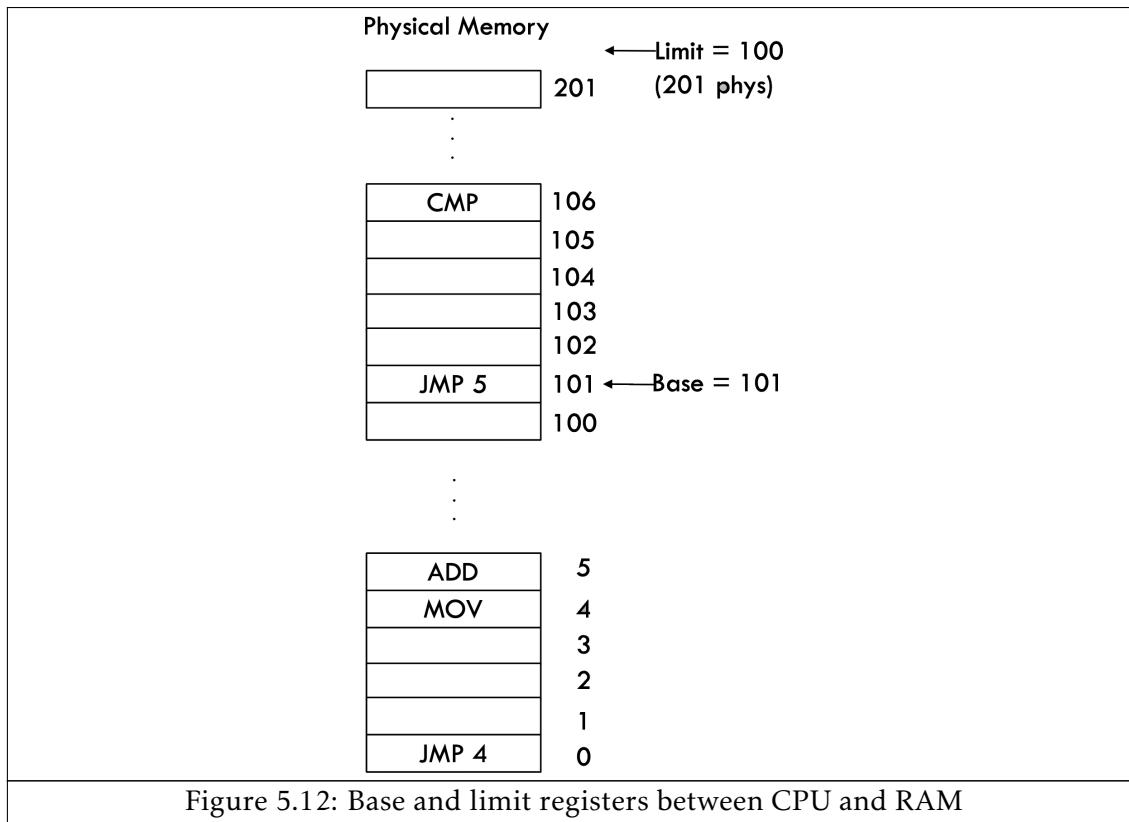
- **Relocation** – The base register value is added to the logical address in order to map the logical address to a physical address.
- **Protection** – Logical addresses that are greater than the limit register value should point to an invalid memory location. This check ensures that the logical address is within the range of the partition in order to ensure that a process is only accessing its own logical address space.

How it works

The address part of a given machine instruction is used as an offset from the base address.

For example, a machine instruction may have

- a base address of 101; and
- the instruction JMP.



The operating system (OS) would be responsible for allowing the process to access the memory location desired by the process after the JMP instruction.

The JMP 5 shows a logical address of 5, while the physical address has a base of 101 and therefore the resulting physical address is 106.

In this scenario, if the machine instruction contained JMP 100 instead, there would be an invalid memory location error. This is because the JMP 100 instruction would result in a physical address of 201. This violates the logical address limit of 100 and the physical address limit of 201 and therefore this error must be thrown for protection.

5.10 Memory partitioning

Definition

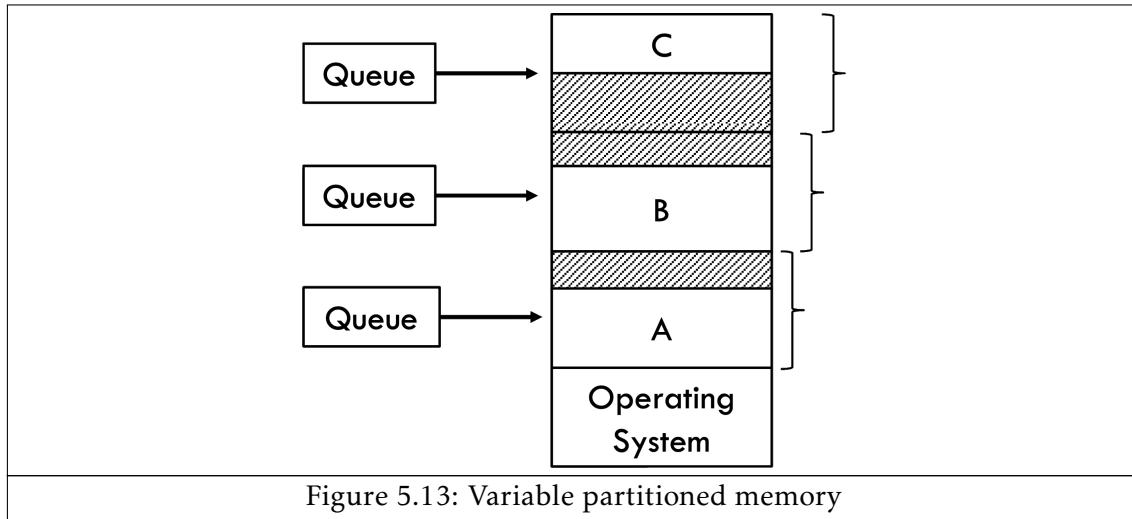
Memory partitioning is the system by which the memory of a computer system is divided into sections for use by the resident programs.

This may be achieved by fixed partitioned memory or variable partitioned memory.

Fixed partitioned memory

How it works

Memory is divided into several partitions of fixed sizes. The partitions may be different sizes from one another, but remain at the same size once created.



As shown above, each process (A, B and C) has its own partition.

Evaluation

Advantages	Disadvantages
Multiprogramming , as it is compatible with process co-existence in memory.	Restricted as there is a fixed number of possible simultaneous processes.
Easy to implement .	Internal fragmentation , as each partition has unused space.
Allows for fast context switching .	Inefficient space usage due to internal fragmentation.
	Issues may occur if a process is too large for any partition.

Variable partitioned memory

How it works

Each process is allocated an exact size of memory space. Processes are loaded into contiguous memory slots until memory is full.



As shown above, when processes A and B terminate, their memory space is deallocated.

Evaluation

Advantages	Disadvantages
No internal fragmentation.	Assumes that the memory manager knows how much memory a process requires.
Efficient space usage as when a process terminates, the space it once occupied in memory is freed up.	External fragmentation may occur, where "holes" outside partitions are introduced.
Adjacent fragmentation can be merged.	
The number of parallel processes is not fixed.	

Compaction

Reducing external fragmentation is achieved by compaction.

This is a process whereby processes are physically moved to close the "holes" created by the deallocation of memory that once belonged to processes.

However, this process has a heavy overhead causing slowdowns and therefore, is not commonly performed.

Dealing with processes of varying size

So far, it is been assumed that processes have a fixed size and that the operating system (OS) allocates that memory.

However, it is possible that some processes may try to grow their memory allocation, such as increasing the size of an array.



As shown in the diagram above, if a hole is adjacent to the process, it can be allocated to the process to allow the process to grow in to that hole.

Otherwise, it is necessary to allocate more memory in a different physical location. If memory is full, swapping may occur or the process needs to be suspended.

5.11 Swapping

Definition

Swapping is a memory reclamation method wherein memory contents not currently in use are swapped to a disk to make the memory available for other applications or processes.

Why is swapping required?

Memory is not an infinite resource, as shown in the memory hierarchy (page n). In practice, the total amount of RAM required by all running processes is often larger than the amount available in the computer system.

Idle processes are typically stored on secondary memory as to provide access to more primary memory for processes that are currently in use.

Memory allocations change as processes are loaded and unloaded from memory.

How it works

Swapping shows a simple strategy of loading each process in to memory in its entirety, executing the process for an amount of time and then returning it to secondary memory. This is a similar concept to context switching.



However, swapping adds a layer of complexity as the operating system must now manage free space.

5.12 Managing free space

When memory is assigned dynamically via swapping, the operating system must manage it. In general terms, there are two ways to keep track of memory usage: bitmaps and free lists.

Bitmaps

How it works

Memory is divided into allocation units as small as a few words and as large as several kilobytes.



Figure 5.17: Using bitmaps to manage free space

As shown above, corresponding to each allocation unit is a bit in the bitmap, which is:

- zero (0) if the unit is free; and
- one (1) if the unit is occupied.

Evaluation

Advantages	Disadvantages
Fast as the bitmap array can be accessed directly and the status of the allocation unit can be determined.	Maintenance overhead as the status of each allocation unit must be maintained in the bitmap.

Free lists

How it works

A list of processes, denoted as P, and holes, denoted as H, are kept in a list.



Figure 5.18: Using free lists to manage free space

In the free list example shown above, there is the following sequence of processes P and holes H in memory:

- a process (P) starting at 0 and occupying 5 spaces;
- a hole (H) starting at 5 and occupying 3 spaces;
- a process (P) starting at 8 and occupying 6 spaces;
- a process (P) starting at 14 and occupying 4 spaces;
- a hole (H) starting at 18 and occupying 2 spaces;
- a process (P) starting at 20 and occupying 6 spaces;
- a process (P) starting at 26 and occupying 3 spaces; and
- a hole (H) starting at 29 and occupying 3 spaces.

Several algorithms, known as placement policies, can be used to allocate memory for new processes.

Evaluation

Advantages	Disadvantages
More compact than using a bitmap as, for example, a sequence of occupied units can be expressed as [P, 0, 5] rather than 11111.	More complex as it is more difficult to manage and placement policies are required.

Placement policies

First fit

In first fit, the memory manager scans along the list of segments until it finds a hole that is large enough.

The hole is then broken up into two pieces, one for the process and one for the unused memory, except in the statistically unlikely case of an exact fit.



Figure 5.19: Loading a process using first fit with size 2KB.



Figure 5.20: Loading a process using first fit with size 4KB.

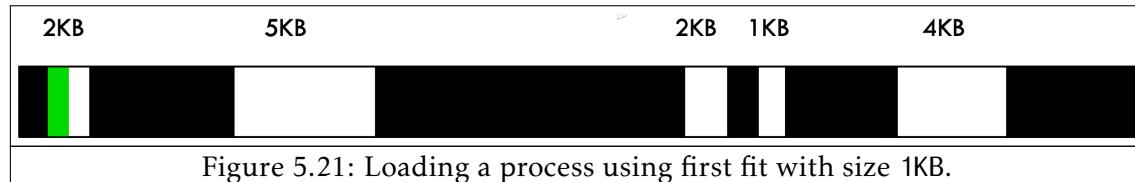


Figure 5.21: Loading a process using first fit with size 1KB.

Advantages	Disadvantages
Fast because it searches as little as possible.	May generate more holes

Next fit

Next fit is a variation of the first fit policy. It works the same way as first fit, except that it keeps track of where it is whenever it finds a suitable hole.

The next time it is called to find a hole, it starts searching the list from the place where it left off last time, rather than the beginning of the list.

This policy was designed to speed up searching by skipping potential tiny holes that cannot fit a process. However, simulations show that next fit gives slightly worse performance than first fit.



Figure 5.22: Loading a process using next fit with size 2KB.



Figure 5.23: Loading a process using next fit with size 4KB.



Figure 5.24: Loading a process using next fit with size 1KB.

Best fit

In best fit, the memory manager searches the entire list, from beginning to end, and finds the hole that is closest to the actual size required rather than breaking up a larger hole that may be required later.



Figure 5.25: Loading a process using best fit with size 2KB.



Figure 5.26: Loading a process using best fit with size 4KB.



Figure 5.27: Loading a process using best fit with size 1KB.

Advantages	Disadvantages
May reduce memory wastage by using better suited holes for processes.	May increase memory wastage as it may generate many tiny holes that are too small to be used by any process. This could be worse than first fit and next fit as they generate larger holes on average that could be used by other processes.
	Slower than first fit and next fit as it requires a search to be completed on the entire list.

Worst fit

Worst fit is the opposite of the best fit policy. It always finds the largest available hole, so that the new hole will be big enough to be useful.

This policy was designed to reduce the number of tiny and unusable holes generated by the best fit policy. However, simulations show that it is not very successful.

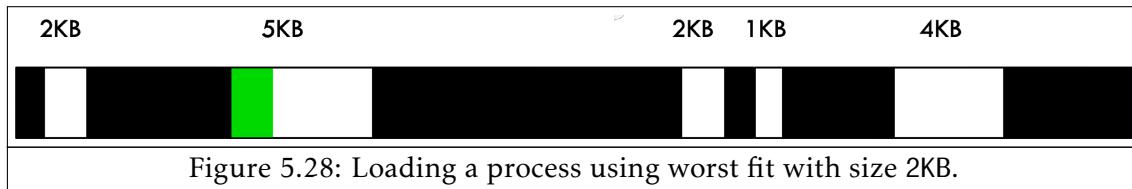


Figure 5.28: Loading a process using worst fit with size 2KB.

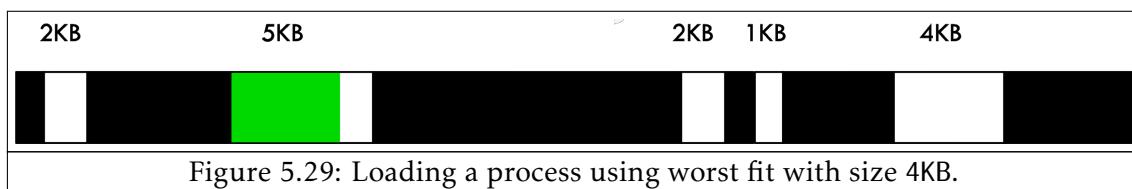


Figure 5.29: Loading a process using worst fit with size 4KB.

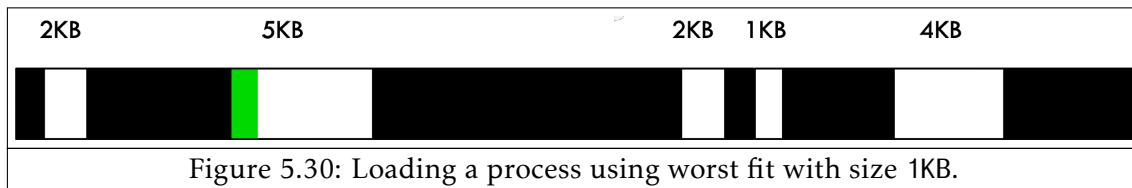


Figure 5.30: Loading a process using worst fit with size 1KB.

Quick fit

In quick fit, the memory manager maintains a table of separate lists for some of the more common sizes requested.

For example, there may be a list for large holes and a list for small holes.

Advantages	Disadvantages
Extremely fast as the lists can be used to quickly find appropriate holes.	Overhead as the lists must be maintained.

6. Virtual memory

6.1	Introduction	99
6.2	Paging	99
6.3	Segmentation	109
6.4	Comparison of paging and segmentation	112
6.5	Paging segments	113

6.1 Introduction

Definition

Virtual memory is a memory management technique that provides an idealised abstraction of the storage resources that are actually available on a given machine.

How it works

Virtual memory involves the use of a partition on a computer system's secondary storage device which acts as a form of main memory for the temporary store of data and instructions used in processing by the CPU; this is required when main memory becomes full and is present to prevent crashing.

Parts of processes are stored on secondary storage and loaded in to main memory when required.

This creates the illusion to users of a very large main memory.

Virtual memory and swapping

While swapping allows multiple processes to run whose total size is larger than overall RAM size, virtual memory additionally allows a single process to run whose size is larger than RAM size.

Virtual memory is implemented using paging and segmentation.

6.2 Paging

Definition

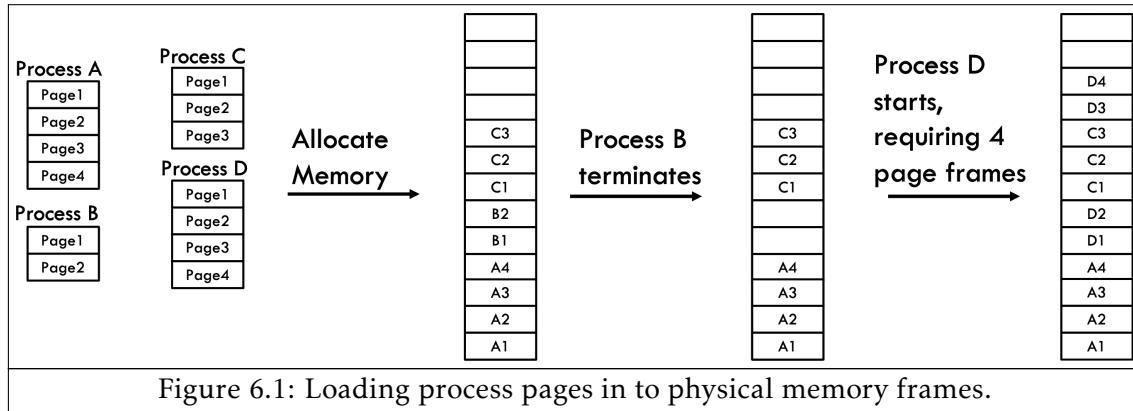
Paging is a memory management mechanism that allows the physical address space of a process to be non-contiguous.

Frames and pages

Physical memory is divided into blocks of equal sizes called **frames**.

In a similar fashion, a process is divided into blocks of the same size, called **pages**.

The pages from the processes are loaded into the available frames in physical memory.



In the diagram above, there are four processes:

- process A which has four pages;
- process B which has two pages;
- process C which has three pages; and
- process D which has four pages.

In the first instance, memory is allocated to processes A-C. Subsequently, process B terminates and its pages are released from the frames.

Later, process D starts. This process has four pages and is therefore split across the two frames available between processes A and C, and after process C.

This is a good solution to address external fragmentation.

Page table

Definition

A **page table** is responsible for mapping virtual pages into page frames when using paging.

Fields

The layout of a page table is highly machine dependent. However, important common fields of a page table entry include:

- page frame – the most important field that outputs the number of the frame;
- present – records whether the page is present in main memory or secondary memory;
- modified (or dirty bit) – records whether the page has been modified since its last loading, if true then it must be copied back to the disk to be saved;
- protection – includes what kind of access is permitted, read/write/execute; and
- referenced – set by the operating system (OS) when the page is used.

Mapping logical program addresses to physical memory addresses

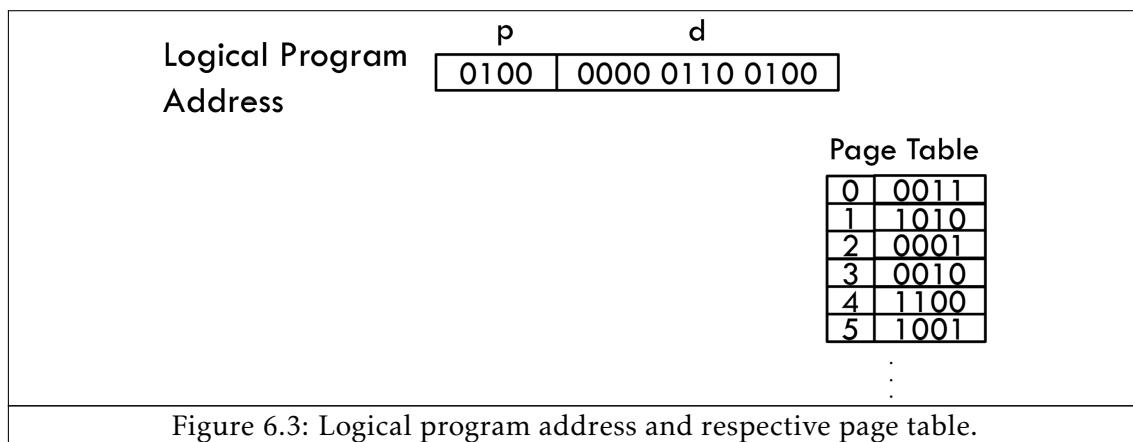
Each logical address is formed as (p, d) where

- p is the process page number; and
- d is the displacement within that page (offset).

p occupies 4 bits and while d occupies 12 bits.



In addition, each process has a page table, which records the memory page frame number for each process page.

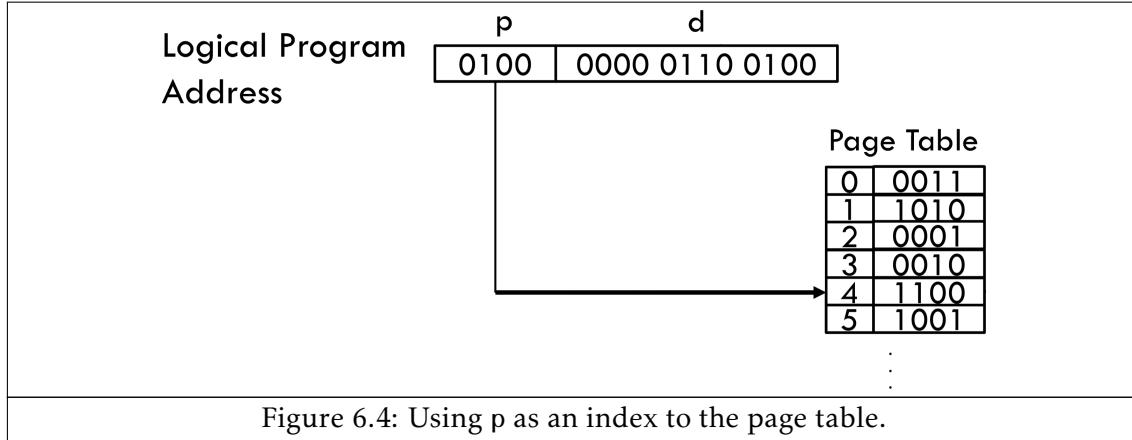


The implementation is performed by the hardware, as such the Memory Management Unit (MMU) is responsible for mapping the logical addresses to physical addresses.

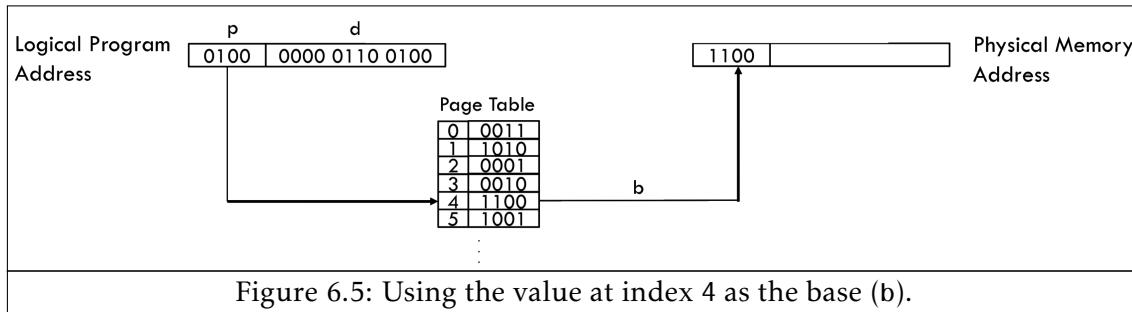
Given a 16-bit program address 0100 0000 0110 0100, it can be deduced that:

$$\begin{aligned} p &= 0100_2 = 4_{10} \\ d &= 0000\ 0110\ 0100_2 = 100_{10} \end{aligned}$$

The process page number (p) provides an index to a location in the page table.



The value at index 4 in the page table can then be used as the base (b) for the physical memory address.



The displacement (d) from the logical program address can then be used to complete the physical memory address.



This can also be shown by performing the following operations:

$$\text{physical memory address} = \text{page table value at index } p + \text{displacement value (d)}$$

$$\text{physical memory address} = 1100\ 0000\ 0000\ 0000 + 0000\ 0110\ 0100$$

$$\text{physical memory address} = 1100\ 0000\ 0110\ 0100_2 = 104_{10}$$

Page faults

Definition

A **page fault** is a type of exception raised by computer hardware when a running program accesses a memory page that is not currently mapped by the memory management unit (MMU) into the virtual address space of a process.

This means that a specific part of an executing process is not in main memory at the moment it is required and therefore the CPU is not able to access this part of the process.

Page fault handling

When a specific part of an executing process is not in main memory when required, a page fault occurs.

↓
The OS reads out hardware registers to determine which virtual address caused the fault.

↓
The OS computes which page is needed and locates that page on disk.

↓
The OS selects an existing available page frame.

↓
The OS checks if that page frame is modified by inspecting the modified field.

↓
If the page frame has been modified, the OS writes the contents of the page frame back to the disk.

↓
The OS fetches a new page from the disk and replaces the old page in the page frame.
This is known as page replacement.

- ↓
The OS updates the mappings and restarts the trapped instruction by:
- marking the virtual page as unmapped by changing the present bit in the page table; and
 - updating the virtual page address with new translation to physical memory.

Figure 6.7: Page fault handling process

When a process exits, the operating system must release its page table, its pages, and the disk space that the pages occupy when they are on disk. If some of the pages are shared with other processes, the pages in memory and on disk can be released only when the last process using them has terminated.

Page replacement

Definition

Page replacement is required during page fault handling when the OS fetches a new page from the disk and replaces the old page in the page frame.

Page replacement algorithms decide which pages to page out, sometimes called swap out, or write to disk, when a page of memory needs to be allocated.

Objective

Page replacement algorithms aim to decide which pages are to be removed from RAM and placed in mass storage such that there are minimal overheads.

This allows thrashing to be avoided, where the CPU is spending more time swapping pages than actually executing a process.

As such, it must be considered that:

- modified pages must first be saved whereas unmodified pages are just overwritten; and
- it is preferable to not replace an often used page in and out of main memory.

These algorithms use information (bits) from the page table.

Optimal page replacement

The **optimal page replacement algorithm** replaces the page that will not be needed for the longest time in the future.

This is not feasible because:

- it is impossible to know the future of a program; and
- it is impossible to know when a given page will be needed next.

However, this is useful for benchmarking page replacement algorithms *a posteriori* (i.e. "what if").

Not recently used (NRU) page replacement

The **not recently used page replace algorithm** uses the reference and modified bits from the page table to collect useful page usage statistics.

Pages are classified based on the contents of their reference and modified bits in the page table.

Class	Reference bit	Modified bit	Meaning
Class 0	0	1	Not referenced and not modified.
Class 1	0	1	Not referenced and modified.
Class 2	1	0	Referenced and not modified.
Class 3	1	1	Referenced and modified.

Figure 6.8: Page classification

Based on the page classifications, the algorithm removes a page from the lowest numbered non-empty class. For example, should pages exist in Class 1, Class 2 and Class 3 but not Class 0, then a page would be removed from Class 1.

A timer interrupt clears the reference bit to distinguish between those pages that have been recently referenced and those which have not been referenced for a given amount of time.

Advantages	Disadvantages
Very low overhead.	Not optimal when compared to the optimal page replacement algorithm.
Easy to implement.	

FIFO page replacement

In the **First In, First Out (FIFO) page replacement algorithm**, main memory maintains a list of all pages currently stored in main memory. In which,

- the most recently arrived page is located at the tail; and
- the least recently arrived page is located at the head.

When page replacement occurs, the page at the head is removed and the new page is added to the tail of the list such that the oldest page is removed.

However, this algorithm is rarely used as the oldest page may still be useful.

Second chance page replacement

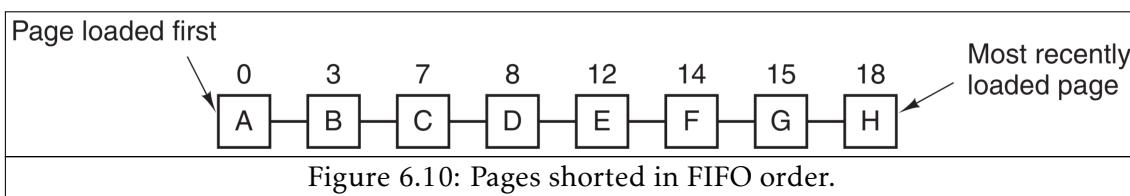
The **second chance page replacement algorithm** is a variation of the FIFO page replacement algorithm. It attempts to avoid the problem of replacing a heavily used page.

The algorithm checks the reference bit of the oldest page first and performs actions based on the page's state deduced from the reference bit.

Reference bit	Page's state	Actions taken
0	Old and unused.	The page is replaced in the list, in the same manner as the FIFO page replacement algorithm.
1	Old but used.	The page is placed at the end of the list of pages, the load time resets as though it has just arrived in memory and the search continues to the next page.

Figure 6.9: Page classification for the second chance page replacement algorithm.

As shown in the figure above, the second chance page replacement algorithm factors in the use of a page as well as its age.



The figure above shows the pages sorted in FIFO order, with the most recently arrived page at the head of the list of pages and the least recently arrived page at the tail of the list of pages.



The figure above shows a "second chance" for page A. If a page fault was to occur at time 20 and page A has a reference bit of 1, the loads times are updated such that it appears that page A has just arrived in memory.

Advantages	Disadvantages
Improvement over the FIFO page replacement algorithm as avoids removing old but useful pages.	Not optimal when compared to the optimal page replacement algorithm as unnecessarily moves pages around the list.

Clock page replacement

The **clock page replacement algorithm** is also a variation fo the FIFO page replacement algorithm and is similar to the second chance page replacement algorithm but maintains the pages on a circular list.



As shown in the figure above, the pages are kept in a circular list in memory, where the hand points to the oldest page.

When a page fault occurs, the page being pointed to by the hand and performs actions based on the page's state deduced from the reference bit.

Reference bit	Page state	Actions taken
0	Old and unused.	The page is replaced in the clock in the same manner as the FIFO page replacement algorithm, and the hand is advanced to the next page.
1	Old but used.	The reference bit for the page is cleared and the hand is advanced to the next page. This is repeated until a page is found with a reference bit of 0.

Figure 6.13: Page classification for the clock page replacement algorithm.

This algorithm is realistic as it avoids the unnecessary movement of pages as seen in the second chance page replacement algorithm.

Least recently used (LRU) page replacement

The **least recently use page replacement algorithm** is based on the ideas that:

- pages that have been heavily used in the last few instructions will probably be heavily used again soon; and
- pages that have not been used for ages will probably remain unused for a long time.

It is necessary for a linked list of all pages to be maintained in memory with:

- the most recently used page at the front; and
- the least recently used page at the back.

As such, the algorithm swaps out pages that have been unused for the longest time.

Advantages	Disadvantages
	Requires time-consuming maintenance of the linked list as the list must be updated on every memory reference such that referenced pages are moved to the front.

Not frequently used (NFU) page replacement

The **not frequently used page replacement algorithm** keeps track of how often a page is used.

A counter is associated with each page and is initially set to 0. At each clock interrupt, the OS will scan all pages in memory. For each page, the reference bit will be added to its associated counter, such that:

- if the reference bit is 0, the counter will remain the same; and
- if the reference bit is 1, the counter will be incremented.

This means that counters can be used to track how often each page has been referenced.

When a page fault occurs, the page with the lowest counter is chosen for replacement.

Summary of page replacement algorithms

Algorithm	Comment
Optimal	Not implementable, but useful as a benchmark.
NRU	Very crude approximation of LRU.
FIFO	May throw out important pages.
Second chance	Big improvement over FIFO.
Clock	Realistic.
LRU	Excellent, but difficult to implement exactly.
NFU	Fairly crude approximation to LRU.

Evaluation of paging

Advantages	Disadvantages
Almost full utilisation of physical memory.	Some internal fragmentation as a process may not use memory in multiples of a page; usually the last page may not use the entire page size.
Can execute programs that have address space larger than physical memory.	Tuning the page size is tricky as <ul style="list-style-type: none"> • a smaller page size leads to less internal fragmentation but more pages are required and therefore the page tables will be larger; and • a larger page size leads to more unused programs to be loaded in to memory but less pages are required and therefore the page tables will be smaller.
No external fragmentation.	Large memory consumption as one page table for each process may consume large amounts of memory.
	Does not support the logical divisions of programs as page sizes are fixed and not based on the actual size of programs.

6.3 Segmentation

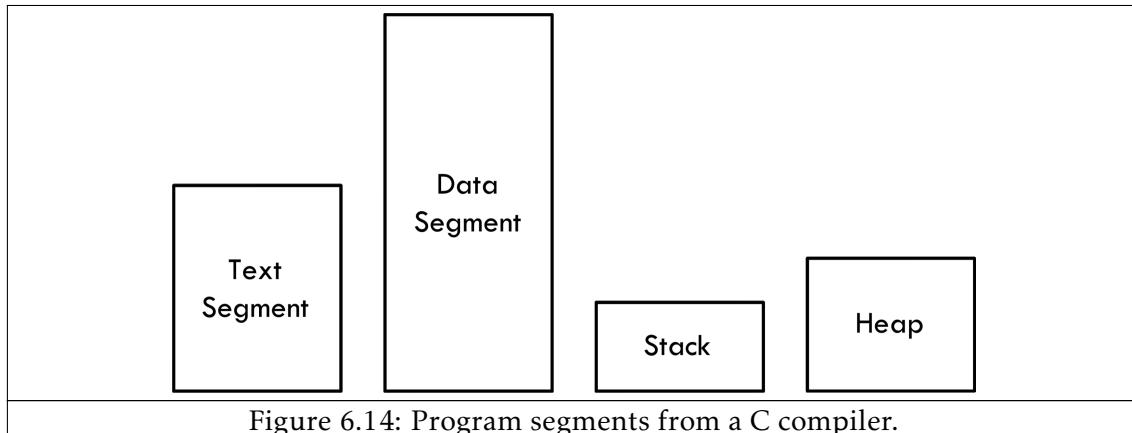
Definition

Segmentation is a memory management scheme that supports the logical user-view of programs.

How it works

Each program is split into variable chunks called segments according to the program's logical structure.

Different compilers for programs may split programs into different segments.



As shown in the figure above, a C compiler may generate four segments with different sizes:

- text segment (main code) and libraries;
- data segment;
- the stack; and
- the heap.

The stack and heap are of dynamic size, as they may shrink and grow over the course of the program's lifetime.

Memory is allocated to processes, segment by segment, in non-contiguous areas of physical memory.

Each segment has its own address space.

Segment table

Definition

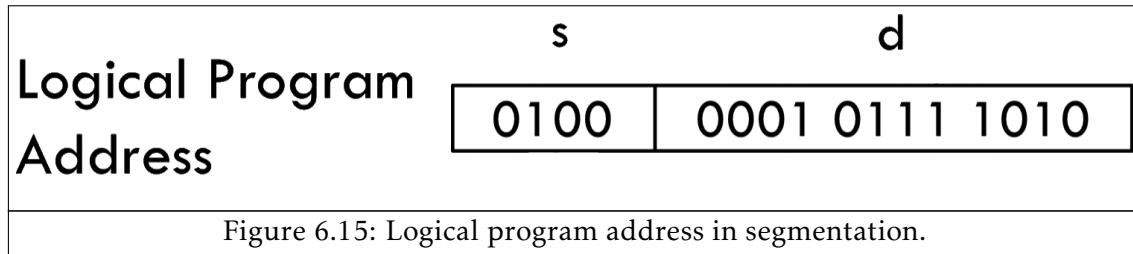
The **segment table** is responsible for mapping logical program addresses to physical memory addresses when using segmentation.

Mapping logical program addresses to physical memory addresses

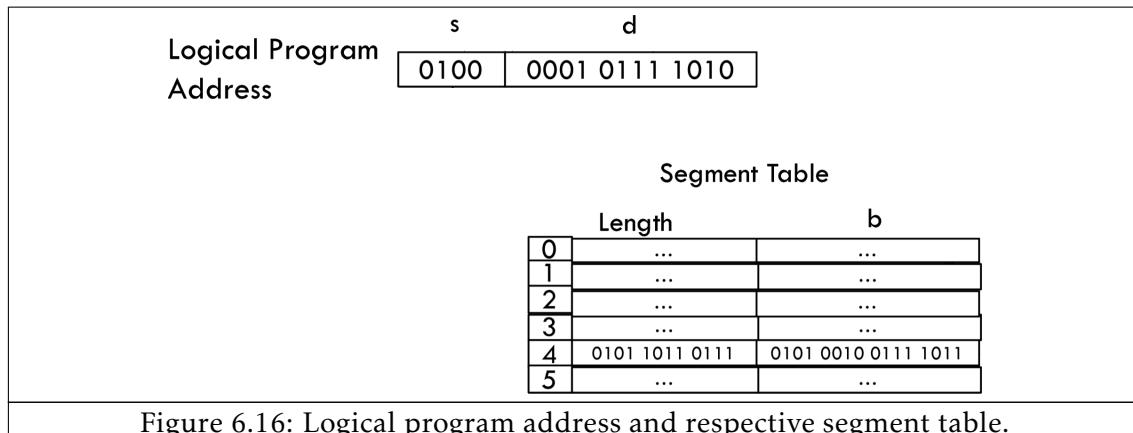
Each logical address is formed as (s, d) where

- s is the segment reference; and
- d is the displacement within that segment (offset).

s occupies 4 bits and while d occupies 12 bits.



In addition, each process has a segment table, which records the base address and the length of the segment.

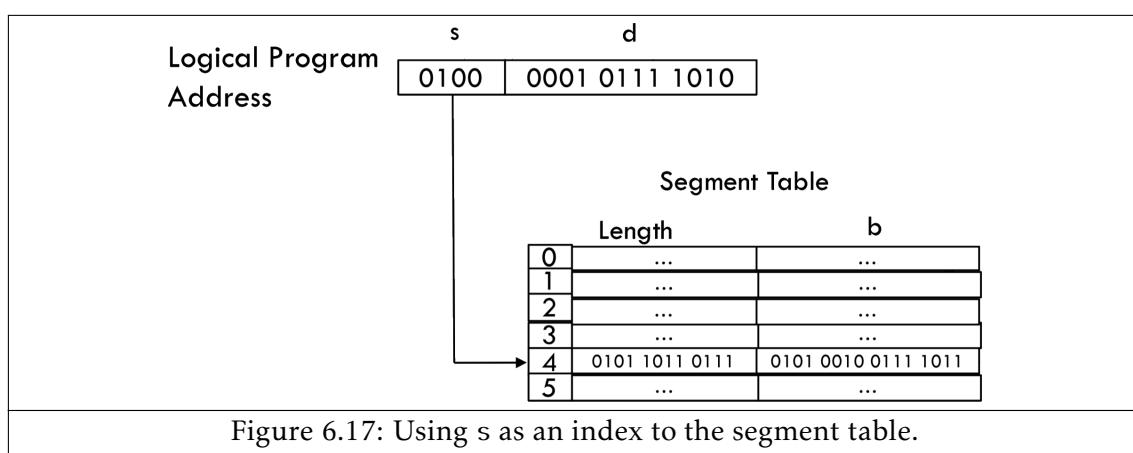


The implementation is performed by the hardware, as such the Memory Management Unit (MMU) is responsible for mapping the logical addresses to physical addresses.

Given a 16-bit program address 0100 0001 0111 1010, it can be deduced that:

$$\begin{aligned} p &= 0100_2 = 4_{10} \\ d &= 0001 0111 1010_2 = 378_{10} \end{aligned}$$

The segment reference (s) provides an index to a location in the segment table.



A check is also performed to ensure that the displacement value (d) is less than the segment length recorded in the segment table.

If the check is passed, the sum of the value of the base address (b) at index 4 in the segment table and the displacement value (d) from the logical program address can then be used as the physical memory address.



This can also be shown by performing the following operations:

$$\text{physical memory address} = \text{base address (b)} + \text{displacement value (d)}$$

$$\text{physical memory address} = 0101\ 0010\ 0111\ 1011 + 001\ 0111\ 1010$$

$$\text{physical memory address} = 0101\ 0011\ 1111\ 0101_2 = 21493_{10}$$

Protection and sharing

Segmentation aids protection as it:

- assigns different modes, such as read, write and execute, to each segment; and
- checks that memory references do not exceed the segment length, therefore preventing a process accessing another process' segments.

Segmentation aids sharing as a shared segment can be referenced by multiple processes, such as libraries.

Evaluation of segmentation

Advantages	Disadvantages
Support the logical divisions of programs as segments are based on the program's attributes.	Wasted space as segments are usually much larger than pages.
Segments can grow and shrink dynamically and independently, such as the stack and heap.	External fragmentation.
Aids sharing and protection.	

6.4 Comparison of paging and segmentation

Attribute	Paging	Segmentation
-----------	--------	--------------

Uses	Useful to have more address space without having to purchase more physical memory.	Useful to allow programs to be broken up into independent logical address spaces.
Memory division	Physical division of memory which is transparent to the user.	Logical division of memory which is visible to the user.
Size	Fixed size as each page is a pre-determined fixed size.	Variable size as segments can be dynamic, such as stack and heap.
Fragmentation	No external fragmentation.	Generates memory holes.
Address space	The total address space can exceed the size of physical memory.	

6.5 Paging segments

Paging segments is a method of combining paging and segmentation by assigning each segment with a page table.

Segments are typically larger than pages and, in the case where a segment does not fit in physical memory, this method must be used.

This allows the advantages of both memory management schemes to be present, such that this method:

- avoids external fragmentation;
- aids sharing and protection; and
- supports the user-view of programs.

Paging segments is roughly the method used in modern systems, such as Intel-based systems.

7. File Management

7.1	Non-volatile storage	114
7.2	Filing system	115
7.3	Files	115
7.4	Directories	123

7.1 Non-volatile storage

Definition

Non-volatile storage is a type of computer memory that can retrieve stored information even after having been power cycled.

Why is it used?

It is desirable that data should persist for later use and further processing. This should happen even after:

- power-down;
- critical errors;
- crashes; and
- forced shut-downs.

Solely relying on RAM would present issues as:

- information in the process address space would be lost when the process terminates; and
- for most applications it is far too small to contain all of the required data and instructions at any given time.

7.2 Filing system

Definition

A **filing system** is a method of organising and retrieving files from a storage medium, such as a hard drive.

Objectives

A filing system aims to:

- facilitate the long-term storage of data;
- allow the creation and deletion of files with automatic management of secondary storage;
- allow for file reference using symbolic names;
- provide access control to protect files against unauthorised access and allow sharing of files when required; and
- protect files against system failure.

7.3 Files

Definition

A **file** is a uniform logical unit of information created by a process.

Characteristics

Files also have an address space. However, the mapping in to a mass storage unit rather than RAM. As such, files can be described as a named collection of related information that is recorded on secondary storage, such as:

- the set of lines in a program; or
- the set of words in a text document.

Files are used for storing large amounts of data in the long-term.



As shown in the figure above, processes are able to access a file's data concurrently.

File naming

Motivation

File naming removes the need for a user to provide numerical addresses to access files. Instead, files can be accessed using a user-friendly name.

File naming conventions

Different OSs enforce different file naming conventions but most follow a common pattern.

Many OSs, such as Windows and Unix-based OSs, support up to approximately 260 characters for file names.

OSs place restrictions on which characters can be used in file names. For example:

- in Windows, the question-mark character (?) is not allowed in file names; while
- in Unix-based OSs, the question-mark character (?) is valid.

Some OSs distinguish between upper-case and lower-case in file names. For example:

- in Windows, files named ABC, Abc and abc would all represent the same file; while
- in Unix-based OSs, files name ABC, Abc and abc would all represent different files as the OS is case sensitive.

File extensions

A **file extension** is a string of characters attached to a filename, usually preceded by a full stop and indicating the format of the file.



The figure above shows the format for a typical file name, consisting of:

- a base name – often indicative of the file's contents; and
- a file extension – informs the user and the OS what type of data the file contains.

In MS-DOS, only three (3) characters were allowed for file extensions. While, in Unix-based OSs, the length of a file extension is down to the discretion of the user or process.

In Unix-based OSs, file extensions are not enforced by the OS. However, for example, a C compiler would require a file with the extension .c to compile.

GUI-based OSs typically attach meanings to extensions and associate applications to file extensions. For example, opening a file with the .docx extension may launch Microsoft Word.

However, this poses an issue as it is easy to "trick" and corrupt files by modifying extensions. This is addressed in macOS, which examines the contents of files in order to determine the file's type.

Extension	Meaning
.bak	Backup file.
.c	C source program.
.gif	Graphical Interchange Format image.
.html	World Wide Web HyperText Markup Language document.
.iso	ISO image of a CD-ROM (for burning to a CD).
.jpg	Still picture encoded with the JPEG standard.
.mp3	Music encoded in MPEG layer 3 audio format.
.mpg	Movie encoded with the MPEG standard.
.o	Object file (compiler output, not yet linked).
.pdf	Portable Document Format file.
.ps	PostScript file.
.tex	Input for the TEX formatting program.
.txt	General text file.
.zip	Compressed archive.

Figure 7.3: Common file extensions.

The figure above shows some common file extensions. However, this list is not exhaustive.

File attributes

Definition

File attributes are metadata associated with computer files that define file system behaviour. Each attribute can have one of two states:

- set; and
- cleared.

Common attributes

Attributes are considered distinct from other metadata.

Attribute	Meaning
Protection	Who can access the file and in what way.
Password	Password needed to access the file.
Creator	ID of the person who created the file.
Owner	Current owner.
Read-only flag	0 for read/write and 1 for read-only.
Hidden flag	0 for normal and 1 for do not display in listings.
System flag	0 for normal file and 1 for system file.
Archive flag	0 for has been backed up and 1 for needs to be backed up.
ASCII/binary flag	0 for ASCII file and 1 for binary file
Random access flag	0 for sequential access only and 1 for random access.
Temporary flag	0 for normal and 1 for delete file on process exit.
Lock flags	0 for unlocked and non-zero for locked.
Record length	Number of bytes in a record.
Key position	Offset of the key within each record.
Key length	Number of bytes in the key field.
Creation time	Date and time the file was created.
Time of last access	Date and time the file was last accessed.
Time of last change	Date and time the file was last changed.
Current size	Number of bytes in the file.
Maximum size	Number of bytes the file may grow to.

Figure 7.4: Common file attributes.

The figure above shows some common file attributes. However, this list is not exhaustive.

File structure

Files can be structured in any of several ways. Three common methods are:

- byte sequence;
- record sequence; and
- tree.

Byte sequence

In the **byte sequence** method, the OS considers a file to be an unstructured sequence of bytes.

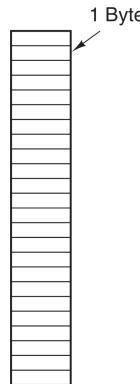


Figure 7.5: Byte sequence for a file.

The OS is not concerned with the file's contents. Instead, it is presented with bytes of data. Any meaning must be imposed by user-level programs.

Both UNIX-based OSs and Windows use this approach.

Record sequence

In the **record sequence** method, a file is a sequence of fixed-length records, each with some internal structure.

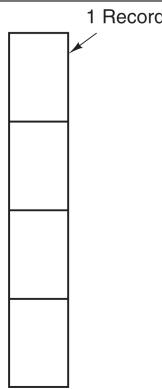


Figure 7.6: Record sequence for a file.

Central to the idea of a file being a sequence of records is the idea that:

- the read operation returns one record; and
- the write operation overwrites or appends one record.

No modern general-purpose system uses this model as its primary file system. However, this was useful when 80-column punched cards and 132-character line printer paper were used on mainframe computer systems.

Tree

In the **tree** method, a file consists of a tree of variable-length records, each containing a key field in a fixed position in the record.



The tree is sorted on the key field, to allow rapid searching for a particular key.

This type of file is different from the unstructured byte streams used in Windows and Unix-based OSs but is used on some large mainframe computers for commercial data processing and database systems.

File access

Definition

File access refers to the way in which the files stored may be accessed.

Sequential access

In **sequential access** a process can read all the bytes or records in a file in order, starting at the beginning.

It is not possible to skip around and read them out of order. However, sequential files can be rewound so they could be read as often as needed.



Figure 7.8: Sequential access of a file.

As shown in the figure above, sequential access allows the operations:

- read next; and
- write next.

Sequential files were convenient when the storage medium was magnetic tape rather than disk, as magnetic tapes are accessed sequentially. However, this method of access is still interesting today due to the locality principle (the tendency of a processor to access the same set of memory locations repetitively over a short period of time).

Random (direct) access

In **random (direct) access** a process can read the bytes or records of a file out of order and access records by key rather than by position.

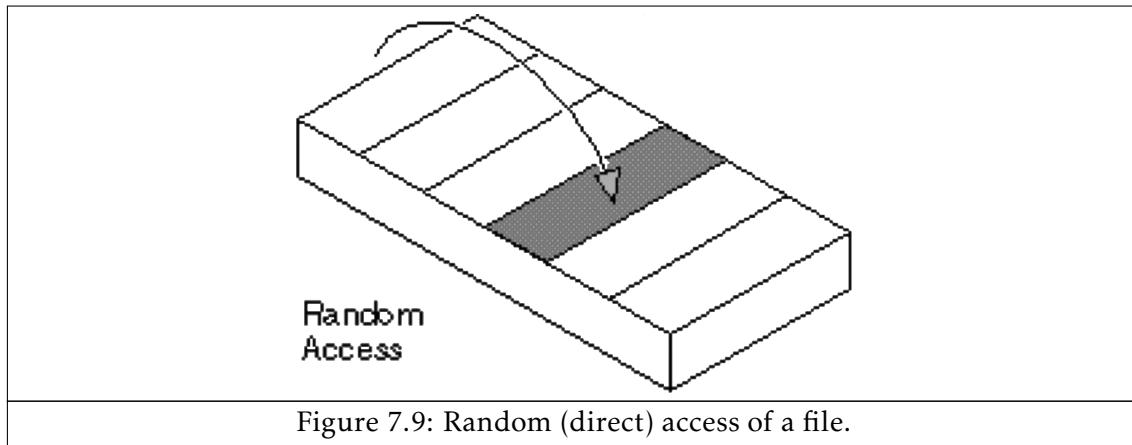


Figure 7.9: Random (direct) access of a file.

As shown in the figure above, random (direct) access allows the operations:

- read n; and
- write n.

n represents a relative block number.

Random (direct) access is essential for most modern applications.

File types

Many OSs support several types of files, as shown below.

File type	Description
Regular	Ordinary files that contain data.
Executable	files that contain program code that can be executed.
Directory / folder	System file containing references to other files.
Character special	Not true files, rather directory entries which refer to character devices and allow communication with I/O devices such as printers.
Block special	Not true files, rather directory entries which refer to block devices and are used for communication with storage devices such as disks and memory.

Figure 7.10: File types

Typical file operations

OSs offer several system calls for file management, as shown below.

Operation	Description
Read	Data are read from file.
Write	Data are written to an existing file.
Append	Restricted form of write. Data can only be added to the end of a file, while maintaining the existing contents of the file
Seek	Used for random access files. Repositions the file pointer to a specific location in the file.
Get attributes	Retrieve the attributes of the file. For example, used by the C compiler.
Set attributes	Some attributes can be set by the user, such as protection-mode.
Rename	Change the name of an existing file.

Figure 7.11: Typical file operations.

7.4 Directories

Definition

A **directory** (or **folder**) is a file system cataloging structure which contains references to other computer files, and possibly other directories.

Characteristics

Most filing systems allow files to be grouped together in to directories, resulting in a more logical organisation.

Directories allow:

- operations to be performed in bulk on groups of files, such as copying files or setting one of their attributes; and
- different files to have the same file name provided that they are in different directories.

File descriptor table

Each directory is managed using a special file containing a file descriptor table.



As shown in the figure above, the file descriptor table contains descriptors for each file under that directory which correspond to specific entries on the global file table.

Directory structure

Definition

The **directory structure** determines how the OS organises the files and directories in the file system.

Single-level directory systems

Single-level directory systems use one directory for all files in the file system.

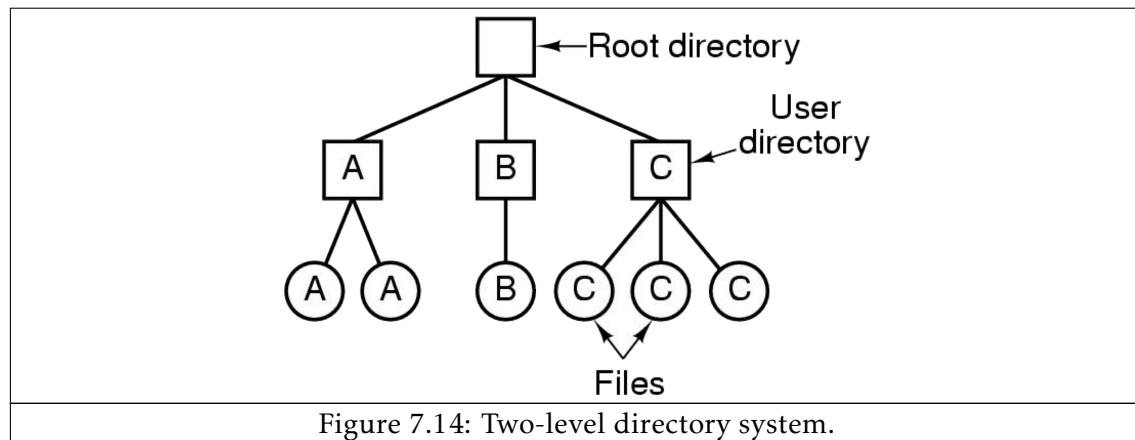


The figure above shows a single-level directory system with four files owned by three different users A, B and C.

Advantages	Disadvantages
Simplicity as the root directory is the only directory to be implemented.	Limited grouping capability as the root directory is the only directory present in the file system.
Ability to quickly find files as only the root directory must be searched.	Naming issues as the root directory is the only directory present in the file system and so no two files can have the same name.

Two-level directory systems

Two-level directory systems use a separate directory for each user.



The figure above shows a two-level directory system with directories for the users A, B and C.

Advantages	Disadvantages
Flexible naming as files can have the same name provided that they are in different user directories.	Limited grouping capability as the user directory is the only directory present in the file system.

Hierarchical directory systems

Hierarchical directory systems maintain directories in a tree-like structure.



The figure above shows a two-level directory system with directories for the users A, B and C.

Advantages	Disadvantages
Grouping capabilities as directories can be created by users to organise files.	Some complexity as a method for browsing and locating files must be implemented.
More flexible naming as files can have the same name provided that they are in different directories.	

The exact implementation may differ between OSs.



The figure above shows the directory tree used by Unix-based OSs.

Directory paths

Browsing files

Two common methods are used for browsing files: absolute path name; and relative path name.

An **absolute path name** is relative to the root directory and is unique. For example:

- in MS-DOS and Windows – \NTU\syssoftware\demo.c; and
- in Unix-based OSs – /NTU/syssoftware/demo.c.

A **relative path name** is relative to the current directory (or present working directory). For example:

- in MS-DOS and Windows – if the current working directory is \NTU\syssoftware, the file whose absolute path is \NTU\syssoftware\demo.c can be referenced as demo.c; and
- in Unix-based OSs – if the current working directory is /NTU/syssoftware, the file whose absolute path is /NTU/syssoftware/demo.c can be referenced as demo.c.

Special entries

Some special entries exist to provide shortcuts:

- . (dot) – shortcut for current directory; and
- .. (dot dot) – shortcut for parent directory.

Example Unix directory commands

Given that it is required to backup the dictionary file located in the lib directory:

- if the current directory is / (root directory), then the command to perform this action would be cp /usr/lib/dictionary /usr/lib/dictionary.bak; and
- if the current directory is /usr/lib/, then the command to perform this action would be cp dictionary dictionary.bak.

Given that it is required to move the dictionary file located in the lib directory to the ast directory:

- if the current directory is / (root directory), then the command to perform this action would be mv /usr/lib/dictionary /usr/ast/dictionary; and
- if the current directory is /usr/ast/, then the command to perform this action would be mv dictionary /usr/ast/dictionary.

Directory file operations

OSs offer several system calls for directories, as shown below.

Operation	Description
Create	An empty directory is created.
Delete	An existing directory can be deleted. Non-empty directories on many Unix-based OSs will require further clarification by using a "force" command.
Opendir	Directories can be read. For example, when listing all of the files in the directory.
Closedir	When a directory has been read, it should be closed to free up space.
Rename	Change the name of an existing directory.
Link	A technique that allows a file to appear in more than one directory.
Unlink	If a file is unlinked, it is only normally present to one directory.

Figure 7.17: Typical directory operations.

There are more systems calls depending on the OS used.

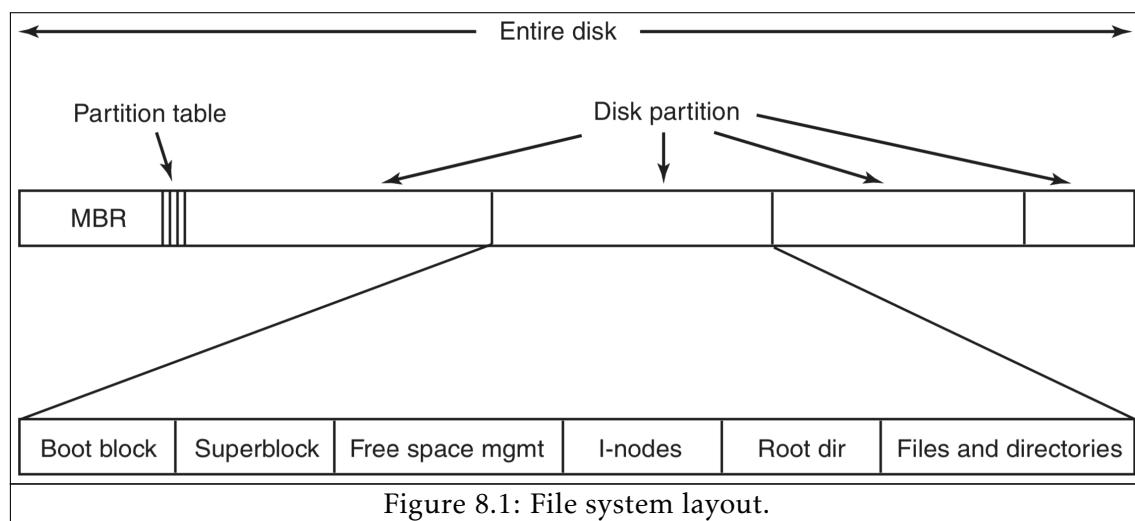
8. File System Implementation

8.1	File system layout	129
8.2	File block allocation methods	130
8.3	Implementation of directories/folders	137
8.4	Locating files using i-nodes	139
8.5	Determining block size	141
8.6	Tracking free space	142
8.7	Consistency	144
8.8	Journaling	145
8.9	Backups	146
8.10	Drives and mount points	146
8.11	RAID	147

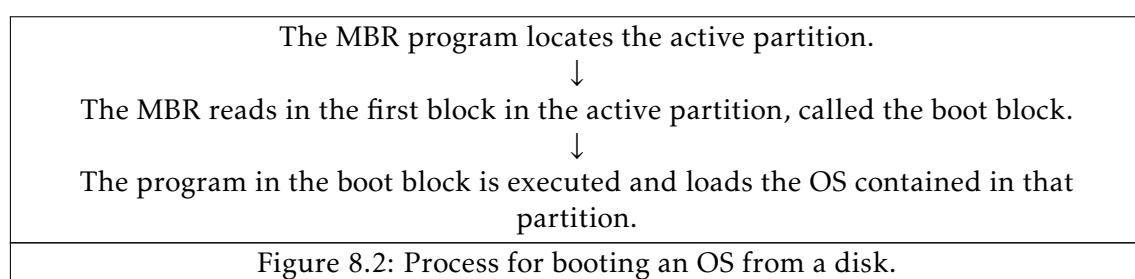
8.1 File system layout

File systems are stored on drives. Most disks can be divided up into one or more **partitions/volumes**, each holding an independent file system.

The terms "drives" and "disk" can be used interchangeably despite referring to different types of mass storage



Sector 0 of the disk is the **Master Boot Record (MBR)**. This is used to boot the computer via a **boot block** from a specified partition, from which the OS is loaded.



The end of the MBR contains the partition table. This table gives the starting and ending addresses of each partition. Every partition starts with a boot block, even if it does not contain a bootable operating system. Besides, it may contain one in the future.

The **superblock** contains key parameters about the file system on the partition and is read in to memory when the computer is booted or the file system is first touched. Typical information in the superblock includes:

- a magic number to identify the file system type;
- the number of blocks in the file system; and
- other key administrative information.

8.2 File block allocation methods

Objective

File block allocation methods aim to keep track of which disk blocks go with which file.

Various methods are used in different OSs, including:

- contiguous allocation;
- linked list allocation (non-contiguous);
- linked list with file allocation table; and
- i-nodes.

Contiguous allocation

How it works

In **contiguous allocation**, disks are split in to blocks of fixed size.

This method is often used for CDs.



The figure above shows how a different number of contiguous blocks may be allocated to different files. For example:

- if the fixed size of the blocks was 1KB, a file with size 50KB would be assigned to 50 consecutive blocks; and
- if the fixed size of the blocks was 2KB, a file with size 50KB would be assigned to 25 consecutive blocks.



As shown in the figure above, removing files frees up the blocks that were once allocated to the deleted file.

Evaluation

Advantages	Disadvantages
Simple implementation as it is only required to store the first block address and its length.	Overhead as it is required to track the size of the files when initially created.
Good performance.	Files cannot grow as their number of blocks is fixed and the size of the blocks is fixed. For example, it would not be possible to edit an existing document and introduce new data.
Allows easy random access.	Internal fragmentation as the data inside a file block may be smaller than the fixed block size.
Resilient to drive faults as damage to a single block results in only localised loss of data.	External fragmentation as when files are deleted, holes may be generated. This can be overcome by performing compaction, however this increases the system overhead.

Linked list allocation (non-contiguous)

How it works

In **linked list allocation**, files are stored as linked lists of blocks.



Figure 8.5: Linked list allocation.

As shown in the figure above, each block contains:

- a pointer to the next block (occupies the first word of each block); and
- data.

The final block contains a null pointer, usually denoted by a zero (0) as no further blocks for the file exist.

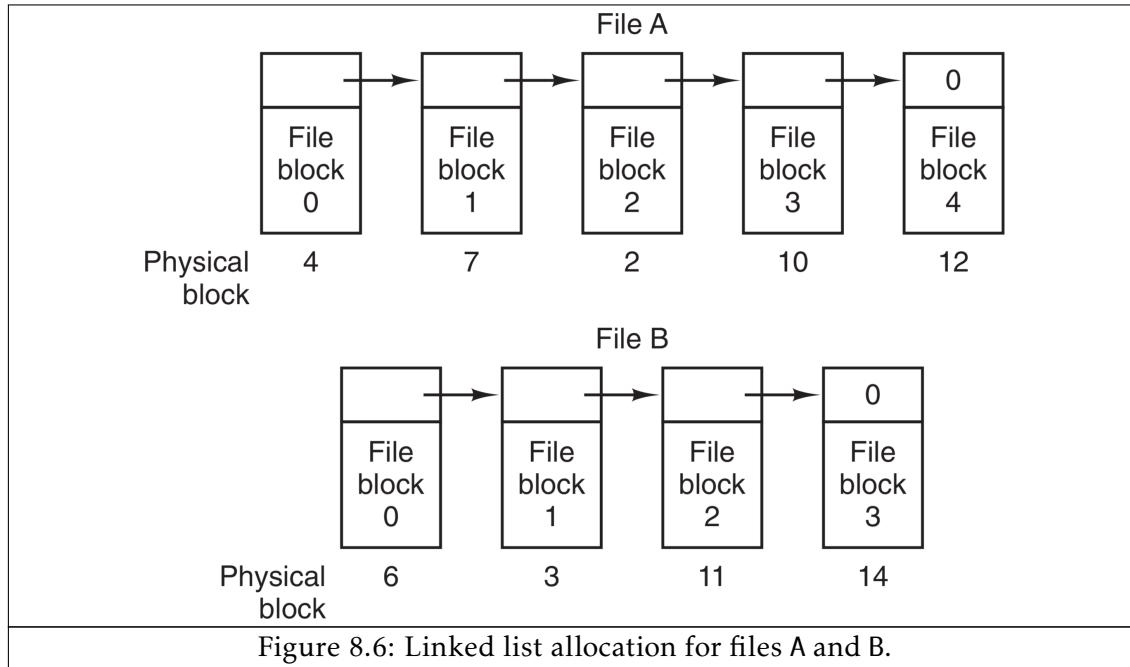
Evaluation

Advantages	Disadvantages
Every block can be used.	Does not support random access as it is very slow due to the necessity to traverse through the chain of blocks.
File size does not have to be known beforehand as files can grow because new blocks can be easily added by changing the null pointer on the last block to point to the newly added block.	Wasted space as some space is lost for useful data within each block due to storage of the pointer.
No external fragmentation.	
No internal fragmentation , except for the last block.	

Linked list with file allocation table (FAT)

When using a **linked list with file allocation table (FAT)**, the FAT is located in main memory and stores the pointer of all of the blocks.

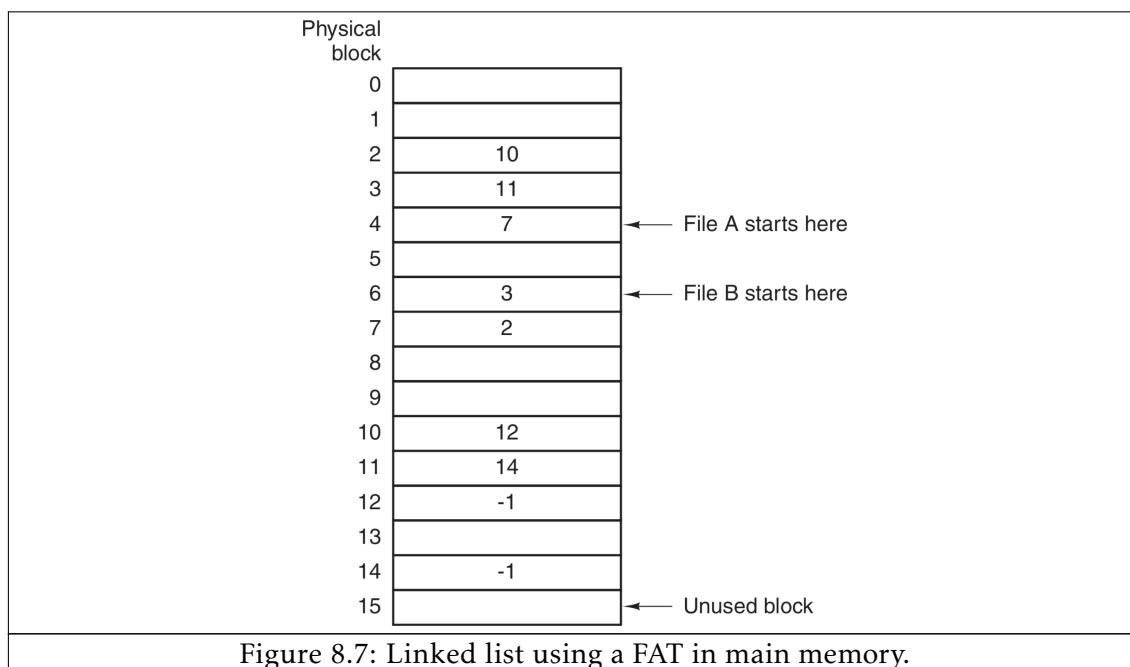
This method is used by MS-DOS and older versions of Windows.



The figure above shows the linked list allocation for files A and B:

- file A uses disk blocks 4, 7, 2, 10, and 12; and
- file B uses disk blocks 6, 3, 11, and 14.

The pointer word from each block is put in to the FAT in main memory.



For example, using the FAT in the figure above:

- for file A, it is possible to start with block 4 and follow the chain to the end at 12; and
- for file B, it is possible to start with block 6 and follow the chain to the end at 14.

Both chains are terminated with a special marker, usually denoted as -1, that is not a valid block number. This is used to signal the end of the blocks for a given file.

Evaluation

Advantages	Disadvantages
No wasted space as the pointer is not stored in the block, as it was in linked list allocation.	Does not scale well to large disks as the FAT may get too large for main memory, especially for drives with large capacities that require more pointers to be stored as there are more files and more blocks.
Random access is much faster than linked list allocation as the FAT is stored in main memory, which is a faster storage medium.	Serious data loss can occur should damage be inflicted on the FAT. However, this can be mitigated by storing several backups of the FAT.
No disk references are required when following a chain in the FAT as it is stored in main memory.	

Index-nodes (I-nodes)

When using **index-nodes** (i-nodes), each file is associated with an i-node which lists:

- all of the attributes of the file's blocks; and
- the disk addresses of the file's blocks.

This method is used by Unix-based OSs.



The figure above shows an example of some i-nodes stored for files on a disk.

Given the i-node, it is then possible to find all the blocks of the file.

Evaluation

Advantages	Disadvantages
More efficient memory usage than FAT as only the i-node of the file must be present in main memory and only when the corresponding file is opened.	Files may grow beyond the limit of the fixed number of addresses in each i-node. However, this can be mitigated by reserving the last disk address for the address of a block containing more disk-block addresses. This can be further improved, as shown in Unix in the section below.
	List disk addresses must point to an address block instead of a data block.

I-nodes in Unix

In Unix, the i-node contains a number of **direct pointers** to disk blocks. There are typically ten direct pointers.

In addition to the direct pointers, there are three **indirect pointers** that point to further address blocks, which eventually lead to a disk data block.



Figure 8.9: An i-node in a Unix-based OS.

As shown in the figure above, the indirect pointer includes:

- a single level of indirections;
- a double indirect pointer; and
- a triple indirect pointer.

These enable the file blocks to be located for a file, without restricting files to a fixed number of addresses.

8.3 Implementation of directories/folders

Directory entries

To open a file, the path name is used to locate its directory entry.

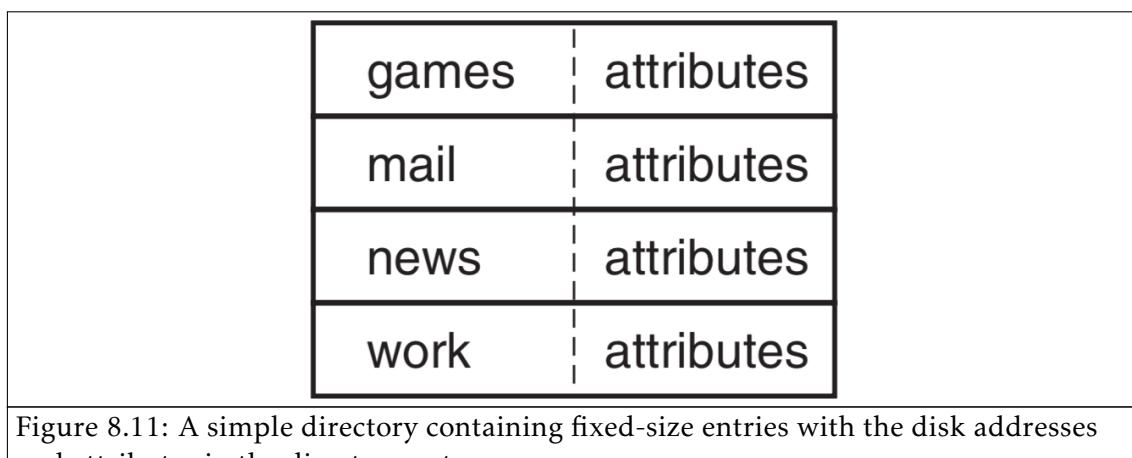
The directory entry provides a mapping from a file name or file descriptor to the disk blocks that contain the data.



As shown in the figure above, the file descriptor table contains descriptors for each file under that directory which correspond to specific entries on the global file table.

The directory entry contains all the information needed to find the disk blocks for a given file:

- in contiguous allocation, the directory entry contains the addresses of the entire file;
- in linked list allocations, the directory entry contains the first disk block address; and
- in an i-node implementation, the directory entry contains the file's respective i-node number.



The figure above shows the state of the file descriptor table if contiguous allocation or

linked list allocations were being used.

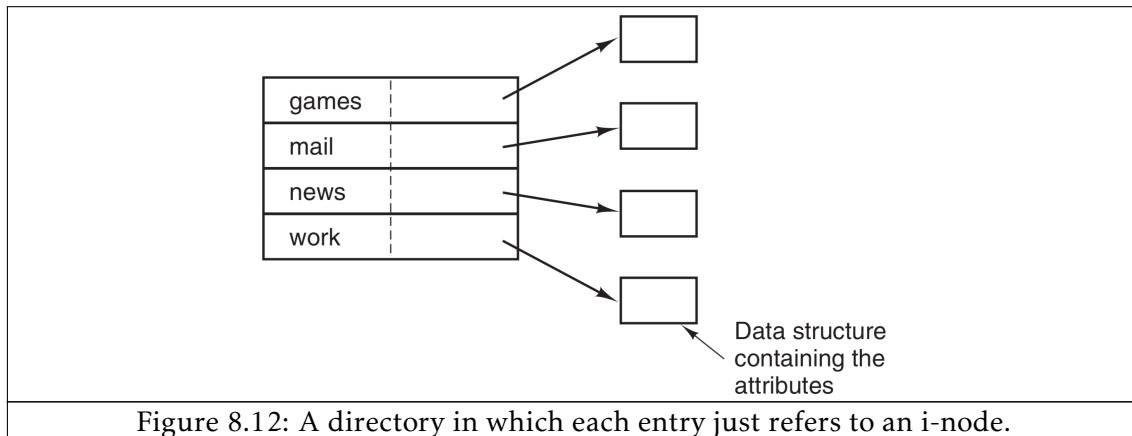


Figure 8.12: A directory in which each entry just refers to an i-node.

The figure above shows the state of the file descriptor table if an i-node implementation was being used.

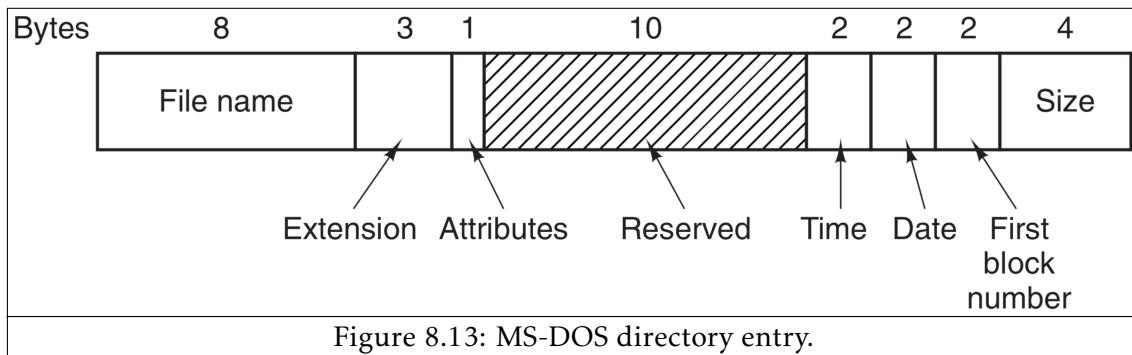


Figure 8.13: MS-DOS directory entry.

As shown in the figure above, in MS-DOS, the directory entry contains the attributes.

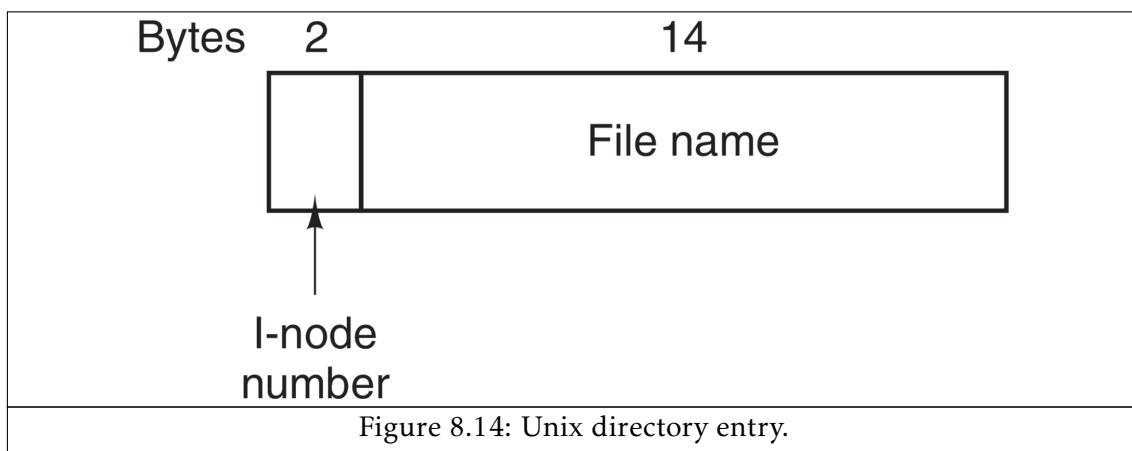


Figure 8.14: Unix directory entry.

As shown in the figure above, in Unix-based OSs, the directory entry contains an i-node number and a file name, where:

- the attributes are stored in the i-node; and
- all i-nodes have a fixed location on the disk, so locating an i-node is simple and fast.

8.4 Locating files using i-nodes



As seen before, the figure above shows that the indirect pointer includes:

- a single level of indirections;
- a double indirect pointer; and
- a triple indirect pointer.

These enable the file blocks to be located for a file, without restricting files to a fixed number of addresses, as shown below.



The process for accessing a relative path name is identical, except that the search begins from the current working directory.

Using this process, the i-node of the mbox file is kept in main memory until the file is closed. However, locating a file in other hierarchical file systems is similar, lookup is completed component by component.

8.5 Determining block size

All allocation methods require the disk to be split in to fixed-size blocks.



Nearly all modern operating systems use the fixed-sized blocks as shown in the figure above.

There is a similar trade off with block size as there was with page size in memory management. As such, tuning the block size is tricky as:

- **a smaller block size may lead to increased overhead**, as a file may occupy several blocks and therefore there will be a longer access time since more blocks have to be located; and
- **a larger block size may lead to wasted space** as a small file of 1KB would occupy a big chunk, such as 32KB, of the disk block.

Typical block sizes include: 512 Bytes, 1KB and 2KB.

8.6 Tracking free space

Linked list method

The **linked list method** is a way of tracking free space on a disk where some of the free blocks are used to hold free disk numbers.



As shown in the figure above, each block holds as many free disk block numbers as possible. The blocks are linked together resulting in a linked list of free blocks, by using one slot in each list for a pointer to the next block.

When considering a 1TB disk, which has about 1 billion disk blocks, storing all these addresses at 32 bit per block requires about 4 million blocks.

Generally, free blocks are used to hold the free list, so the storage is essentially free.

Bitmap method

The **bitmap method** is another way of tracking free space on a disk where a bit map is stored that records whether particular blocks are free or not.

	1001101101101100
	0110110111110111
	1010110110110110
	0110110110111011
	1110111011101111
	1101101010001111
	0000111011010111
	1011101101101111
	1100100011101111
≈	≈
	0111011101110111
	1101111101110111

Figure 8.19: Tracking free space using a bitmap.

As shown in the figure above, the bitmap records:

- a zero (0) if the disk block is used; and
- a one (1) if the disk block is free.

A disk with n blocks requires a bitmap with n bits.

When considering a 1TB disk, which requires around 130,000 1KB blocks to store the information, it is clear to see that a bitmap occupies less space than a linked list. Only if the disk is nearly full (i.e. has few free blocks) will the linked list method require fewer blocks than the bitmap.

Comparison of linked list method and bitmap method

Criteria	Linked list method	Bitmap method
Occupied space	32 bits per block.	1 bit per block.
Number of blocks	The number of blocks required to track free space becomes less as the drive gets fuller.	The number of blocks required to track free space is constant.

8.7 Consistency

A file system may be in any one of the following states:

- consistent;
- missing block;
- duplicate block in free list; or
- duplicate data block.

Many file systems read blocks, modify their contents and write them out later.

If the system crashes before all the modified blocks have been written out, the file system can be left in an inconsistent state. This problem is especially critical if some of the blocks that have not been written out are:

- i-node blocks;
- directory blocks; or
- blocks containing the linked list.

To deal with inconsistent file systems, most OSs provide a utility program that checks for file system consistency. For example:

- Windows provides chkdsk; and
- Unix-based OSs provide fsck.

These utilities can be run whenever the system is booted, especially after a crash.

8.8 Journaling

Definition

A **journaling file system** is a file system that keeps track of changes not yet committed to the file system's main part by recording the intentions of such changes in a data structure known as a "journal", which is usually a circular log.

Why is journaling required?

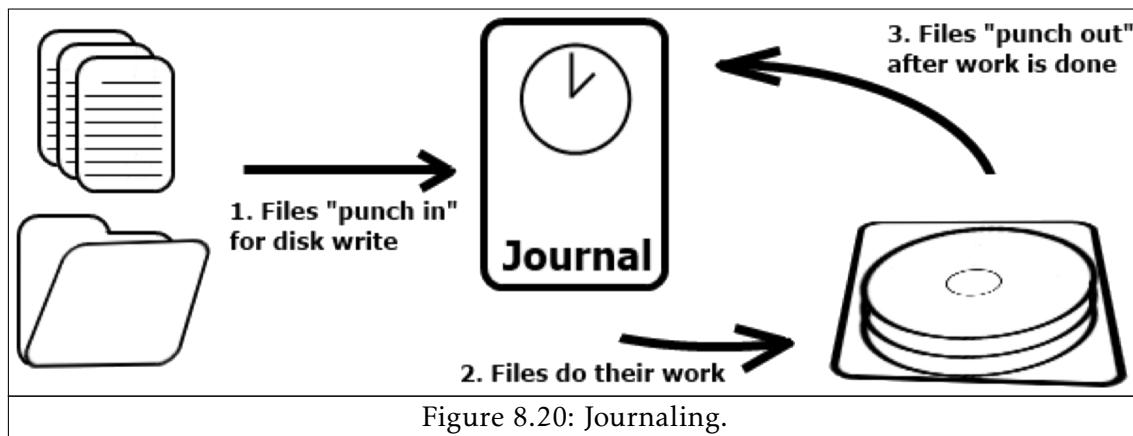
A single file operation may involve multiple actual writes to the disk. For example, in a Unix-based OS, removing a file requires the following three operations to be performed:

- remove the file from its directory;
- release the i-node to the pool of free i-nodes; and
- return all the disk blocks to the pool of free disk blocks.

This process is similar in Windows.

In the event that a system crashes between any one of the three tasks above, the i-nodes and file blocks will not be accessible from any file or available for reassignment. This would require time consuming requires by using the consistency utilities.

How it works



The journaling file system uses a special disk area to make a log entry. The log entry lists the actions to be completed.

In the event of failure, the log is used to bring the disk back in to a consistent state by completing all pending actions.

Log entries are erased once the operations complete successfully.

8.9 Backups

Definition

A **backup**, or **data backup**, refers to the process of copying data on the disk in to an archive file so that it may be used to restore the original after a data loss event.

Why are backups used?

Backups help to solve potential problems by allowing:

- recovery from disasters; and
- recovery from mistakes.

Issues

The process for **backing up is slow** as all files must be copied. However, this can be mitigated by performing incremental backups where only those files that have been changed since the previous backup are copied. This reduces the time taken for a backup, apart from the initial backup.

Backups occupy large amounts of space, especially if the stored data on the disk is larger. This can be mitigated by compressing the backup data. However, a corrupted part of the compressed data can render an entire backup unreadable.

The **system must be inactive** while a backup is being performed. This makes it difficult to continue using the system as the file system is actively being backed up.

Must ensure **physical security of backup media** as physical corruption of backups can occur on CD-ROMS, disks and magnetic tapes etc.

8.10 Drives and mount points

In Windows, multiple drives appear with distinct drive letters such as C, D, E etc. When a new drive is attached, it is assigned a unique letter.

In Unix-based OSs, a uniform file system is presented with a single root where new drives usually appear in the /media/ directory.

8.11 RAID

Definition

Redundant Array of Inexpensive Disks (RAID) is a data storage virtualisation technology that combines multiple physical disk drive components into one or more logical units for the purposes of data redundancy, performance improvement, or both.

RAID provides another layer of abstraction as data is distributed across several physical disks which appear to be a single logical disk.

Why is RAID required?

Disk drives remain the bottleneck in computer systems as their access speeds are much slower than the CPU's clock speed. As such, RAID is often used for speed enhancements along with other benefits, as seen below.

Types of RAID

There are multiple types of RAID that are suited to different use cases.

Striping, mirroring and parity

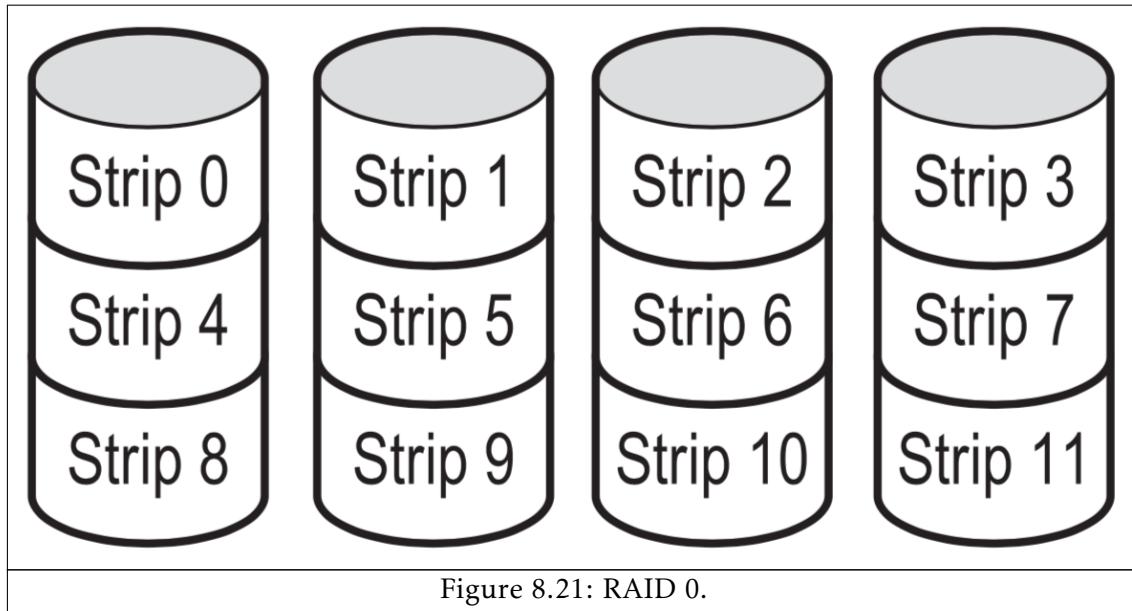
Different types of RAID use different techniques including:

- striping – distributing the data across several disks in a way which gives improved speed and full capacity;
- mirroring – uses more than one disk which store the same data; and
- parity – uses an additional disk to provide fault tolerance where a single data bit is added to the end of a data block to ensure the number of bits in the message is either odd or even.

Striping provides no security on its own but gives the best speed and capacity.

Mirroring degrades speed and capacity but provides security.

RAID 0



Striping	YES
Mirroring	NO
Parity	NO

The capacity of a RAID 0 volume is the sum of the capacities of the disks in the set. However, striping distributes the contents of each file among all disks in the set and therefore the failure of any disk causes all files to be lost.

RAID 1



Striping	NO
Mirroring	YES
Parity	NO

Data is written identically to two drives, thereby producing a "mirrored set" of drives. Therefore, any read request can be serviced by any drive in the set.

If a request is broadcast to every drive in the set, it can be serviced by the drive that ac-

cesses the data first (depending on its seek time and rotational latency), improving performance.

Sustained read throughput, if the controller or software is optimised for it, approaches the sum of throughputs of every drive in the set, just as for RAID 0.

Actual read throughput of most RAID 1 implementations is slower than the fastest drive.

Write throughput is always slower because every drive must be updated, and the slowest drive limits the write performance.

The array continues to operate as long as at least one drive is functioning.

RAID 1+0



Striping	YES
Mirroring	YES
Parity	NO

RAID 1+0 takes advantage of both RAID 0's striping and RAID 1's mirroring. As such, if striping is applied to two disks, an additional two disks are required to mirror the striped drives.

9. Input/Output (I/O)

9.1 Input/Output (I/O) devices

150

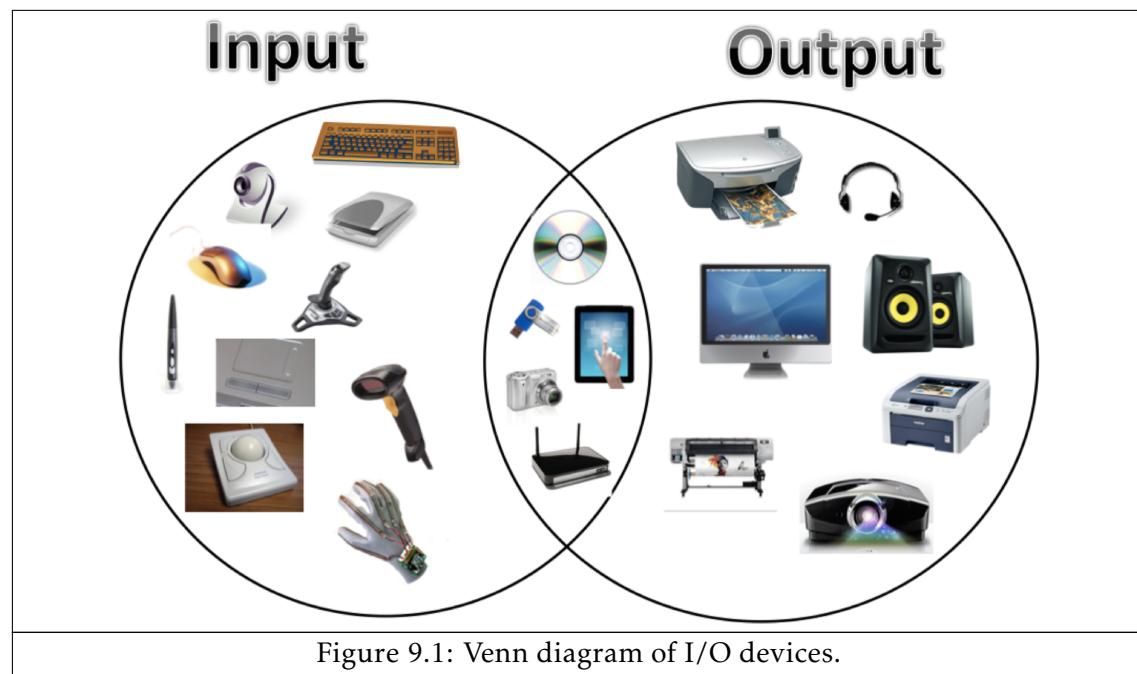
9.2 Input/Output (I/O) system

152

9.1 Input/Output (I/O) devices

Defintion

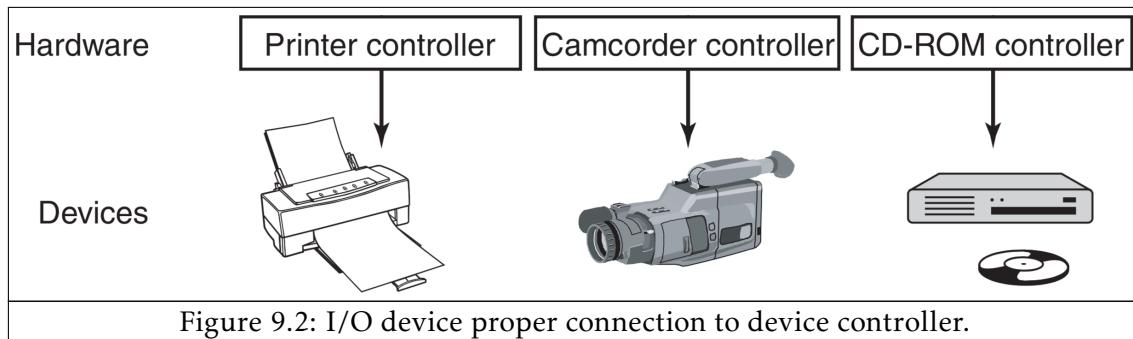
An **Input/Output (I/O) device** is a hardware device that has the ability to accept inputted, outputted or other processed data. It also can acquire respective media data as input sent to a computer or send computer data to storage media as storage output.



An I/O device can be responsible for input and output as shown in the cross section in the Venn diagram above.

I/O devices typically consists of:

- the **device proper** – an electro-mechanical component generally located outside the computer, such as a peripheral; and
- the **device controller** – an electrical component often in the form of a chip on the mother board or a card that can be inserted in to an expansion slot, such as on the Peripheral Component Interconnect Express (PCIe) slot.



As shown in the figure above, the I/O device proper is connected using a cable to the device controller for communication to and from the outer world.

I/O controllers

Each I/O controller has a few registers that are used for communicating with the CPU:

- write operation – the OS can command the device to deliver/accept data or switch on or off itself, or perform any other action; and
- read operation – the OS can determine the device's current state, such as whether it is prepared to accept a new command.

Typical characteristics and specifications

Evaluation

Device	Data rate
Keyboard	10 bytes/sec
Mouse	100 bytes/sec
56K modem	7 KB/sec
Scanner at 300dpi	1 MB/sec
Digital camcorder	3.5 MB/sec
4x Blu-ray disc	18 MB/sec
802.11n Wireless	37.5 MB/sec
USB 2.0	60 MB/sec
FireWire 800	100 MB/sec
Gigabit Ethernet	125 MB/sec
SATA 3 disk drive	600 MB/sec
USB 3.0	625 MB/sec
SCSI Ultra 5 bus	640 MB/sec

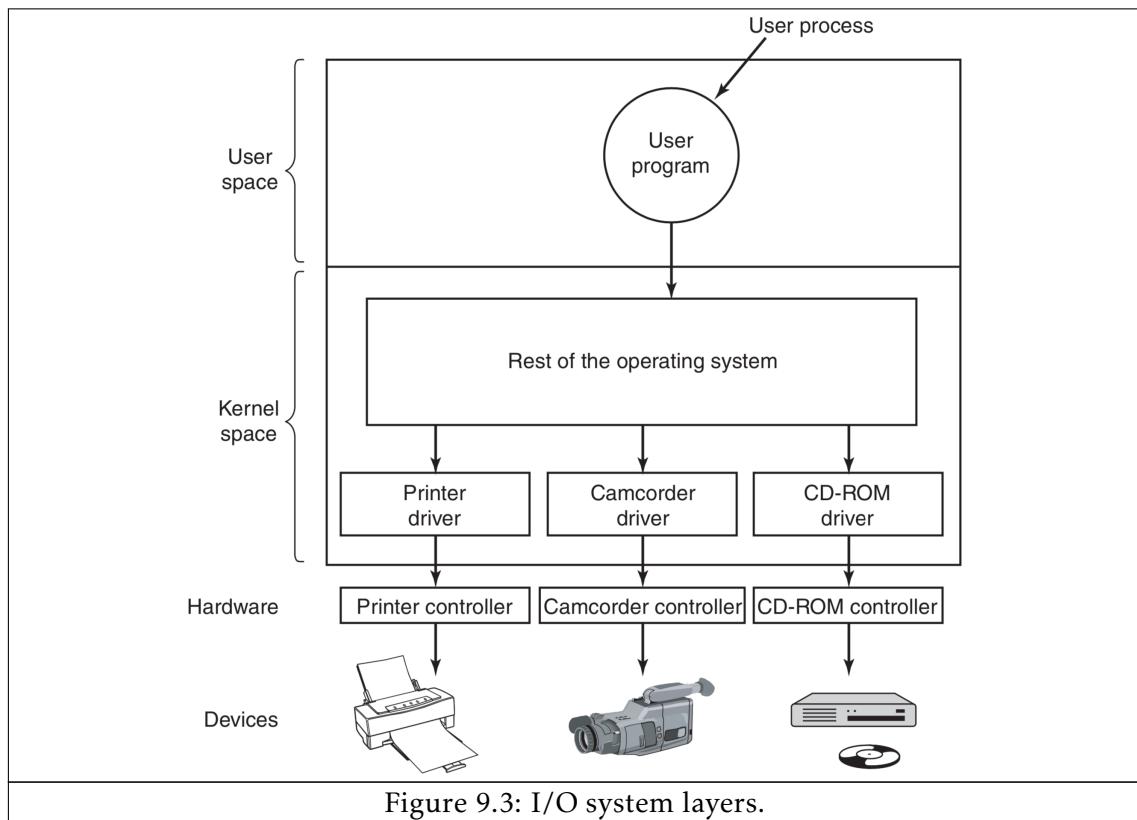
Single-lane PCIe 3.0 bus	985 MB/sec
Thunderbolt 2 bus	2.5 GB/sec
SONET OC-786 network	5 GB/sec

9.2 Input/Output (I/O) system

Definition

The **Input/Output (I/O) system** includes all the agents that are involved in I/O operations, including:

- the device(s) themselves;
- the device controller(s);
- the OS mediating the operations; and
- the user process requesting the operations.

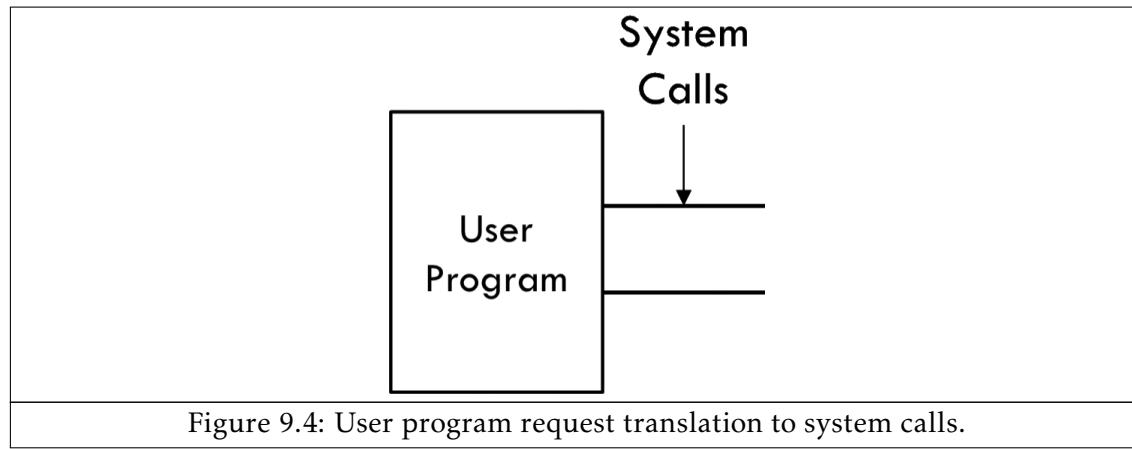


The figure above shows the several layers on which the I/O system depends, with the agents residing in each layer.

Structure

User program

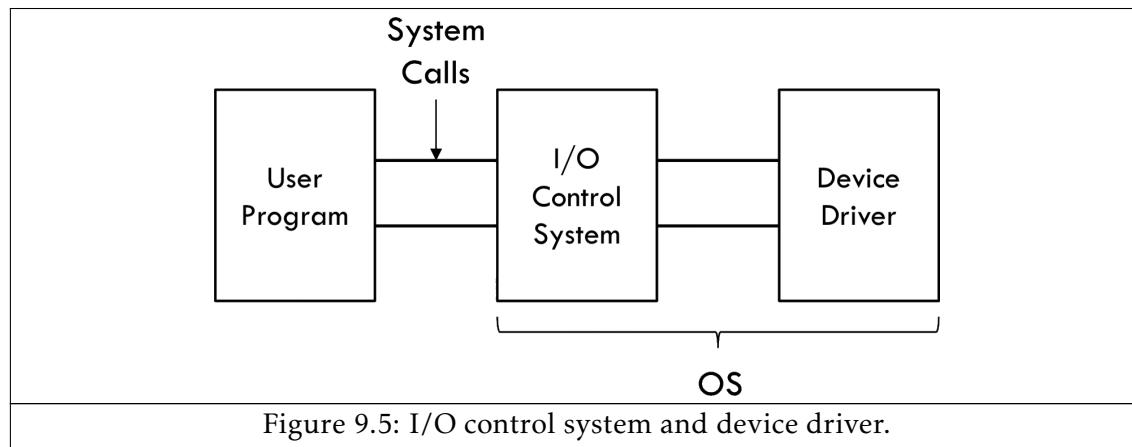
The user program makes an I/O request in a high level language.



As shown in the figure above, the request is translated and relayed by system calls.

I/O control system

The I/O control system is a piece of device-independent software that deals with I/O-related system calls.



As shown in the figure above, the I/O control system is a platform for installed device drivers to communicate.

Device driver

A device driver is a software module that manages the communication with a specific I/O device. This is written by the manufacturer of the piece of hardware and is specifically designed for that piece of hardware.

Specific device drivers are provided by hardware manufacturers on installation discs or Internet download links. Hardware manufacturers must create device drivers for each operating system for which the hardware must be compatible.

Generic device drivers, or **non-specific device drivers**, are often provided in operating systems. These may not be as good as specific device drivers because they may not be supported by the device or provide all of the necessary features to allow the hardware device to operate properly.

The exact detail of which device driver is required by the operating system is stored in a file and the drivers for the hardware connected are loaded when the computer system is booted.

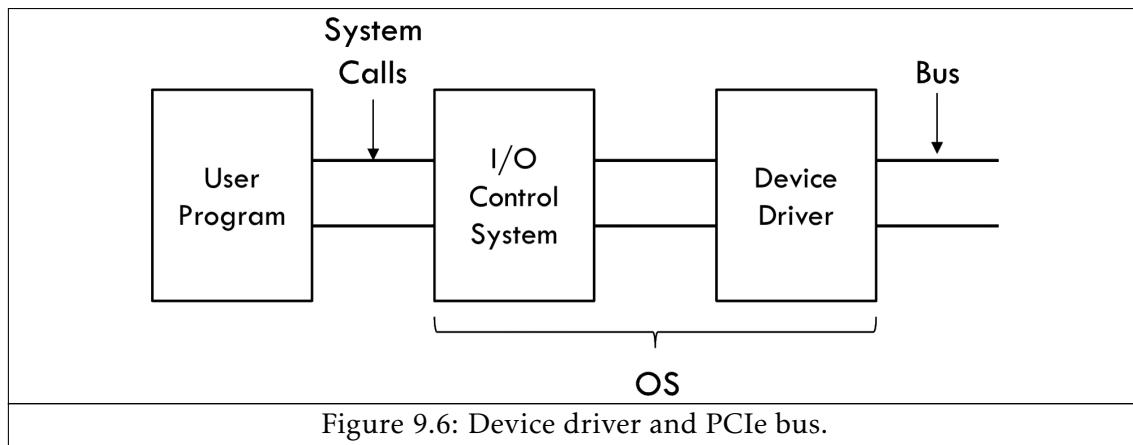
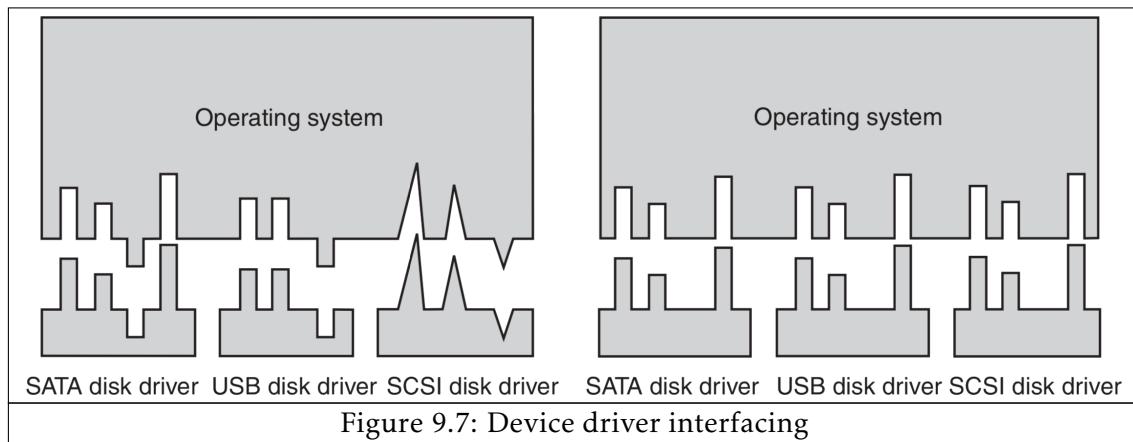


Figure 9.6: Device driver and PCIe bus.

As shown in the figure above, the device driver converts the logical I/O requests into specific commands directed to that device. This is communicated using the bus, to which the PCIe bus is connected.

A problem arises as many different manufacturers create device drivers for their devices.



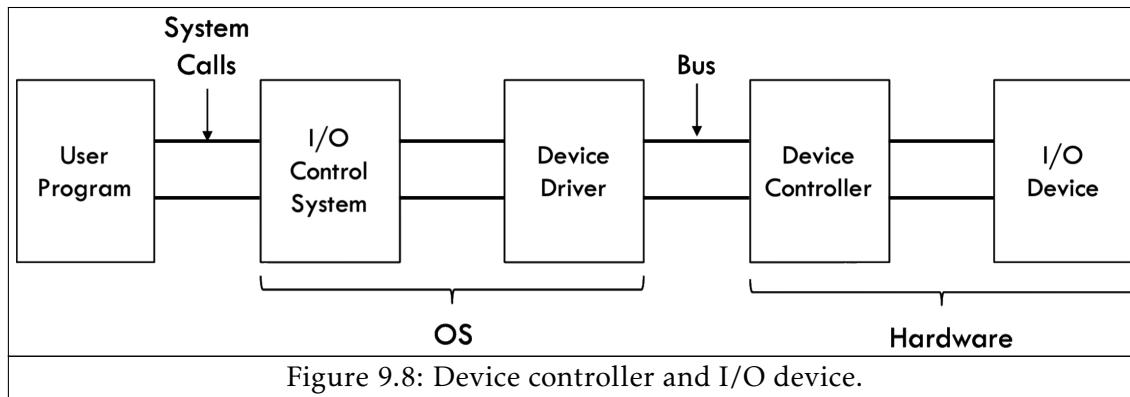
In the figure above, the left image shows a case where a standard driver interface is not used. This means that the operating system would need to be modified to be compatible

with each device driver after a manufacturer develops a new device. This would yield a high amount of maintenance and code required in the operating system.

In the figure above the right image shows a case where a standard driver interface is used. This allows the operating system to expose an application programming interface (API) to device manufacturers. The device manufacturers must adhere to the API and therefore, all drivers follow this API definition so the operating system code does not have to be modified.

Device controller

The device controller is a hardware unit attached to the bus.

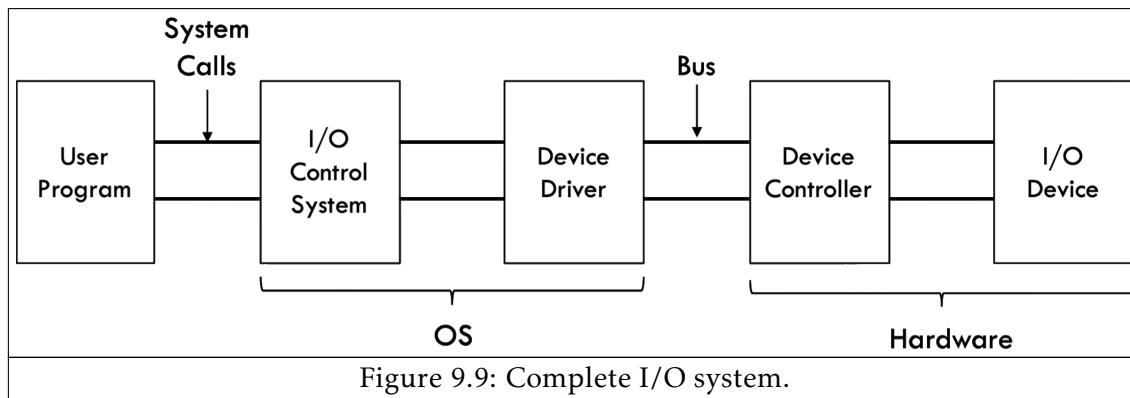


As shown in the figure above, the device controller provides an interface between the computer system and the I/O device using electrical signals.

I/O device

The I/O device may be:

- a **block device** that transfers data in groups of characters; or
- a **character device** that transfers data one character at a time.



The figure above shows the complete I/O system.

Each I/O device has a **device descriptor** containing information about the device. For example:

- the device identification;
- the instructions that operate the device;
- the current status; and
- the current user process (if any).

The information in the device descriptor is used by the I/O controller.

Objectives of the I/O system

The I/O system aims to be **efficient** as:

- I/O devices are much slower than the CPU; and
- I/O devices should be used efficiently as to avoid blocking the CPU.

In addition, the I/O system aims to allow **independence and uniformity** as:

- I/O devices are quite complex and diverse;
- the user should not be concerned with the details of I/O devices;
- the user should be able to treat I/O devices in a uniform way; and
- the user should deal with the virtual devices presented by the OS, rather than the physical devices.

Finally, the I/O system should provide **error handling** that should be:

- handled as close to the hardware as possible; and
- transparent to the user.

I/O interrupts

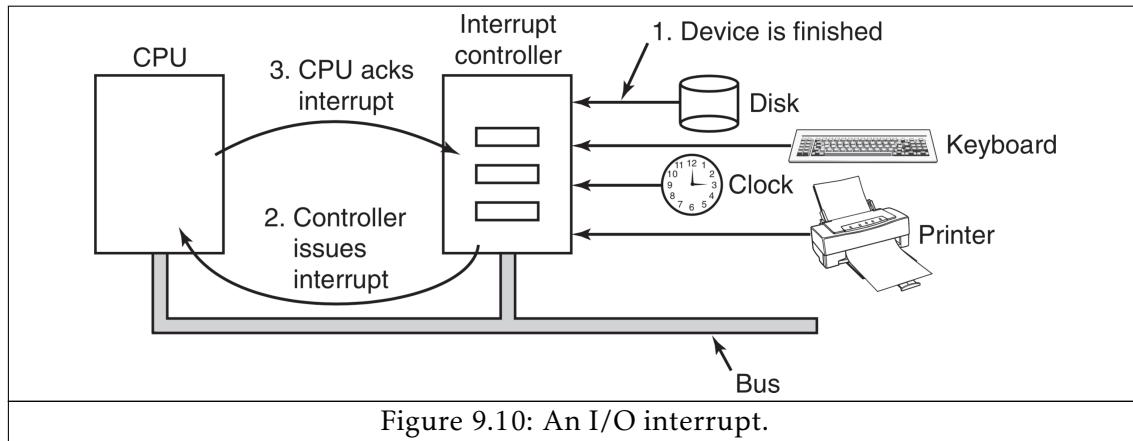
Definition

An **interrupt** is a signal sent to the processor indicating that an event caused by hardware or software requires immediate attention.

In I/O, interrupts are used as a way of controlling I/O activity in which a peripheral or terminal that needs to make or receive a data transfer.

How it works

Interrupts are very important and are not trivial to handle.

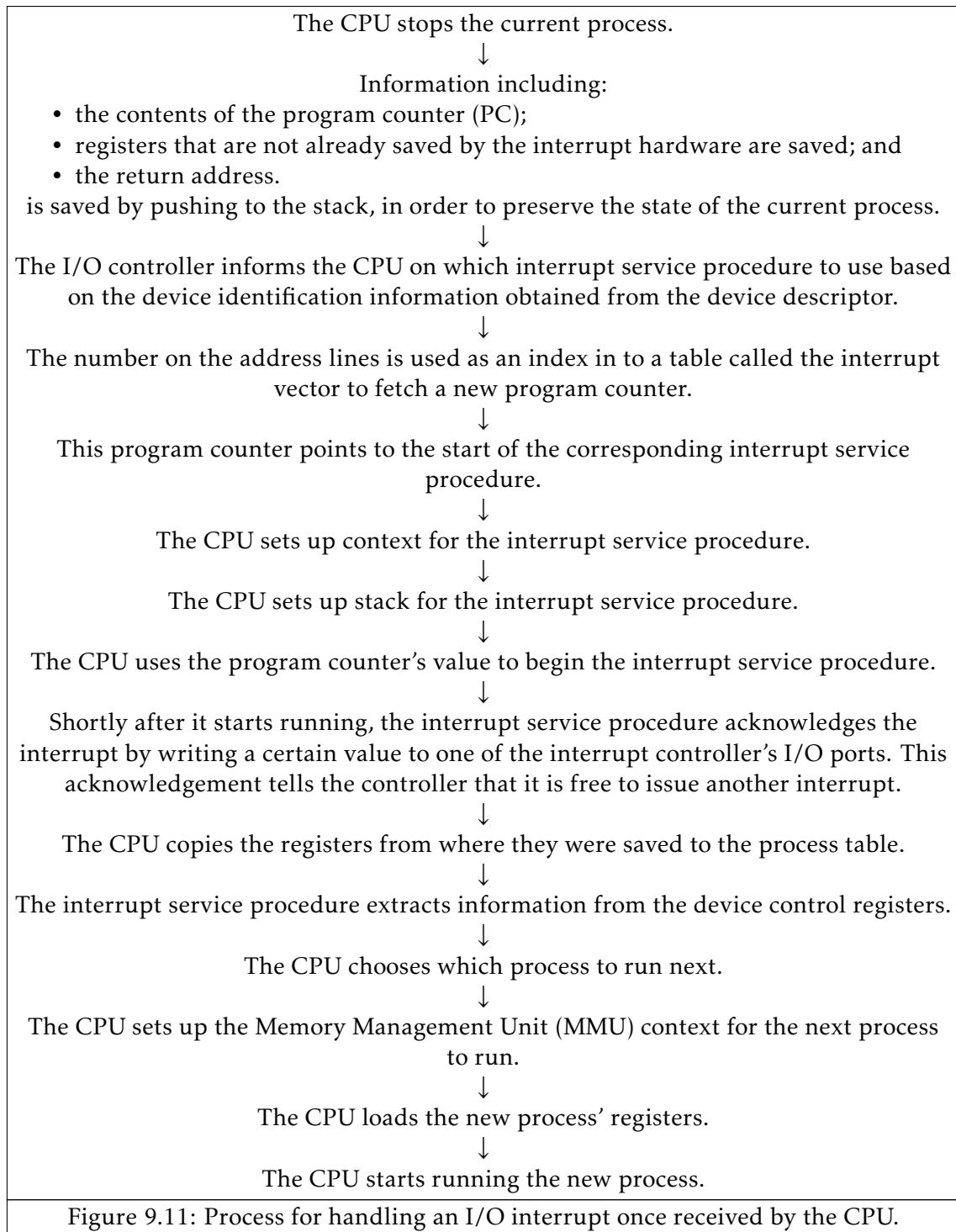


The figure above shows that:

- an I/O interrupt is generated when a particular I/O device is finished with a current process;
- the interrupt controller issues the interrupt to the CPU by asserting a signal on its assigned bus line;
- the CPU's interrupt controller chip detects the signal from the interrupt controller and performs some subsequent actions based on the conditions; and
- the CPU acknowledges the interrupt to the interrupt controller.

The connections between the devices and the controller actually use interrupt lines on the bus rather than dedicated wires.

If no other interrupts are pending, the interrupt controller handles the interrupt immediately. However, if another interrupt is in progress, or another device has made a simultaneous request on a higher-priority interrupt request line on the bus, the device is just ignored for the moment. In this case it continues to assert an interrupt signal on the bus until it is serviced by the CPU.



The figure above shows the process for handling an I/O interrupt once it has been received by the the CPU.

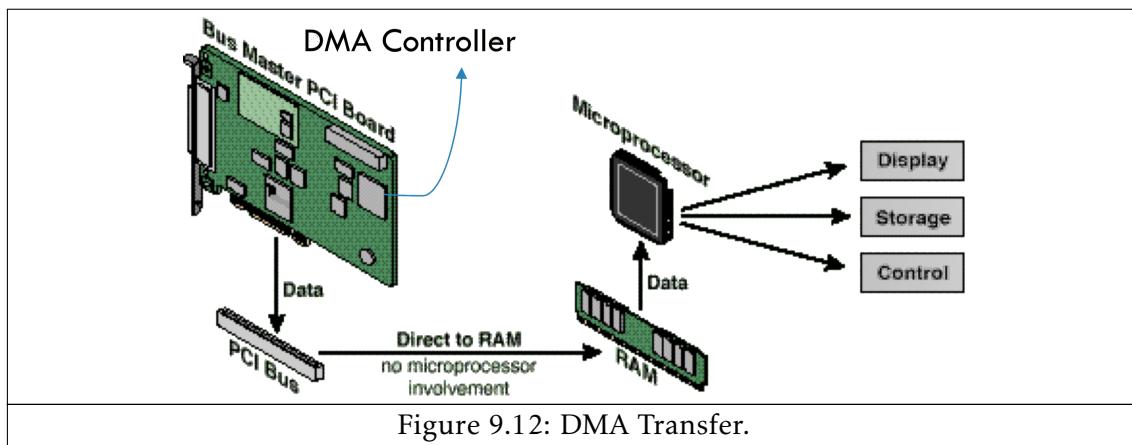
Direct Memory Access (DMA)

Definition

Direct Memory Access (DMA) is an alternative method of handling I/O operations that allows an I/O device to send or receive data directly to or from the main memory, bypassing the CPU to speed up memory operations.

How it works

The computer system and peripherals communicate directly with each other using memory buses, removing intervention of the CPU.

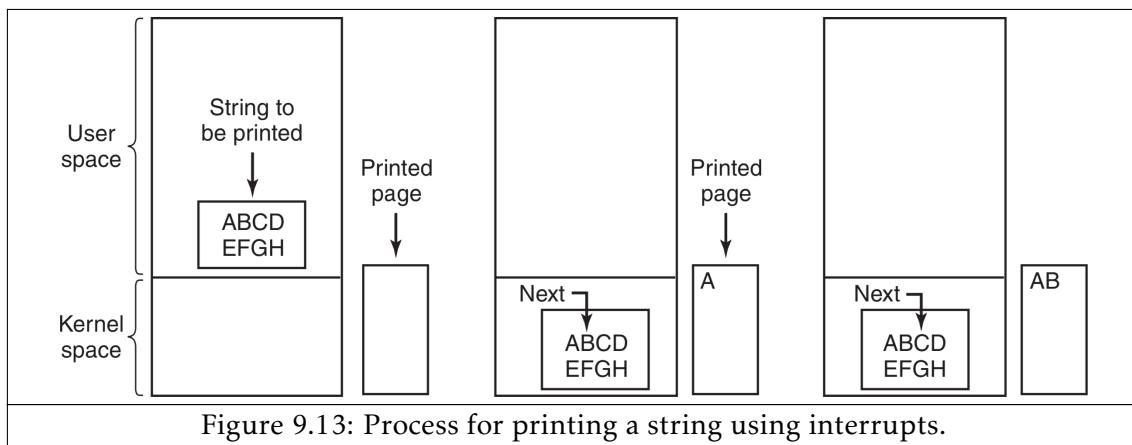


During DMA, the CPU is idle and has no control over the memory buses.

The DMA controller takes over the buses to manage the transfer directly between the I/O devices and the memory unit.

Benefits over interrupts

When considering a case where a user process requests sending an eight character string "ABCDEFGH" to a printer peripheral, an interrupt would be caused for each character sent, as shown in the figure below.



This is not very efficient as the data transfer between the device and memory unit is limited by the speed of the CPU. Whereas, DMA allows **burst transfer** where a block sequence of words in memory is transferred in a continuous burst by the DMA controller.

Buffering

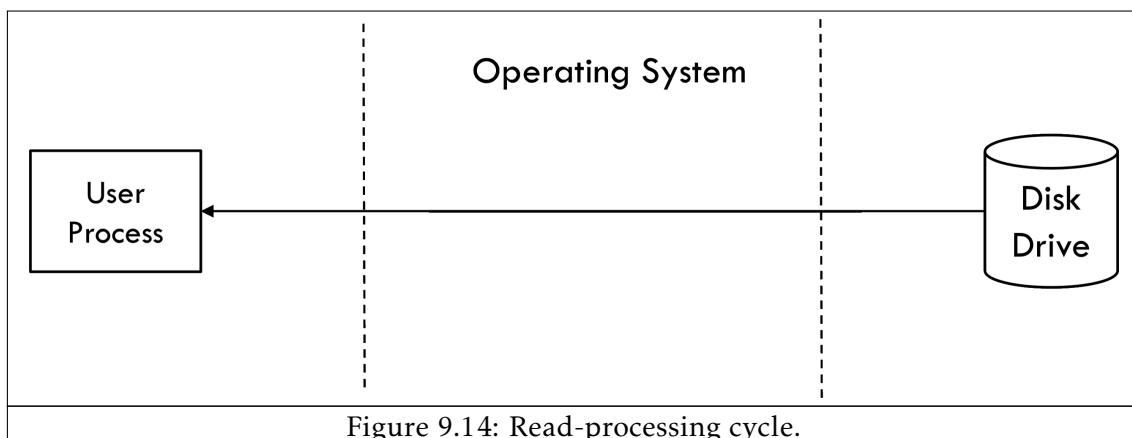
Definition

A **buffer** is an interim memory area that holds data in transit between process's working RAM and the I/O device.

Why is buffering required?

If a process does repeated I/O, then it will be repeatedly waiting for the I/O to complete. For example, in a disk read operation:

- the process issues some system calls;
- the process waits to read the first block of data to memory; and
- the process performs some processing on the data.

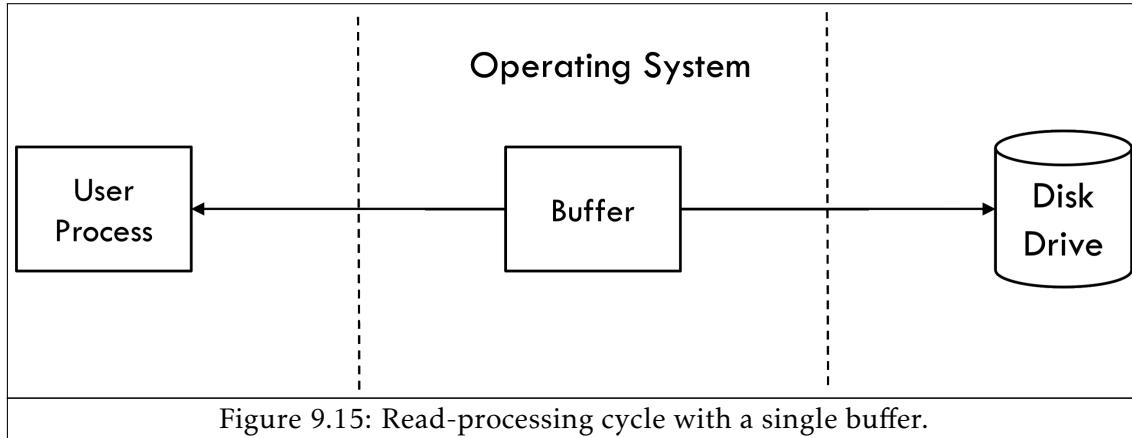


The read-processing cycle shown in the figure above presents some issues as:

- the CPU is idle for most of the time as it is waiting for data transfer to finish and will issue interrupts; and
- the disk drive is not running continuously.

This shows that with no buffering, neither the CPU or disk drive are being used to their maximum capabilities as allowing a process to run many times for short runs is inefficient.

Input and output buffering



As shown in the figure above, the buffer holds the data in transit between process's working RAM and the I/O device in order to:

- cope with the speed mismatch between the devices as the I/O devices are generally slower than the CPU; and
- perform a checksum validation to ensure that the data is valid.

When a user process is attempting to obtain input data from an I/O device, **input buffering** is used where:

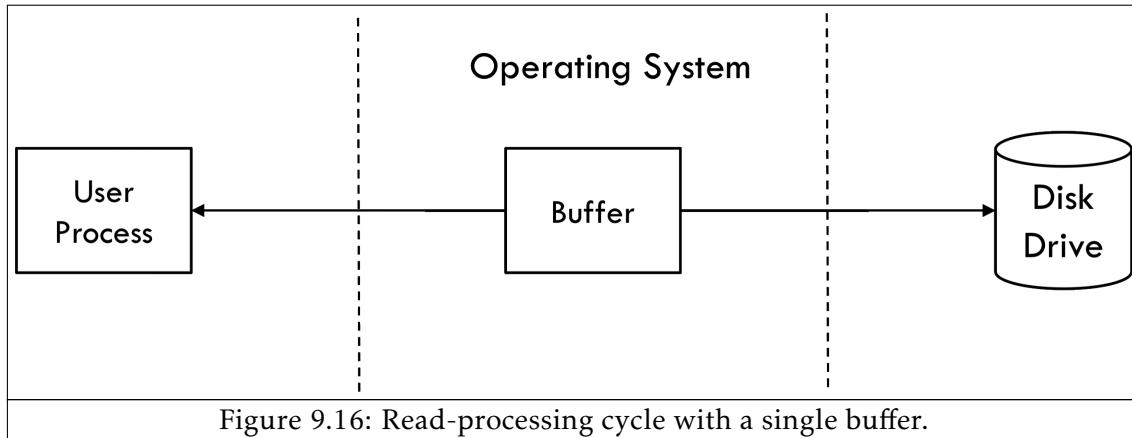
- input data is transferred to an area of memory called the input buffer;
- the process gets data from the input buffer;
- the OS fills the input buffer when it is empty; and
- the process is blocked if it attempts to get data from an empty input buffer.

When a user process is attempting to output data to an I/O device, **output buffering** is used where:

- output data is transferred to an area of memory called the output buffer;
- the OS empties the output buffer when it is full; and
- the process is blocks if it attempts to output before the system has emptied the output buffer.

Single buffering

When using single buffering, blocks are read in to the buffer and then moved to the user's work area.



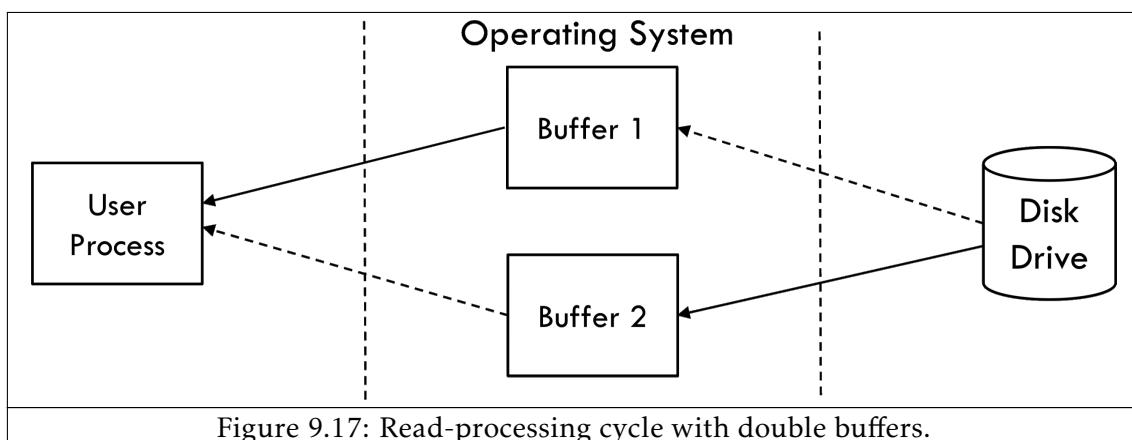
When blocks are moved from the buffer, processing can be completed in parallel with the transfer of the next blocks.

Advantages	Disadvantages
Interrupts are reduced as the CPU is not always waiting for the data transfer to finish.	The I/O device is idle for some time as the CPU must wait for the buffer to be filled.

Double buffering

When using double buffering:

- the process reads/writes to one buffer; while
- the OS empties/fills the other buffer.

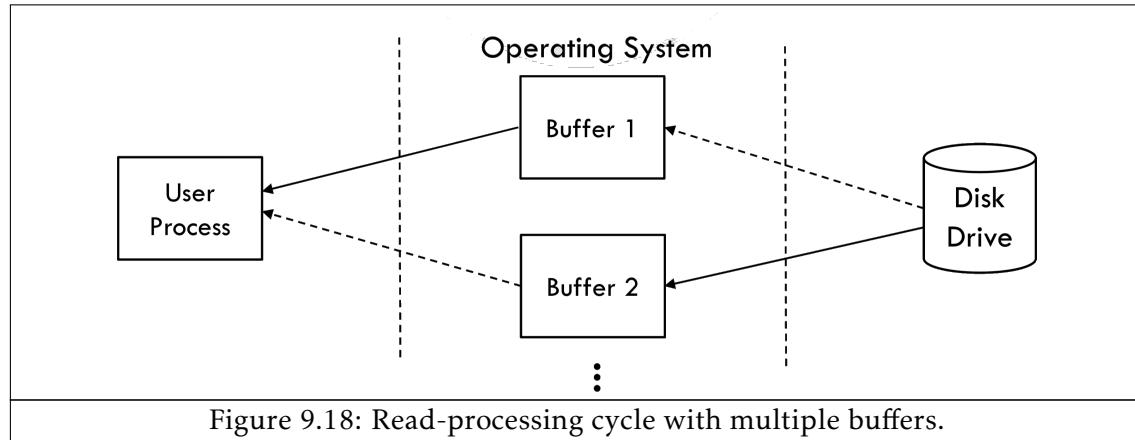


Advantages	Disadvantages
Transferring data may be continuous if the processing time is less than the transferring time.	The I/O device may still be idle for some time if the transferring time is less than the processing time.

Increases utilisation of the I/O device as it is idle for less time than in single buffering.	
---	--

Multiple buffering

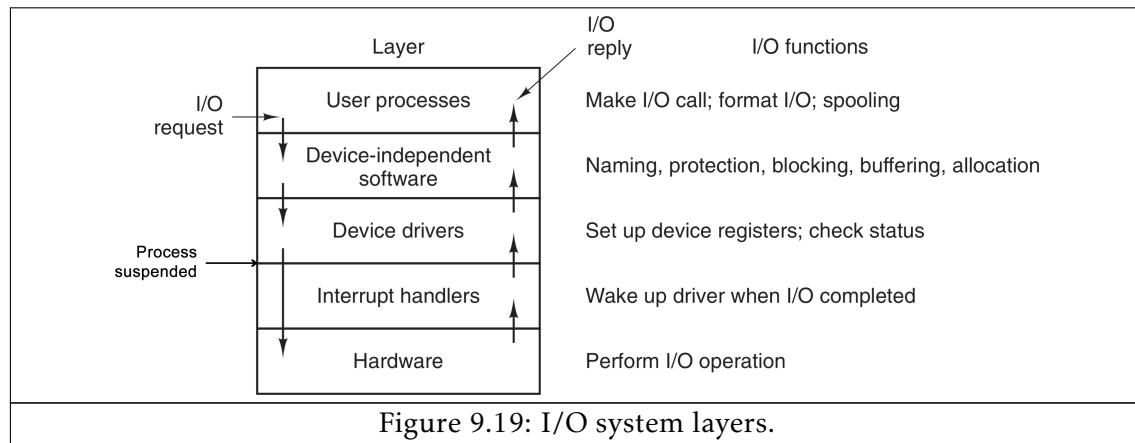
In order to achieve continuous I/O, double buffering can be extended to multiple buffering.



As shown in the figure above, multiple buffering can be implemented with any number, n, buffers in order to achieve continuous I/O.

Advantages	Disadvantages
Can achieve continuous I/O if enough buffers are used.	Adds more complexity and overhead as more buffers must be managed.
Maximum utilisation of the I/O device as there is continuous I/O.	Increased memory usage as more buffers are stored in memory.

Overview of the I/O system



The figure above shows an overview of the I/O system which contains multiple layers where each layer has a main function.

