
Introduction to Computer Architecture

**Neil Willis
and Jon Kerridge**

Introduction to Computer Architecture

**Neil Willis
and Jon Kerridge**

Department of Computer Studies
Sheffield City Polytechnic

Pitman

Preface

The aim of this book is to provide a student text which presents the ideas and concepts of computer architecture in a well organized and straightforward manner. It develops the fundamental concepts from the register transfer level upwards, not concerning itself with the logic level at all.

It does not assume any prerequisite knowledge of computers although we feel that the material is best taught in parallel with a course in assembler level or machine code programming. Indeed, the assembler level programming might well form the practical or tutorial part of such a course. Such practical sessions will certainly reinforce the understanding of Chapters 4, 5, 6, 7 and 8.

The choice of an illustrative computer for a text such as this is a much debated topic. Our basic premise is that to use examples of real machines is ideal and consequently examples are taken throughout the text from a variety of machines such as an IBM 370, PDP-11, ICL 2900, Hewlett Packard 2000 and HP 3000, Prime 750, Harris 800 and Motorola M6800.

These examples will complement the lecturer's illustrations of facilities available on the computer being used for the practical work

However, if it is not possible to illustrate a particular technique on the computer available, Appendix 2 of the text describes a computer system simulator. This is a software package written in Pascal and is available to lecturers who would like to use it, at a nominal materials charge. Some of the problems in the text relate to the use of the simulator although they can obviously be solved on a real computer.

The general structure of the book is as follows. Chapter 1 introduces all the basic concepts of computing systems giving definitions of processors, memory, peripherals, backing store, communications and distributed systems. Because a distributed system is a collection of single processors, this topic is discussed further in Chapter 9 after the details of a single processor system have been described

Chapter 2 introduces descriptions of basic peripheral devices and gives the reader an appreciation of the physical characteristics, relative costs and flexibility of a wide range of peripheral devices.

Chapter 3 describes a range of backing store devices such as magnetic disc, drum, tape, floppy disc, bubble store, charge couple devices, content addressable file store.

Chapter 4 provides a functional description of memory together with a brief description of the technologies used to implement memory. The functional description includes interleaving and associative techniques as well as the Von Neumann structure.

Chapter 5 gives an outline of the structure and requirements of machine code and a review of addressing techniques.

Chapter 6 is a functional description of a processing unit including synchronous and asynchronous operation. Discussion of information flow in the processing unit leads naturally to the requirements for micro instructions and to descriptions of hardware and software micro instruction techniques.

Chapter 7 introduces the need for interrupts and discusses situations in which they may occur. A review of interrupt handling techniques follows.

Chapter 8 is concerned with the transfer of information to and from peripheral devices. As well as input/output to 'close' devices the chapter includes a review of data transmission techniques.

Chapter 9 shows how computers may be linked to form a distributed system, from both the hardware and software aspects.

Chapter 10 describes architectural developments which have been used to increase processor throughput.

Appendix 1 is a brief description of binary, octal, decimal and hexadecimal number systems. It describes the representation of characters and fixed and floating point numbers.

Appendix 2 is a specification of a computer system simulator together with a number of systems which can be used for teaching purposes.

This book has been based on a lecture course forming part of a BSc in Computer Studies and much of the text and many of the problems have been 'class tested'. We are very grateful to many students who have commented on their difficulties and for suggestions they have made for improvements to the material.

We also wish to express our thanks to all the people who have helped during the preparation of this book. In particular we are very grateful to Dan Simpson for his contribution of Chapter 9 and to Doug Bell for his patient reading of the entire text and his detailed and very helpful comments and suggestions.

Finally we wish to thank Mrs J M Kerridge and Mrs M Wilson for their very careful typing.

Neil Willis
Jon Kerridge
Sheffield, 1982

Contents

Preface

1 Basic concepts of computer systems	1
1.1 Computing machines and computer systems	1
1.1.1 What is a computing machine?	1
1.1.2 The difference between computing machines and computer systems	2
1.2 Basic components of a computing machine	3
1.2.1 Central processing unit	3
1.2.2 Memory	4
1.2.3 File store	4
1.2.4 Peripherals	5
1.3 Putting it all together	6
1.3.1 Interconnection	6
1.3.2 Operating systems	6
1.3.3 Distribution	7
1.4 Summary	8
2 Peripherals	9
2.1 Input peripherals	9
2.1.1 Paper tape input	10
2.1.2 Punched cards	11
2.1.3 Key to tape/disk	12
2.1.4 Point of sale (POS) systems	12
2.1.5 Optical character recognition (OCR)	13
2.1.6 Magnetic ink character recognition (MICR)	14
2.2 Output peripherals	15
2.2.1 Line printers	16
2.2.2 Dot matrix printer	16
2.2.3 Daisy-wheel printer	17
2.2.4 Thermal printers	17
2.2.5 Ink jet printers	18
2.2.6 Drum and flatbed plotters	18
2.2.7 Computer output microfilm	18

2.3	Input/output devices	19
2.3.1	Visual display unit and keyboard	19
2.3.2	Typewriter terminals	20
2.4	Television technology	20
2.5	Specific transducers	21
2.6	Summary	22
3	File store	23
3.1	Files records and items	23
3.1.1	Sequential files	24
3.1.2	Indexed-sequential files	25
3.1.3	Random files	25
3.1.4	Filing systems and privacy	26
3.2	Storage media	27
3.2.1	Tape storage	27
3.2.2	Disk storage	29
3.2.3	Bubble memories	31
3.2.4	Charge coupled devices	33
3.2.5	Content addressable file store	34
3.3	Error checking	35
3.4	Summary	36
3.5	Problems	36
4	Memory	38
4.1	The basic building blocks	38
4.2	Dedicated registers	40
4.3	Von Neumann memory organization	41
4.4	Associative memories	42
4.5	Cache memories	45
4.6	Interleaved and segmented memories	46
4.6.1	Segmented memories	48
4.7	Tagged memories	49
4.8	Paged memories	49
4.9	Read only memories	51
4.10	Memory technology	52
4.11	Summary	53
4.12	Problems	54
5	Machine codes and addressing techniques	56
5.1	Instruction formats	56
5.2	Instruction sets	58
5.2.1	Memory reference instructions	58
5.2.2	Branching	60
5.2.3	Subroutines	61
5.2.4	Non memory reference and input/output instructions	63

5.3	Memory addressing	64
5.3.1	Direct (absolute) address	64
5.3.2	Indirect address	64
5.3.3	Immediate address	65
5.3.4	Index addressing	65
5.4	Addressing mechanisms	66
5.4.1	Base and current page addressing mechanism	67
5.4.2	Base addressing mechanism	67
5.4.3	Base and displacement addressing mechanism	69
5.5	Stack computers	70
5.5.1	The organization of a stack	70
5.5.2	Stack instructions	71
5.5.3	Use of stacks	73
5.6	Instruction sets in microprocessors	74
5.7	Summary	74
5.8	Problems	75
6	The central processing unit	79
6.1	Component parts	79
6.1.1	Internal registers	79
6.1.2	The arithmetic and logic unit	81
6.1.3	Floating point unit	82
6.2	Execution of a complete instruction	83
6.2.1	Execution of a non memory reference instruction	83
6.2.2	Execution of a memory reference instruction	83
6.2.3	Branching instructions	85
6.2.4	Indirect and indexed addressing	85
6.3	Architectural considerations	86
6.3.1	Synchronous and asynchronous processors	86
6.3.2	Interconnection methods	88
6.4	Control unit	91
6.4.1	Hardwired control	91
6.4.2	Microprogrammed control	92
6.5	Summary	93
6.6	Problems	93
7	Interrupts	96
7.1	Causes of interrupts	96
7.1.1	Hardware fault	96
7.1.2	Program errors	97
7.1.3	Real time conditions	98
7.1.4	Input/output	98
7.2	Interrupt handling	99
7.2.1	The state of the CPU	99

7.2.2	A simple interrupt system	100
7.3	Identification of interrupt source	101
7.3.1	Identification by software	102
7.3.2	Identification by hardware	103
7.4	Priority systems	104
7.4.1	Nested interrupts	105
7.4.2	Enabling and disabling interrupts	106
7.4.3	Software allocation of priorities	106
7.4.4	Hardware allocation of priorities	107
7.4.5	Interrupt masking	108
7.5	Other interrupt facilities	108
7.6	Examples of real interrupt systems	109
7.6.1	The interrupt mechanism of the IBM 370 series of machines	109
7.6.2	The interrupt mechanism of the Hewlett Packard HP21MX	111
7.6.3	The interrupt mechanism of the Motorola M6800 microprocessor	111
7.7	Summary	111
7.8	Problems	112
8	Data transfers	114
8.1	Input/output bus and the input/output interface	114
8.2	Programmed input/output	116
8.2.1	'Wait for flag' programmed input/output	116
8.2.2	Interrupt programmed input/output	118
8.3	Autonomous input/output—DMA	119
8.4	Input/output channels	121
8.4.1	Selector channels	122
8.4.2	Multiplexer channel	123
8.4.3	Block multiplexer channel	123
8.4.4	Channel programming	124
8.5	Front end processors	126
8.6	Introduction to line transfers	127
8.7	Transmission codes	127
8.7.1	The Baudot code	128
8.7.2	Binary coded decimal	129
8.7.3	The EBCDIC code	129
8.7.4	The ASCII code	129
8.8	Modulation	132
8.9	Transmission techniques	134
8.9.1	Data transmission arrangements	134
8.9.2	Parallel transmission	135
8.9.3	Serial transmission—asynchronous	135

8.9.4	Serial transmission—synchronous	136
8.10	Line sharing	136
8.10.1	Multiplexers	136
8.10.2	Multidrop lines	137
8.11	Line provision	138
8.12	Summary	138
8.13	Problems	139
 Distributed computing 141		
9.1	History	142
9.2	System requirements	142
9.3	Some simple examples	143
9.4	Basic configurations	145
9.5	System design considerations	146
9.6	Hardware	148
9.7	Software	149
9.8	Protocols	149
9.9	Order codes for distributed computing	151
9.10	Costs	152
9.11	Summary	156
9.12	Problems	156
 10 System throughput 158		
10.1	Simple measures of performance	158
10.2	Techniques for increasing throughput	159
10.2.1	Memory bus width	160
10.2.2	Interleaving	160
10.2.3	Cache memory	160
10.2.4	Lookahead processors	164
10.2.5	Pipeline processors	165
10.2.6	Parallel machines	166
10.3	Summary	169
10.4	Problems	169
 Appendix 1 172		
A1.1	Number systems	173
A1.2	Conversion between the number systems	173
A1.2.1	From binary to octal	173
A1.2.2	From binary to hexadecimal	174
A1.2.3	From binary to decimal	174
A1.3	The representation of numeric information	175
A1.3.1	Integer numbers	175
A1.3.1.1	Sign and magnitude	176
A1.3.1.2	Two's complement	176

x Contents

A1.3.2 Floating point numbers	177
A1.4 Problems	179
Appendix 2 A computer system simulator	182
A2.1 Introduction	182
A2.2 The specification of the simulator	183
A2.2.1 Hardware components	183
A2.2.2 Micro-instructions	184
A2.2.3 Output from the simulator	186
A2.3 Error handling	186
A2.4 Implementation	186
A2.5 Use of a simulator in a teaching environment	187
A2.6 Predefined computers	187
A2.6.1 A one address machine with indirect addressing and indexing	188
A2.6.2 A simple stack machine	189
A2.7 Self-defined computers	189
 <i>Index</i>	 191

1 Basic concepts of computer systems

Computer systems are constructed from many different component parts and by connecting these parts together in different ways computer systems with different facilities can be built. However, there is within this variation much that is common. This chapter introduces you to the basic concepts that are common to all computing machines. Having acquired this fundamental knowledge you will then be able to go on and expand upon the framework provided.

1.1 Computing machines and computer systems

The first computing machines were built in the late 1940s. In retrospect they were somewhat limited in that they were designed to solve one particular problem or a small class of problems. Throughout the 1950s and 1960s developments were made which allowed the same computing machine to be used for many different classes of problem. During this time there were technological changes from the valve to the transistor which allowed more complex and more reliable computing machines to be built. It was also during this period that software was developed which opened up the use of the computer to the less specialized users. This was achieved by using programming languages which did not require the user to have intimate knowledge of the construction of the computing machine. There were also developed software systems which made the computing machine more flexible and this gave rise to the term computer system.

During the 1970s there were large developments in the electronics technology industry. These developments enabled the price of computing machines to fall dramatically. Together this meant that it was possible to build computer systems which were constructed from more than one computing machine.

1.1.1 What is a computing machine?

A computing machine is a piece of electronic equipment which, when given some data will process that data in some pre-defined way to produce the required results. It is fundamental to the concept that the operations which

2 Introduction to Computer Architecture

take place do so in a finite time. A computing machine requires some means of accepting data, somewhere to store the sequence of operations to be carried out, some means of performing the calculations required and finally a method of communicating the results. These facilities are common to any computing machine regardless of the task to which it will be put.

1.1.2 The difference between computing machines and computer systems

A computer system takes the basic computing machine and by the addition of software makes the computer system more flexible so that it can be used to solve many different types of problem.

In the simplest case there is a one to one correspondence between a computing machine and a computer system. That is, the computer system uses only one computing machine.

However, the tendency is to develop computer systems that consist of several computing machines. There are two important reasons for this. Firstly, the nature of the environment—may best be supported by providing many localized computing machines—which also provide flexibility by communicating with each other. Secondly, there may be a security aspect which requires certain information held by the computer system to be maintained separate from other parts of the system. Such systems are called distributed systems. An example of such an environment is provided by a medical information system for a hospital. Each ward requires rapid information about patients on that ward but when a patient is moved within the hospital the information relating to that patient should still be available. Also by distributing patient information it is inherently more secure.

These variations of requirement arise, for example, when many computing machines form a computer system but where some of the machines are general purpose and others are dedicated machines.

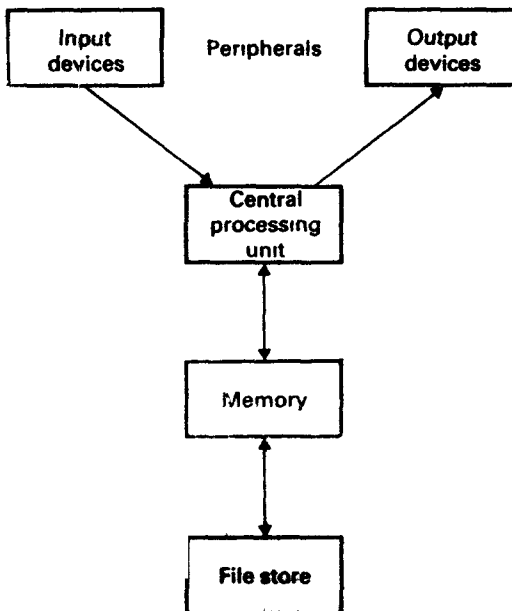
A general purpose computing machine is one which can be used to solve many types of problem. A dedicated computer is one used for one task only: for example, controlling a process in manufacturing industry. In a computer system constructed from such components the general purpose computing machines could be used to extract information maintained by the dedicated machine.

A dedicated computing machine is the same as a general purpose machine in its fundamental construction. It is only dedicated by the connections made between it and the environment in which it operates. An illustration would be a computerized railway system where there will be specialized sensors able to detect the presence of a train and others to check that a set of points have moved as well as the electronically controlled equipment which moves the points. The computer which controls such a system will be dedicated to the task of controlling the

movement of trains and will not be used for any other task. These specialized sensors and the equipment which causes mechanical operations to take place are called transducers.

1.2 Basic components of a computing machine

Computing machines, whether dedicated or general purpose can be represented by the following block diagram (Fig. 1.1). Obviously in



In particular instances some of the components may assume greater importance.

1.2.1 Central processing unit

The central processing unit (CPU) provides central control of the functions of the whole computing machine. In order to function correctly a computing machine requires all necessary data to be made available to the component parts of the machine at exactly the right time. If this is not done the user of the computing machine would find that the results produced by the computing machine would be wrong. As well as providing control

4 Introduction to Computer Architecture

functions the CPU also provides the arithmetic and logic capability of the computing machine.

1.2.2 Memory

The memory of a computing machine is used to hold two distinct types of information: first the data which is to be processed by the computing machine and secondly, the finite sequence of operations which will process that data to produce desired results. The sequence of operations is called an algorithm. The algorithm is said to operate upon a data structure, the data structure being a description of the form of the data rather than the particular data values. The combination of an algorithm and its associated data structure is called a program.

The memory of a computing machine is divided into many similar cells or locations each of which is individually addressable. To say that a cell is addressable means that a CPU can uniquely identify a location within a program so that it can hold a specified data item value. Similarly the address can be used by the algorithmic part of a program to alter the sequence of execution of operations. Instructions are stored in the memory of the computer in sequence, one location after another. Computing machines are extremely good at carrying out repetitive sequences of operations because they are fast and accurate.

This situation is very similar to the way in which knitting patterns are described. The knitter is told the sequence of stitches required to produce a desired effect and then told to repeat that sequence until the desired length has been achieved.

In order that repetitions can be generated in algorithms it is necessary for the sequence to be able to go back over itself, thus generating the necessary repetitions. The addresses of the memory locations are utilized by some of the operations within the computing machine to achieve this effect.

Most memories are unable to distinguish from the contents of a location whether it contains data, algorithm or irrelevant information. It is up to the user of a computer to ensure that appropriate locations of the memory are filled with sensible data and an algorithm before commencing any processing.

1.2.3 File store

The file store holds bulk data and its storage capacity is very large. A single unit of file store may contain ten or a hundred times more locations than the memory itself. The file store can then be constructed from many such units. In order that this large amount of storage space can be more easily managed it is logically split into smaller sections called files. It is then usual

to associate some of the files with a particular program. However, in most computer systems there is nothing to stop any file being processed by any program provided the data is stored in the correct format according to the data structure in the program.

Memory is one of the more expensive components of a computer and it is, therefore, not an infinite resource. Hence, within one computing machine it would not take long before the memory was filled by programs which although required were not currently being executed. The file store provides a means of keeping copies of these programs available to the computing machine in a state of readiness. When such a stored program is required it can be brought from the file store and placed in memory so that it can be executed. On completion, the program will be returned to the file store until required again.

Probably a more important use of file store is for the recording of bulk data. The program describes the data structure upon which it will operate but not the actual data values to be manipulated. These data values have to be stored somewhere and presented to the algorithm as required. In the majority of data processing applications there is far too much data to be processed for it all to be in the memory at the same time. This bulk data is stored in the file store and then transferred to the program as necessary.

It is usual to have a means of identifying users of the computer system so that only those with appropriate authorization can process specific data files with particular programs.

1.2.4 Peripherals

The majority of data processed by a computer together with the algorithms has to be placed in the computer by some method of communication. A peripheral provides this communication with the outside world.

Peripherals can be divided into three classes; input devices, output devices and i/o devices which are capable of communicating in both modes. In a general purpose computer system there will be many different peripherals, so that the user has a degree of flexibility and the use of peripherals is thus tailored to the needs of the application.

In a dedicated computer system it is often the peripherals which make the system dedicated. Up to this stage the computer could have been used for many different applications. The connection of certain peripherals in the form of specific transducers dedicates the computer system. Such dedication is found in process control.

In this case the computer system has to be able to respond rapidly to the requirements of the process being controlled. It would not be sufficient for this data to be transferred by a human operator. Thus it is necessary for the information collected from the parameter being sensed to be entered directly into the computer.

1.3 Putting it all together

Up to this stage it may have appeared that the only place for variation in computing machines and computer systems is in the nature of peripherals and in the amount of memory and file store available. This is not the case and it is the role of the computer architect to design a computing machine so that it will best respond to the environment in which it will be used. To this end, the computer architect has to make a great many decisions about the way in which the basic building blocks will be constructed and then how they will be connected together to form a unified solution for the problem environment. There is no one correct solution for a single environment because each architect has his own ideas and there is usually a financial compromise to be taken into account. Hence, given the same initial building blocks, it is possible to construct many different computing machines, some of which will be more appropriate to particular areas.

We will discuss briefly some of these design criteria at this stage but they will be expanded upon in later chapters.

1.3.1 Interconnection

The architect or designer has initially many ways in which he can connect the component parts of the computer together. By incorporating a large number of specific connections he can make the computing machine process extremely quickly. This approach is, however, costly because it requires a lot of expensive electronic components.

Conversely the components of the computer can be connected together in a far more generalized way which saves on component costs, as not so many are used. In that event the computer system will process information more slowly because it is typical of such systems that the same components are used for many different data transfers. Hence some transfers may have to wait until a previous transfer has been completed. It is thus the job of the architect to choose the right compromise between speed and cost so that the computer system satisfies its design criteria.

1.3.2 Operating systems

Besides producing the computing machine a computer manufacturer usually provides some means of controlling the complete computer system. In very small systems this may be left to the purchaser. The computer system control mechanism is provided by means of a program called the operating system. The operating system is a program always available to the computer which allows the computing machine to be operated at its maximum efficiency. The operating system is related to a particular environment and provides the user with three basic resources namely:

processing time, memory to store algorithms and data, and finally use of peripherals. One measure of operating system efficiency is to see how effectively it manages these resources so as to ensure that each resource is used to the maximum. One way is to keep several programs simultaneously in a partial state of execution within the computer. Each is given a share of the machine in turn. Obviously this requires a very complex operating system. Such an operating system will maintain a list of programs waiting to be executed for the first time. From this list the operating system will select that program which allows it to make most effective use of the resources.

An alternative approach is to initiate a program and then let it execute until processing is completed. This will obviously be wasteful of resources, but the operating system will be much simpler and quicker in operation.

The computer architect has to be aware of the sort of operating system to be used with the computing machine. He has, therefore, to design the computing machine and take account of the operating system. It is possibly true to say that there is more cost involved in the analysis, design and implementation of the operating system than there is of the underlying computing machine. Hence the computing machine and the operating system are normally designed in parallel.

1.3.3 *Distribution*

Increasingly the computer architect has had to become more aware of the problems caused by distributed computer systems.

As explained earlier, a distributed computer system is one which is constructed from several computing machines. Each of the computing machines needs to communicate with its immediate neighbours in the system. So that this may be done in an orderly fashion facilities have to be provided either by programs within the operating system or by the CPUs of the machines. Such facilities are called protocols.

If copies of the same data are stored on several different computing machines then there has to be a mechanism to ensure that all copies have been updated before any user has access to any copy of the data.

Another major problem arises when deciding what facilities of the computer system are to be distributed. The facilities to be considered are processing capability and data storage. The spectrum ranges from a distributed system which has a central file store with many computing machines connected, to a system which has processing and data storage capabilities uniformly distributed.

The computer architect has to ensure that either he constructs computing machines that can operate in all these modes or a range of computing machines each of which is suited to a particular role within the distributed system.