Mr. Abdul Rehman

email id: abdul@northern.edu.pk
(Week 16) Lecture 31-32

Whatsapp# 0308-7792217

Learning objectives:

Review of the last lecture

Modes of data transfer continue....

Direct Memory Access (DMA)

• Input output processor (IOP)

• Memory introduction.

Resources: Beside these lecture handouts, this lesson will draw from the following

Text Book: Computer System Architecture by Morris Mano (3rd Edition) and

Reference book: Computer Architecture, by William Stallings (4th Edition).

Lecture:

Direct Memory Access (DMA)

In the Direct Memory Access (DMA) the interface transfer the data into and out of the memory

unit through the memory bus. The transfer of data between a fast storage device such as

magnetic disk and memory is often limited by the speed of the CPU. Removing the CPU from

the path and letting the peripheral device manage the memory buses directly would improve the

speed of transfer. This transfer technique is called Direct Memory Access (DMA).

During the DMA transfer, the CPU is idle and has no control of the memory buses. A DMA

Controller takes over the buses to manage the transfer directly between the I/O device and

memory.

The CPU may be placed in an idle state in a variety of ways. One common method extensively

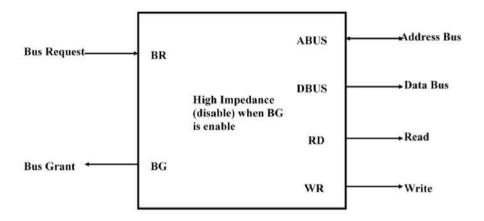
used in microprocessor is to disable the buses through special control signals such as:

• Bus Request (BR)

• Bus Grant (BG)

1

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217
These two control signals in the CPU that facilitates the DMA transfer. The Bus Request (BR) input is used by the DMA controller to request the CPU. When this input is active, the CPU terminates the execution of the current instruction and places the address bus, data bus and read write lines into a high Impedance state. High Impedance state means that the output is disconnected.



The CPU activates the Bus Grant (BG) output to inform the external DMA that the Bus Request (BR) can now take control of the buses to conduct memory transfer without processor.

When the DMA terminates the transfer, it disables the Bus Request (BR) line. The CPU disables the Bus Grant (BG), takes control of the buses and return to its normal operation.

The transfer can be made in several ways that are:

- DMA Burst
- Cycle Stealing

DMA Burst:

In DMA Burst transfer, a block sequence consisting of a number of memory words is transferred in continuous burst while the DMA controller is master of the memory buses.

Cycle Stealing

Cycle stealing allows the DMA controller to transfer one data word at a time, after which it must returns control of the buses to the CPU.

DMA Controller:

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217

The DMA controller needs the usual circuits of an interface to communicate with the CPU and I/O device. The DMA controller has three registers:

- Address Register
- Word Count Register
- Control Register

Address Register

Address Register contains an address to specify the desired location in memory.

Word Count Register

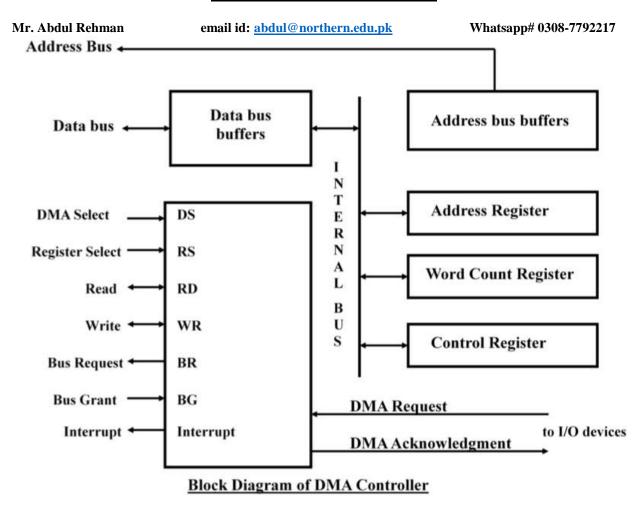
WC holds the number of words to be transferred. The register is incre/decre by one after each word transfer and internally tested for zero.

Control Register

Control Register specifies the mode of transfer

The unit communicates with the CPU via the data bus and control lines. The registers in the DMA are selected by the CPU through the address bus by enabling the DS (DMA select) and RS (Register select) inputs. The RD (read) and WR (write) inputs are bidirectional.

When the BG (Bus Grant) input is 0, the CPU can communicate with the DMA registers through the data bus to read from or write to the DMA registers. When BG =1, the DMA can communicate directly with the memory by specifying an address in the address bus and activating the RD or WR control.



DMA Transfer

The CPU communicates with the DMA through the address and data buses as with any interface unit. The DMA has its own address, which activates the DS and RS lines. The CPU initializes the DMA through the data bus. Once the DMA receives the start control command, it can transfer between the peripheral and the memory.

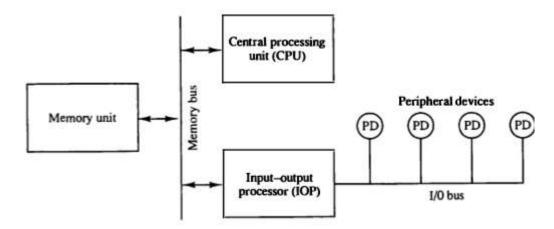
When BG = 0 the RD and WR are input lines allowing the CPU to communicate with the internal DMA registers. When BG=1, the RD and WR are output lines from the DMA controller to the random access memory to specify the read or write operation of data.

Input-Output Processor (IOP)

- It is a processor with direct memory access capability that communicates with IO devices.
- IOP is similar to CPU except that it is designed to handle the details of IO operation.

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217

- Unlike DMA which is initialized by CPU, IOP can fetch and execute its own instructions.
- IOP instruction are specially designed to handle IO operation.

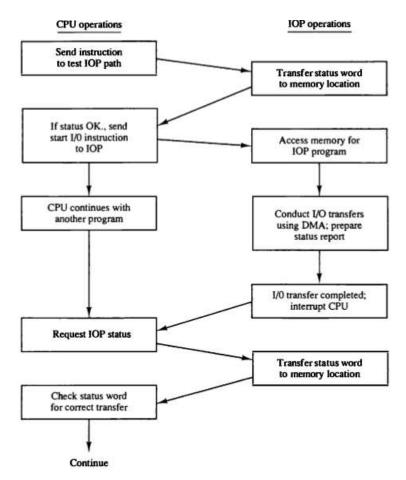


The IOP must structure data words from many different sources. For example, it may be necessary to take four bytes from an input device and pack them into one 32-bit word before the transfer to memory. Data are gathered in the IOP at the device rate and bit capacity while the CPU is executing its own program. After the input data are assembled into a memory word, they are transferred from IOP directly into memory by "stealing" one memory cycle from the CPU. Similarly, an output word transferred from memory to the IOP is directed from the IOP to the output device at the device rate and bit capacity. The communication between the IOP and the devices attached to it is similar to the program control method of transfer. Communication with the memory is similar to the direct memory access method. The way by which the CPU and IOP communicate depends on the level of sophistication included in the system. In very-large-scale computers, each processor is independent of all others and any one processor can initiate an operation. In most computer systems, the CPU is the master while the IOP is a slave processor. The CPU is assigned the task of initiating all operations, but I/O instructions are executed in the IOP. CPU instructions provide operations to start an I/O transfer and also to test I/O status conditions needed for making decisions on various 110 activities. The IOP, in turn, typically asks for CPU attention by means of an interrupt. It also responds to CPU requests by placing a status word in a prescribed location in memory to be examined later by a CPU program. When an I/O operation is desired, the CPU informs the IOP where to find the I/O program and then leaves the transfer details to the IOP.

Mr. Abdul Rehman CPU-IOP Communication email id: abdul@northern.edu.pk

Whatsapp# 0308-7792217

The communication between CPU and IOP may take different forms, depending on the particular computer considered. In most cases the memory unit acts as a message center where each processor leaves information for the other. To appreciate the operation of a typical IOP, we will illustrate by a specific example the method by which the CPU and IOP communicate. This is a simplified example that omits many operating details in order to provide an overview of basic concepts.



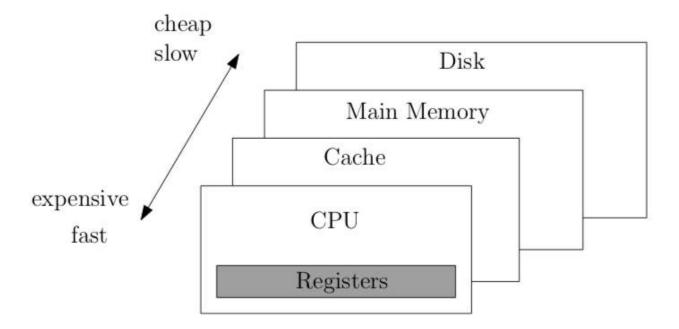
Memory Hierarchy

While studying CPU design in the previous chapters, we considered memory at a high level of abstraction, assuming it was a hardware component that consists of millions of memory cells, which can be individually addressed, for reading or writing, in a reasonable time (i.e., one CPU clock cycle). In this chapter, we will learn that memory is, in fact, built hierarchically, in different layers. This is because the ultimate goals in memory design are to

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217

- have lots of it (gigabytes, terabytes, etc., enough to contain the entire address space),
- make it fast (as fast as CPU registers),
- make it affordable (not too expensive).

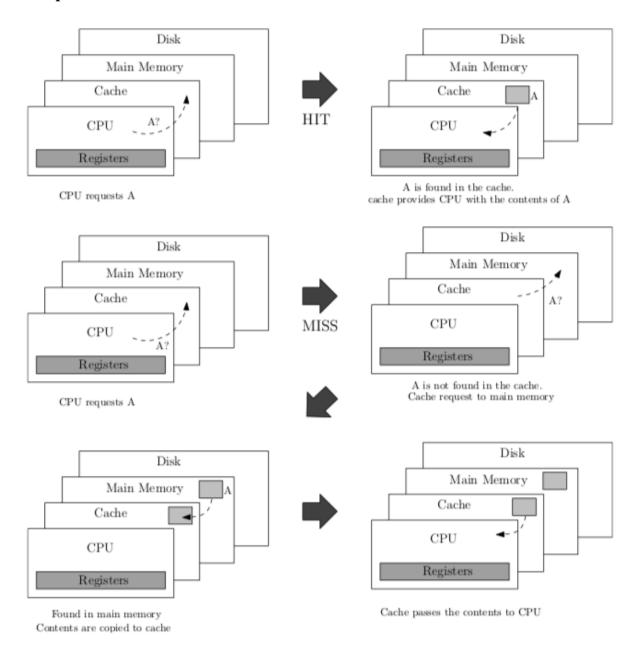
To achieve all of the three design goals, hardware designers combine a small amount of expensive, fast memory and large amounts of inexpensive, slow memory in such a way that the combination of the two behaves as if large amounts of fast memory were available (and that, at an affordable price). To create this illusion of lots of fast memory, we create a hierarchical memory structure, with multiple levels. An example of a structure with 4 levels is shown in Figure. Studying such hierarchical structure in more detail is the topic of this chapter.



In the CPU, registers allow to store 32 words, which can be accessed extremely fast. If information is not present in one of the 32 registers, the CPU will request information from memory, by providing the address of the location where the required information is stored. First, the cache will verify whether it has the requested information available, or not. The cache is located close to the CPU and composed of a relatively small amount of fast and expensive memory. So, if the requested information is available in the cache, it can retrieved quickly. If not, main memory, which is significantly larger and composed of slower and cheaper, is accessed. If the requested information is in the main memory, it is provided to the cache, which

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217 then provides it to the CPU. If not, the hard drive, which contains all information that is stored in the machine, is accessed. The hard drive offers a vast amount of storage space, at an affordable price, however, accessing it is slow. So, fundamentally, the closer to the CPU a level in the memory hierarchy is located, the faster, smaller, and more expensive it is.

Concept of hit and miss



Mr. Abdul Rehman Main Memory email id: abdul@northern.edu.pk

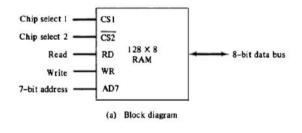
Whatsapp# 0308-7792217

The main memory is the central storage unit in a computer system. It is a relatively large and fast memory used to store programs and data during the computer operation. The principal technology used for the main memory is based on semiconductor integrated circuits. Integrated circuit RAM chips are available in two possible operating modes, **static** and **dynamic**.

The static RAM consists essentially of internal flip-flops that store the binary information. The stored information remains valid as long as power is applied to the unit. The dynamic RAM stores the binary information in the form of electric charges that are applied to capacitors. The capacitors are provided inside the chip by MOS transistors. The stored charge on the capacitors tend to discharge with time and the capacitors must be periodically recharged by refreshing the dynamic memory. Refreshing is done by cycling through the words every few milliseconds to restore the decaying charge. The dynamic RAM offers reduced power consumption and larger storage capacity in a single memory chip. The static RAM is easier to use and has shorter read and write cycles.

RAM and ROM Chips

A RAM chip is better suited for communication with the CPU if it has one or more control inputs that select the chip only when needed. Another common feature is a bidirectional data bus that allows the transfer of data either from memory to CPU during a read operation, or from CPU to memory during a write operation.



CSI	CS2	RD	WR	Memory function	State of data bus				
0	0	×	×	Inhibit	High-impedance				
0	1	×	×	Inhibit	High-impedance				
1	0	0	0	Inhibit	High-impedance				
1	0	0	1	Write	Input data to RAM				
1	0	1	×	Read	Output data from RAM				
1	1	×	×	Inhibit	High-impedance				

(b) Function table

Mr. Abdul Rehman email id: abdul@northern.edu.pk Whatsapp# 0308-7792217

A ROM chip is organized externally in a similar manner. However, since a ROM can only read, the data bus can only be in an output mode. The block diagram of a ROM chip is shown in Fig..

For the same-size chip, it is possible to have more bits of ROM than of RAM, because the internal binary cells in ROM occupy less space than in RAM. For this reason, the diagram specifies a 512-byte ROM, while the RAM has only 128 bytes.

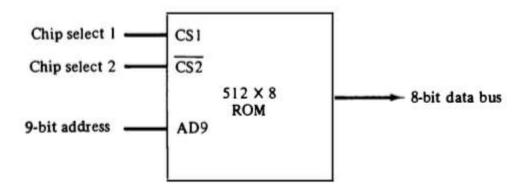


Figure 12-3 Typical ROM chip.

The designer of a computer system must calculate the amount of memory required for the particular application and assign it to either RAM or ROM. The interconnection between memory and processor is then established from knowledge of the size of memory needed and the type of RAM and ROM chips available. The addressing of memory can be established by means of a table that specifies the memory address assigned to each chip. The table, called a memory address map, is a pictorial representation of assigned address space for each chip in the system.

TABLE 12-1 Memory Address Map for Microprocomputer

	Hexadecimal	Address bus									
Component	address	10	9	8	7	6	5	4	3	2	1
RAM 1	0000-007F	0	0	0.	х	х	x	x	х	x	x
RAM 2	0080-00FF	0	0	1	x	x	x	x	x	x	X
RAM 3	0100-017F	0	1	0	x	x	x	x	x	x	X
RAM 4	0180-01FF	0	1	1	x	x	x	x	x	x	х
ROM	0200-03FF	1	x	x	x	x	x	x	x	x	X

email id: abdul@northern.edu.pk Whatsapp# 0308-7792217 Mr. Abdul Rehman CPU Address bus 16-11 9 7-1 RD WR Data bus Decoder CSI CS2 128 × 8 Data RD RAM I WR AD7 CS1 CS2 128×8 Data RD RAM 2 WR AD7 CS1 CS2 128 × 8 RAM 3 Data RD WR AD7 CSI CS2 128 × 8 RAM 4 Data RD WR AD7 CS1 CS2 128×8 Data ROM