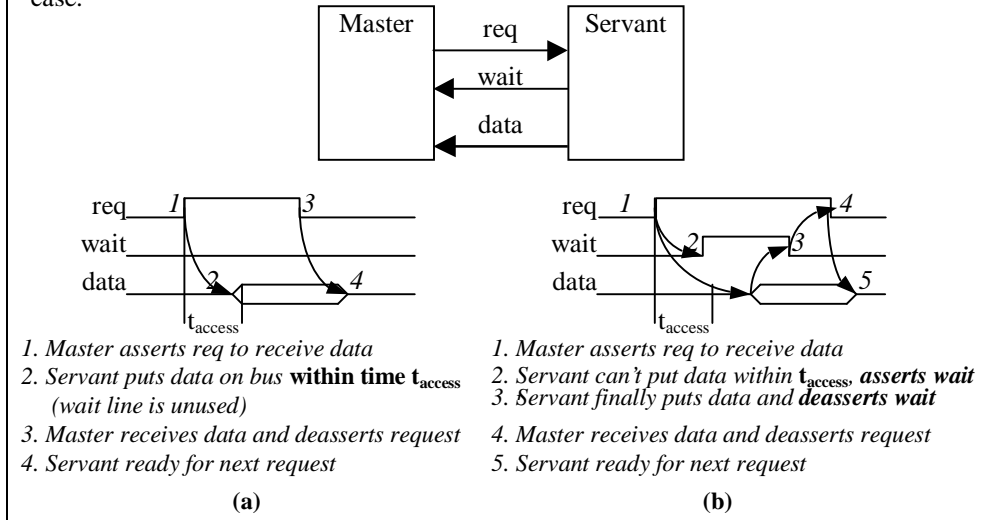


Figure 6.4: A strobe/handshake compromise: (a) fast-response case, (b) slow-response case.



servants) with varying response times, unlike a strobe protocol. However, when response time is known, a handshake protocol may be slower than a strobe protocol, since it requires the master to detect the acknowledgement before getting the data, possibly requiring an extra clock cycle if the master is synchronizing the bus control signals. A handshake also requires an extra line for acknowledgement.

To achieve both the speed of a strobe protocol and the varying response time tolerance of a handshake protocol, a compromise protocol is often used, as illustrated in Figure 6.4. In the case, when the servant can put the data on the bus within time t_{access} , the protocol is identical to a strobe protocol, as shown in Figure 6.4(a). However, if the servant cannot put the data on the bus in time, it instead tells the master to wait longer, by asserting a line we've labeled *wait*. When the servant has finally put the data on the bus, it deasserts the wait line, thus informing the master that the data is ready. The master receives the data and deasserts the request line. Thus, the handshake only occurs if it is necessary. In our boss-employee analogy, the boss tells the employee "I want that report on my desk in an hour; if you can't finish by then, let me know that and then let me know when it's ready."

Example: A simple bus protocol

A protocol for a simple bus (ISA) will be described.

6.4 Interfacing with a general-purpose processor

Perhaps the most common communication situation in embedded systems is the input and output (I/O) of data to and from a general-purpose processor, as it

communicates with its peripherals and memories. I/O is relative to the processor: input means data comes into the processor, while output means data goes out of the processor. We will describe three processor I/O issues: addressing, interrupts, and direct memory access. We'll use the term microprocessor in this section to refer to a general-purpose processor.

6.4.1 I/O addressing

A microprocessor may have tens or hundreds of pins, many of which are control pins, such as a pin for clock input and another input pin for resetting the microprocessor. Many of the other pins are used to communicate data to and from the microprocessor, which we call processor I/O. There are two common methods for using pins to support I/O: ports, and system buses.

A *port* is a set of pins that can be read and written just like any register in the microprocessor; in fact, the port is usually connected to a dedicated register. For example, consider an 8-bit port named P0. A C-language programmer may write to P0 using an instruction like: `P0 = 255`, which would set all 8 pins to 1's. In this case, the C compiler manual would have defined P0 as a special variable that would automatically be mapped to the register P0 during compilation. Conversely, the programmer might read the value of a port P1 being written by some other device, by saying something like `a=P1`. In some microprocessors, each bit of a port can be configured as input or output by writing to a configuration register for the port. For example, P0 might have an associated configuration register called CP0. To set the high-order four bits to input and the low-order four bits to output, we might say: `CP0 = 15`. This writes 00001111 to the CP0 register, where a 0 means input and a 1 means output. Ports are often bit-addressable, meaning that a programmer can read or write specific bits of the port. For example, one might say: `x = P0.2`, giving x the value of the number 2 connection of port P0. Port-based I/O is also called *parallel I/O*.

In contrast to a port, a *system bus* is a set of pins consisting of address pins, data pins, and control pins (for strobing or handshaking). The microprocessor uses the bus to access memory as well as peripherals. We normally consider the access to the peripherals as I/O, but don't normally consider the access to memory as I/O, since the memory is considered more as a part of the microprocessor. A microprocessor may use one of two methods for communication over a system bus: standard I/O or memory-mapped I/O.

In *memory-mapped I/O*, peripherals occupy specific addresses in the existing address space. For example, consider a bus with a 16-bit address. The lower 32K addresses may correspond to memory addresses, while the upper 32K may correspond to I/O addresses.

In *standard I/O* (also known as I/O-mapped I/O), the bus includes an additional pin, which we label M/IO, to indicate whether the access is to memory or to a peripheral (i.e., an I/O device). For example, when M/IO is 0, the address on the address bus corresponds to a memory address. When M/IO is 1, the address corresponds to a peripheral.

Example: HM6264 and 27C256 RAM/ROM memory devices

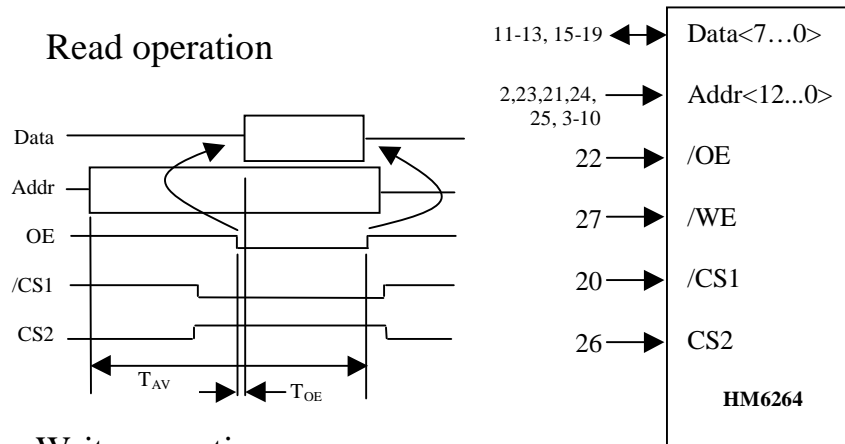
In this example, we introduce a pair of low-cost low-capacity memory devices commonly used in 8-bit micro-controller based embedded systems. The first two numeric digits in these devices indicate whether the device is *random-access memory* (RAM), 62, or *read-only memory* (ROM), 27. Subsequent digits give the memory capacity in K bits. Both these devices are available in 4, 8, 16, 32, and 64K bytes, i.e., part numbers 62/27 followed by 32, 64, 128, 256, or 512. The following table summarizes some of the characteristics of these devices.

| Device | Access Time (ns) | Standby Pwr. (mw) | Active Pwr. (mw) | Vcc Voltage (V) |
|--------|------------------|-------------------|------------------|-----------------|
| HM6264 | 85-100 | .01 | 15 | 5 |
| 27C256 | 90 | .5 | 100 | 5 |

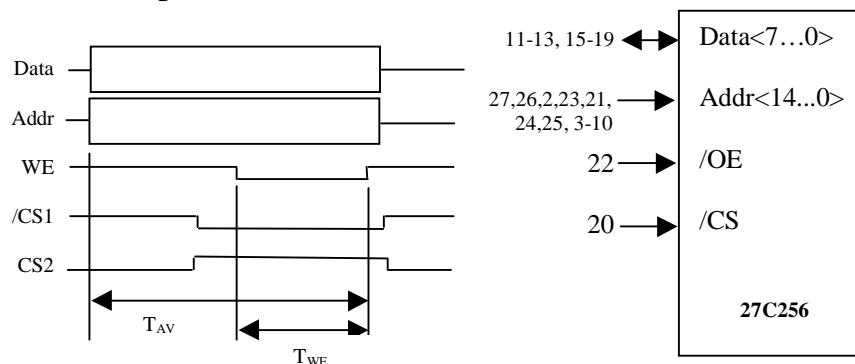
Source Hitachi HM6264B and Microchip 27C256 datasheets.

Memory access to and from these devices is performed through an 8-bit parallel protocol. Placing a memory address on the address-bus and asserting the read signal *output-enable* (OE) performs a read operation. Placing some data and a memory address on the data/address busses and asserting the write signal *write-enable* (WE) performs a write operation. Read and write timing diagrams are presented here.

Read operation



Write operation



In the next example, we demonstrate interfacing of these devices with an Intel 8051 micro-controller.

* $T_{WE}=40\text{ns}$, $T_{OE}=75$, $T_{AV}=40-148$

An advantage of memory-mapped I/O is that the microprocessor need not include special instructions for communicating with peripherals. The microprocessor's assembly instructions involving memory, such as `MOV` or `ADD`, will also work for peripherals. For example, a microprocessor may have an `ADD A, B` instruction that adds the data at address *B* to the data at address *A* and stores the result in *A*. *A* and *B* may correspond to memory locations, or registers in peripherals. In contrast, if the microprocessor uses standard I/O, the microprocessor requires special instructions for reading and writing peripherals. These instructions are often called *IN* and *OUT*. Thus, to perform the same addition of locations *A* and *B* corresponding to peripherals, the following instructions would be necessary:

```
IN R0, A
IN R1, B
ADD R0, R1
OUT A, R0
```

Advantages of standard I/O include no loss of memory addresses to use as I/O addresses, and potentially simpler address decoding logic in peripherals. Address decoding logic can be simplified with standard I/O if we know that there will only be a small number of peripherals, because the peripherals can then ignore high-order address bits. For example, a bus may have a 16-bit address, but we may know there will never be more than 256 I/O addresses required. The peripherals can thus safely ignore the high-order 8 address bits, resulting in smaller and/or faster address comparators in each peripheral.

Situations often arise in which an embedded system requires more ports than available on a particular microprocessor. For example, one may desire 10 ports, while the microprocessor only has 4 ports. An *extended parallel I/O* peripheral can be used to achieve this goal.

Similarly, a system may require parallel I/O but the microprocessor may only have a system bus. In this case, a *parallel I/O* peripheral may be used. The peripheral is connected to the system bus on one side, and has several ports on the other side. The ports are connected to registers inside the peripheral, and the microprocessor can read and write those registers in order to read and write the ports.

Example: memory mapped and standard I/O

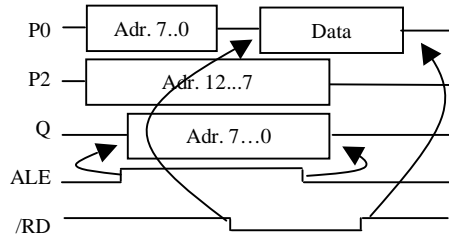
Two examples of real processor I/O to be added.

6.4.2 Interrupts

Another microprocessor I/O issue is that of interrupt-driven I/O. In particular, suppose the program running on a microprocessor must, among other tasks, read and process data from a peripheral whenever that peripheral has new data; such processing is called *servicing*. If the peripheral gets new data at unpredictable intervals, then how can the program determine when the peripheral has new data? The most straightforward approach is to interleave the microprocessor's other tasks with a routine that checks for new data in the peripheral, perhaps by checking for a 1 in a particular bit in a register of

Example: A basic memory protocol

In this example, we illustrate how to interface 8K of data and 32 K of program code memory to a micro-controller, specifically the *Intel 8051*. The Intel 8051 uses separate memory address spaces for data and program code. Data or code address space is limited to 64K, hence, addressable with 16 bits through ports P0 (LSBs) and P2 (MSBs). A separate signal, called PSEN (program strobe enable), is used to distinguish between data/code. For the most part, the I8051 generates all of the necessary signals to perform memory I/O, however, since port P0 is used for both LSB address bits and data flow into and out of the RAM an 8-bit latch is required to perform the necessary multiplexing. The following timing diagram illustrates a memory read operation. Memory write operation is performed in a similar fashion with data flow reversed and RD (read) replaced with WR (write).



Memory read operation proceeds as follows. The micro-controller places the source address, i.e., the memory location to be read, on ports P2 and P0. P2, holding the 8-MSB bits of the address, retains its value throughout the read operation. P1, holding the 8-LSB bits of the address is stored inside an 8-bit latch. The ALE signal (address latch enable), is used to trigger the latching of port P0. Now, the micro-controller asserts high impedance on P0 to allow the memory device to drive it with the requested data. The memory device outputs valid data as long as the RD signal is asserted. Meanwhile, the micro-controller reads the data and de-asserts its control and port signals. The following figure gives the interface schematic.

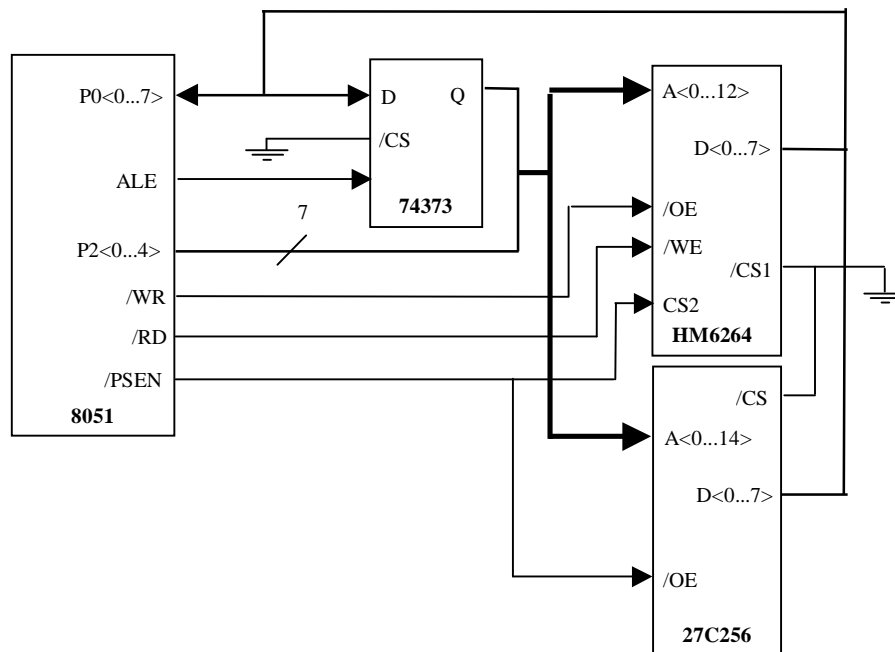
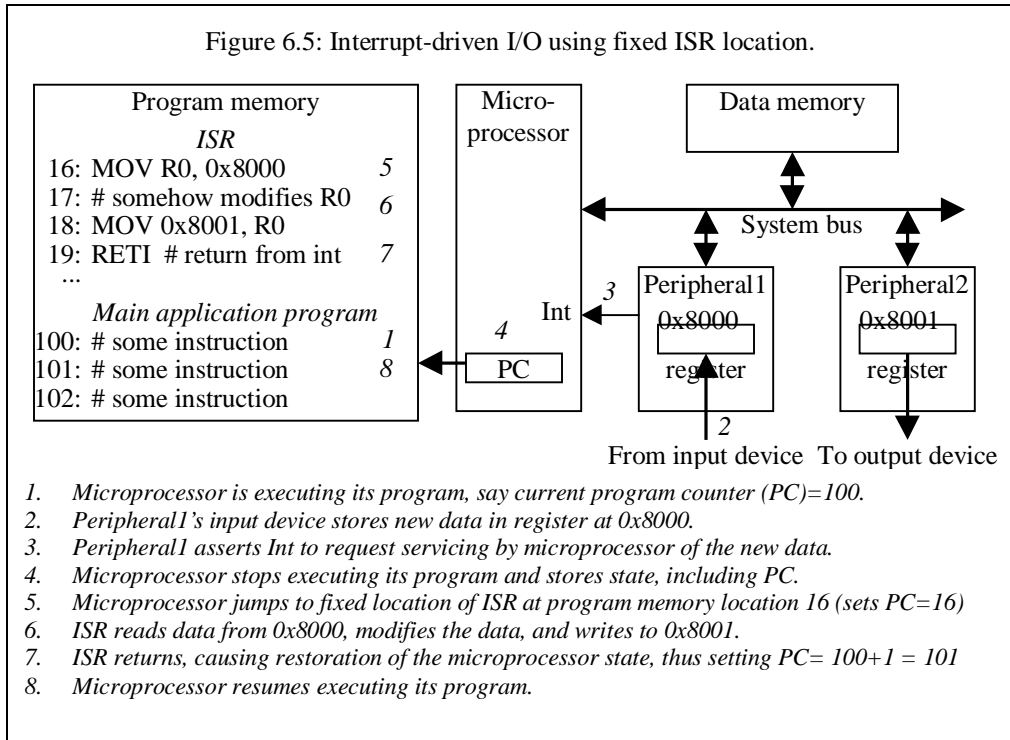


Figure 6.5: Interrupt-driven I/O using fixed ISR location.



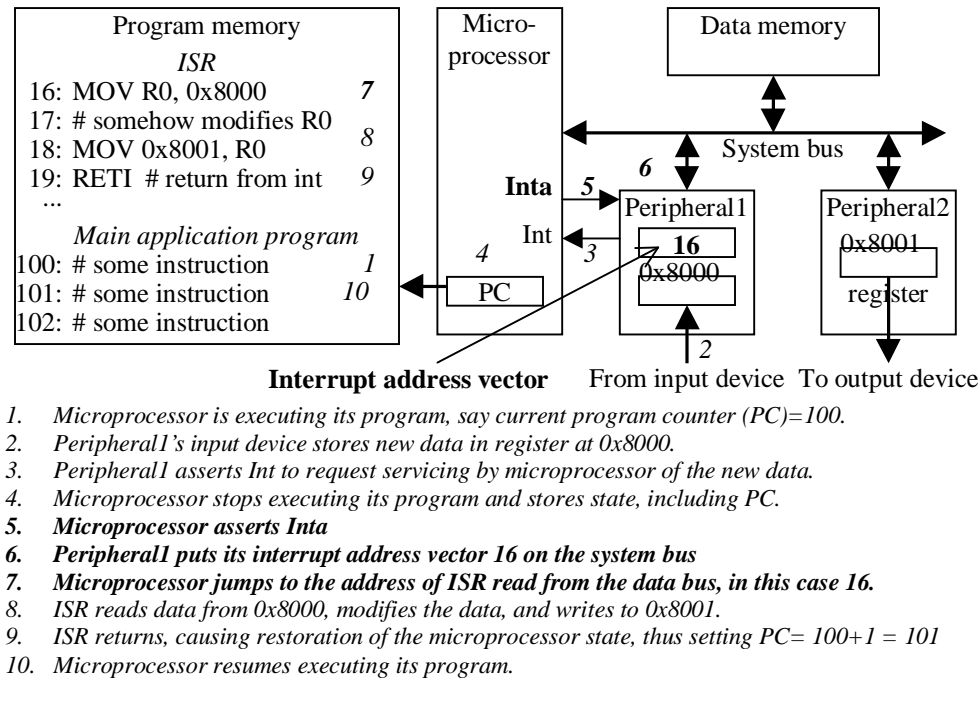
the peripheral. This repeated checking by the microprocessor for data is called *polling*. Polling is simple to implement, but this repeated checking wastes many clock cycles, so may not be acceptable in many cases, especially when there are numerous peripherals to be checked. We could check at less-frequent intervals, but then we may not process the data quickly enough.

To overcome the limitations of polling, most microprocessors come with a feature called *external interrupt*. A microprocessor with this feature has a pin, say *Int*. At the end of executing each machine instruction, the processor's controller checks *Int*. If *Int* is asserted, the microprocessor jumps to a particular address at which a subroutine exists that services the interrupt. This subroutine is called an *Interrupt Service Routine*, or *ISR*. Such I/O is called *interrupt-driven I/O*.

One might wonder if interrupts have really solved the problem with polling, namely of wasting time performing excessive checking, since the interrupt pin is "polled" at the end of every microprocessor instruction. However, in this case, the polling of the pin is built right into the microprocessor's controller hardware, and therefore can be done simultaneously with the execution of an instruction, resulting in no extra clock cycles.

There are two methods by which a microprocessor using interrupts determines the address, known as the *interrupt address vector*, at which the ISR resides. In some processors, the address to which the microprocessor jumps on an interrupt is *fixed*. The assembly programmer either puts the ISR there, or if not enough bytes are available in

Figure 6.6: Interrupt-driven I/O using vectored interrupt.



that region of memory, merely puts a jump to the real ISR there. For C programmers, the compiler typically reserves a special name for the ISR and then compiles a subroutine having that name into the ISR location. In microprocessors with fixed ISR addresses, there may be several interrupt pins to support interrupts from multiple peripherals. For example, Figure 6.5 provides an example of interrupt-driven I/O using a fixed ISR address. In this example, data received by Peripheral1 must be read, transformed, and then written to Peripheral2. Peripheral1 might represent a sensor and Peripheral2 a display. Meanwhile, the microprocessor is running its main program, located in program memory starting at address 100. When Peripheral1 receives data, it asserts *Int* to request that the microprocessor service the data. After the microprocessor completes execution of its current instruction, it stores its state and jumps to the ISR located at the fixed program memory location of 16. The ISR reads the data from Peripheral1, transforms it, and writes the result to Peripheral2. The last ISR instruction is a return from interrupt, causing the microprocessor to restore its state and resume execution of its main program, in this case executing instruction 101.

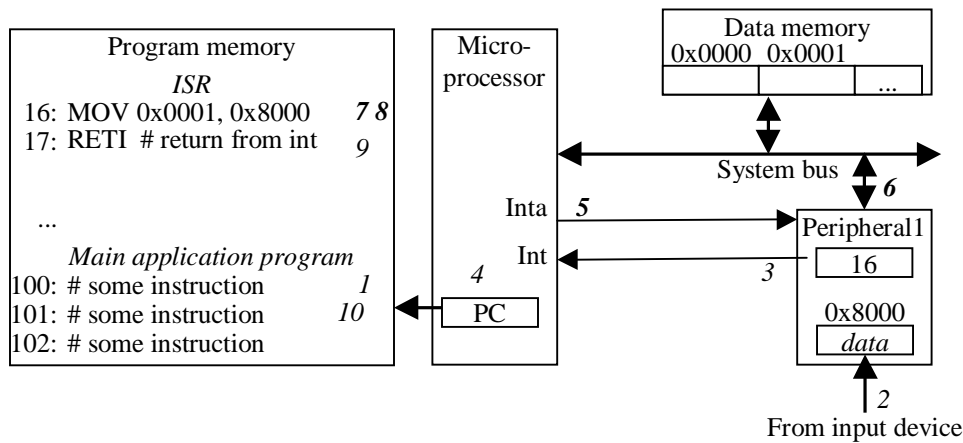
Other microprocessors use *vectored interrupt* to determine the address at which the ISR resides. This approach is especially common in systems with a system bus, since there may be numerous peripherals that can request service. In this method, the

microprocessor has one interrupt pin, say *IntReq*, which any peripheral can assert. After detecting the interrupt, the microprocessor asserts another pin, say *IntAck*, to acknowledge that it has detected the interrupt and to request that the interrupting peripheral provide the address where the relevant ISR resides. The peripheral provides this address on the data bus, and the microprocessor reads the address and jumps to the ISR. We discuss the situation where multiple peripherals simultaneously request servicing in a later section on arbitration. For now, consider an example of one peripheral using vectored interrupt, shown in Figure 6.6, which represents an example very similar to the previous one. In this case, however, the ISR location is not fixed at 16. Thus, Peripheral1 contains an extra register holding the ISR location. After detecting the interrupt and saving its state, the microprocessor asserts *Inta* in order to get Peripheral1 to place 16 on the data bus. The microprocessor reads this 16 into the PC, and then jumps to the ISR, which executes and completes in the same manner as the earlier example.

As a compromise between the fixed and vectored interrupt methods, we can use an *interrupt address table*. In this method, we still have only one interrupt pin on the processor, but we also create in the processor's memory a table that holds ISR addresses. A typical table might have 256 entries. A peripheral, rather than providing the ISR address, instead provides a number corresponding to an entry in the table. The processor reads this entry number from the bus, and then reads the corresponding table entry to obtain the ISR address. Compared to the entire memory, the table is typically very small, so an entry number's bit encoding is small. This small bit encoding is especially important when the data bus is not wide enough to hold a complete ISR address. Furthermore, this approach allows us to assign each peripheral a unique number independent of ISR locations, meaning that we could move the ISR location without having to change anything in the peripheral.

External interrupts may be *maskable* or *non-maskable*. In maskable interrupt, the programmer may force the microprocessor to ignore the interrupt pin, either by executing a specific instruction to disable the interrupt or by setting bits in an interrupt configuration register. A situation where a programmer might want to mask interrupts is when there exists time-critical regions of code, such as a routine that generates a pulse of a certain duration. The programmer may include an instruction that disables interrupts at the beginning of the routine, and another instruction re-enabling interrupts at the end of the routine. Non-maskable interrupt cannot be masked by the programmer. It requires a pin distinct from maskable interrupts. It is typically used for very drastic situations, such as power failure. In this case, if power is failing, a non-maskable interrupt can cause a jump to a subroutine that stores critical data in non-volatile memory, before power is completely gone.

In some microprocessors, the jump to an ISR is handled just like the jump to any other subroutine, meaning that the state of the microprocessor is stored on a stack, including contents of the program counter, datapath status register, and all other registers, and then restored upon completion of the ISR. In other microprocessors, only a few registers are stored, like just the program counter and status registers. The assembly programmer must be aware of what registers have been stored, so as not to overwrite non-stored register data with the ISR. These microprocessors need two types of assembly instructions for subroutine return. A regular return instruction returns from a regular

Figure 6.7: Peripheral to memory transfer *without* DMA.

1. Microprocessor is executing its program, say current program counter (PC)=100.
2. Peripheral 1's input device stores new data in register at 0x8000.
3. Peripheral 1 asserts *Int* to request servicing by microprocessor of the new data.
4. Microprocessor stops executing its program and stores its state, including PC.
5. **Microprocessor asserts *Inta***
6. **Peripheral 1 puts its interrupt address vector 16 on the system bus.**
7. **Microprocessor jumps to the address of ISR read from the data bus, in this case 16.**
8. ISR reads data from 0x8000 and writes that data to 0x0001 in data memory.
9. ISR returns, causing restoration of the microprocessor state, thus setting PC= 100+1 = 101, and terminating handshake with peripheral.
10. Microprocessor resumes executing its program.

subroutine, which was called using a subroutine call instruction. A return from interrupt instruction returns from an ISR, which was jumped to not by a call instruction but by the hardware itself, and which restores only those registers that were stored at the beginning of the interrupt. The C programmer is freed from having to worry about such considerations, as the C compiler handles them.

The reason we used the term "external interrupt" is to distinguish this type of interrupt from internal interrupts, also called traps. An internal interrupt results from an exceptional condition, such as divide-by-0, or execution of an invalid opcode. Internal interrupts, like external ones, result in a jump to an ISR. A third type of interrupt, called software interrupts, can be initiated by executing a special assembly instruction.

Example: Interrupts

Two examples of real processor interrupt handling to be added, one using fixed interrupt, the other vectored.

6.4.3 Direct memory access

Commonly, the data being accumulated in a peripheral should be first stored in memory before being processed by a program running on the microprocessor. Such temporary storage to await processing is called buffering. For example, packet-data from an Ethernet card is stored in main memory and is later processed by the different software layers (such as IP stacks). We could write a simple interrupt service routine on the microprocessor, such that the peripheral device would interrupt the microprocessor whenever it had data to be stored in memory. The ISR would simply transfer data from the peripheral to the memory, and then resume running its application. For example, Figure 6.7 shows an example in which Peripheral1 interrupts the microprocessor when receiving new data. The microprocessor jumps to ISR location 16, which moves the data from 0x8000 in the peripheral to 0x0001 in memory. Then the ISR returns. However, recall that jumping to an ISR requires the microprocessor to store its state (i.e., register contents), and then to restore its state when returning from the ISR. This storing and restoring of the state may consume many clock cycles, and is thus somewhat inefficient. Furthermore, the microprocessor cannot execute its regular program while moving the data, resulting in further inefficiency.

The I/O method of direct memory access (DMA) eliminates these inefficiencies. In DMA, we use a separate single-purpose processor, called a DMA controller, whose sole purpose is to transfer data between memories and peripherals. Briefly, the peripheral requests servicing from the DMA controller, which then requests control of the system bus from the microprocessor. The microprocessor merely needs to relinquish control of the bus to the DMA controller. The microprocessor does not need to jump to an ISR, and thus the overhead of storing and restoring the microprocessor state is eliminated. Furthermore, the microprocessor can execute its regular program while the DMA controller has bus control, as long as that regular program doesn't require use of the bus (at which point the microprocessor would then have to wait for the DMA to complete). A system with a separate bus between the microprocessor and cache may be able to execute for some time from the cache while the DMA takes place.

We set up a system for DMA as follows. As shown in Figure 6.8, we connect the peripheral to the DMA controller rather than the microprocessor. Note that the peripheral does not recognize any difference between being connected to a DMA controller device or a microprocessor device; all it knows is that it asserts a request signal on the device, and then that device services the peripheral's request. We connect the DMA controller to two special pins of the microprocessor. One pin, which we'll call Hreq (bus Hold REQuest), is used by the DMA controller to request control of the bus. The other pin, which we'll call Hlda (HoLD Acknowledge), is used by the microprocessor to acknowledge to the DMA controller that bus control has been granted. Thus, unlike the peripheral, the microprocessor must be specially designed with these two pins in order to support DMA. The DMA controller also connects to all the system bus signals, including address, data, and control lines.

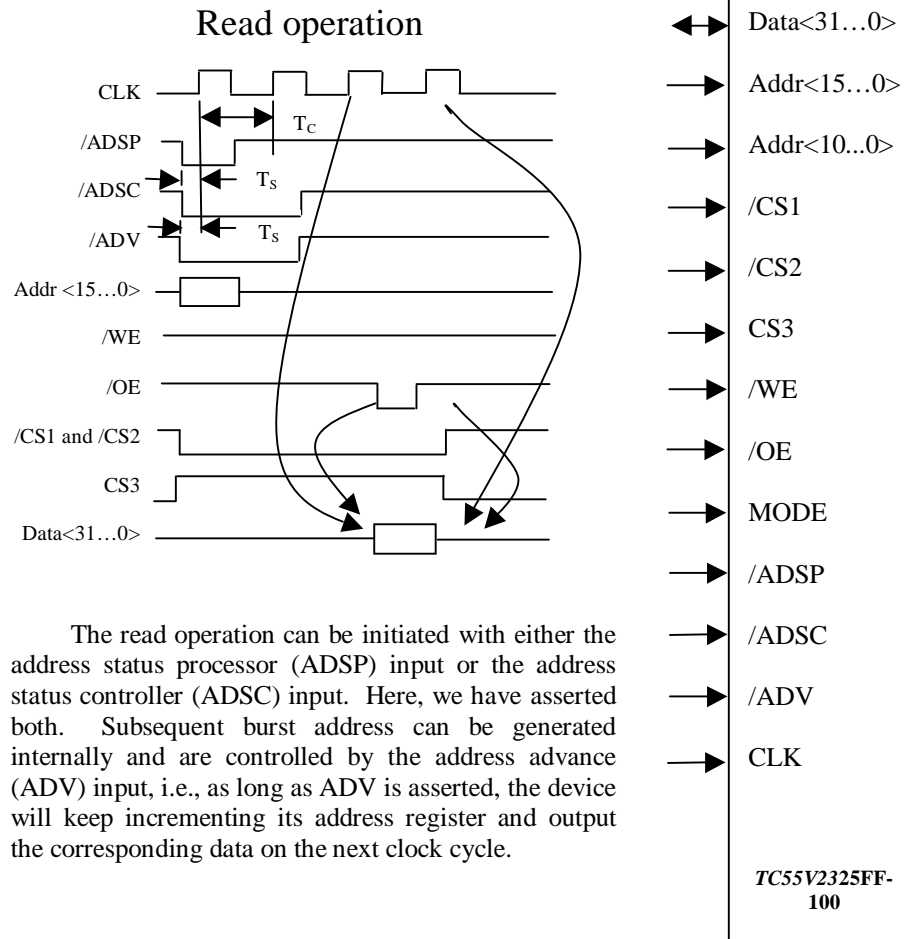
Example: TC55V2325FF-100 memory device

In this example, we introduce a 2M bit synchronous pipelined burst SRAM memory device designed to be interfaced with 32 bit processors. This device, made by Toshiba Inc., is organized as $64K \times 32$ bits. The following table summarizes some of its characteristics.

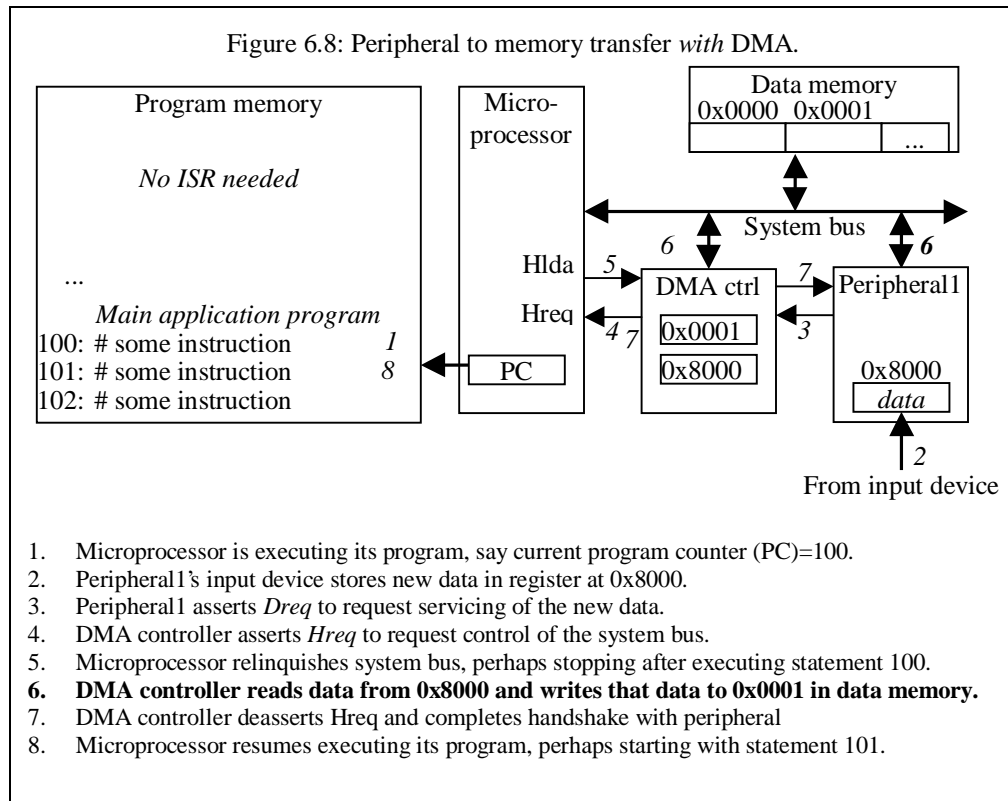
| Device | Access Time (ns) | Standby Pwr. (mw) | Active Pwr. (mw) | Vcc Voltage (V) |
|-----------------|------------------|-------------------|------------------|-----------------|
| TC55V2325FF-100 | 10 | na | 1200 | 3.3 |

Source Toshiba TC59S6432CFT datasheets.

Here, we present the block and timing diagram for a single read operations. Write operation is similar. This device is capable of sequential, fast, reads and writes as well as singles byte I/O. Interested reader should refer to the manufacturer's datasheets for complete pinout and complete timing diagrams.



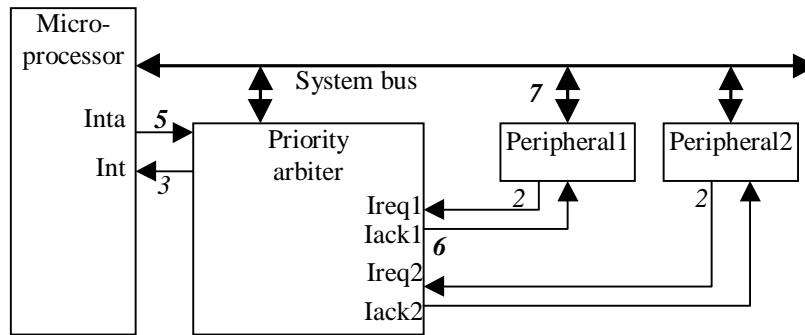
* $T_C = 10\text{ns}$, $T_S = 3\text{ns}$



To achieve the above, we must have configured the DMA controller to know what addresses to access in the peripheral and the memory. Such setting of addresses may be done by a routine running on the microprocessor during system initialization. In particular, during initialization, the microprocessor writes to configuration registers in the DMA controller just as it would write to any other peripheral's registers. Alternatively, in an embedded system that is guaranteed not to change, we can hardcode the addresses directly into the DMA controller. In the example of Figure 6.8, we see two registers in the DMA controller holding the peripheral register address and the memory address.

During its control of the system bus, the DMA controller might transfer just one piece of data, but more commonly will transfer numerous pieces of data (called a block), one right after other, before relinquishing the bus. This is because many peripherals, such as any peripheral that deals with storage devices (like CD-ROM players or disk controllers) or that deals with network communication, send and receive data in large blocks. For example, a particular disk controller peripheral might read data in blocks of 128 words and store this data in 128 internal registers, after which the peripheral requests servicing, i.e., requests that this data be buffered in memory. The DMA controller gains

Figure 6.9: Arbitration using a priority arbiter.



1. Microprocessor is executing its program.
2. Peripheral1 needs servicing so asserts *Ireq1*. Peripheral2 also needs servicing so asserts *Ireq2*.
3. Priority arbiter sees at least one *Ireq* input asserted, so asserts *Int*.
4. Microprocessor stops executing its program and stores its state.
5. **Microprocessor asserts *Inta*.**
6. **Priority arbiter asserts *Iack1* to acknowledge Peripheral1.**
7. **Peripheral1 puts its interrupt address vector on the system bus**
8. **Microprocessor jumps to the address of ISR read from data bus, ISR executes and returns (and completes handshake with arbiter).**
9. Microprocessor resumes executing its program.

control of the bus, makes 128 peripheral reads and memory writes, and only then relinquishes the bus. We must therefore configure the DMA controller to operate in either single transfer mode or block transfer mode. For block transfer mode, we must configure a base address as well as the number of words in a block.

DMA controllers typically come with numerous channels. Each channel supports one peripheral. Each channel has its own set of configuration registers.

Example: Intel 8237

Description of the Intel 8237 DMA controller and an example of its use.

6.5 Arbitration

In our discussions above, several situations existed in which multiple peripherals might request service from a single resource. For example, multiple peripherals might share a single microprocessor that services their interrupt requests. As another example,