

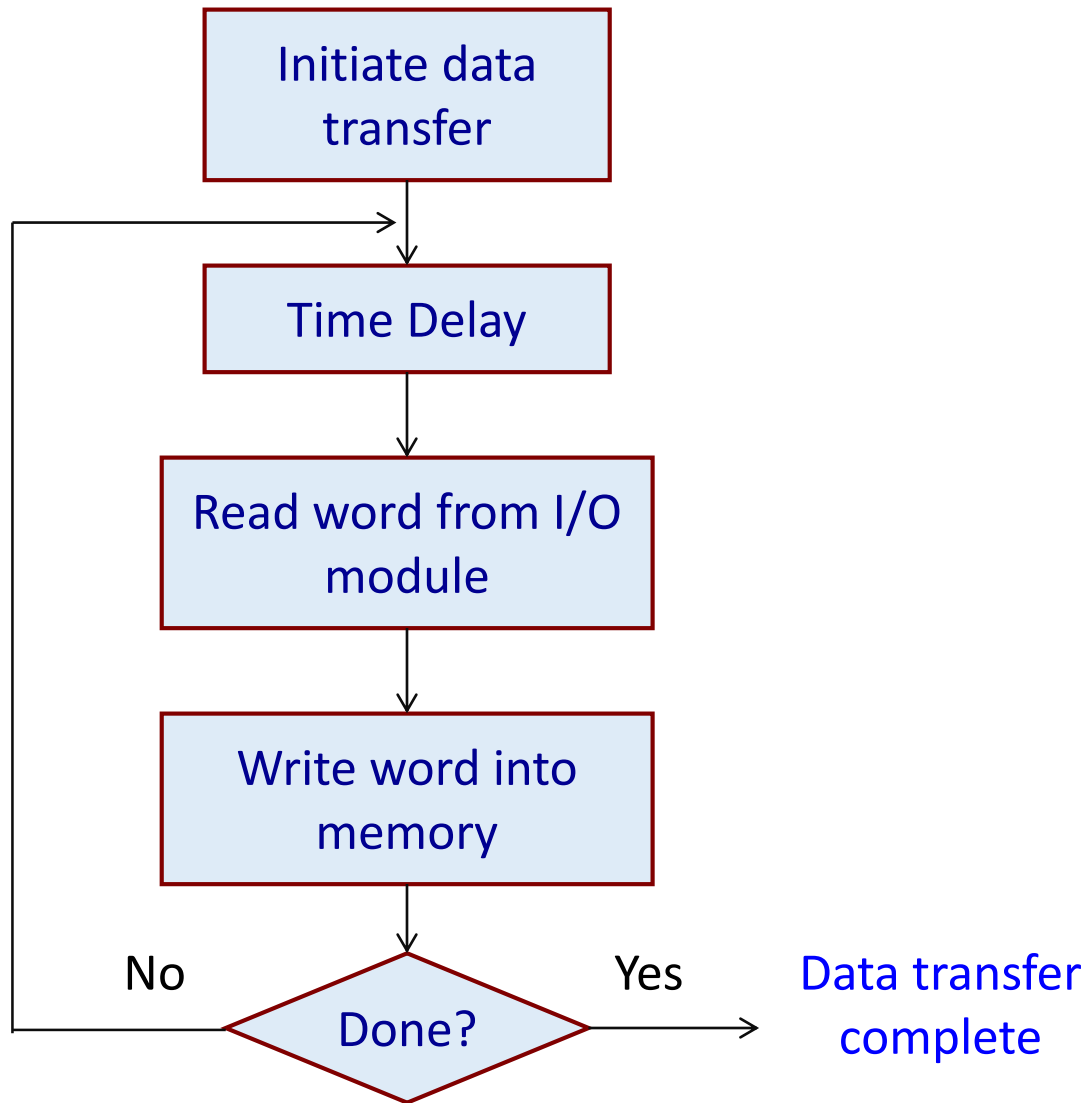
Data Transfer Techniques

Data Transfer Techniques

- 1) Programmed:** CPU executes a program that transfers data between I/O device and memory.
 - a) Synchronous
 - b) Asynchronous
 - c) Interrupt-driven
- 2) Direct Memory Access (DMA):** An external controller directly transfers data between I/O device and memory without CPU intervention.

(a) Synchronous Data Transfer

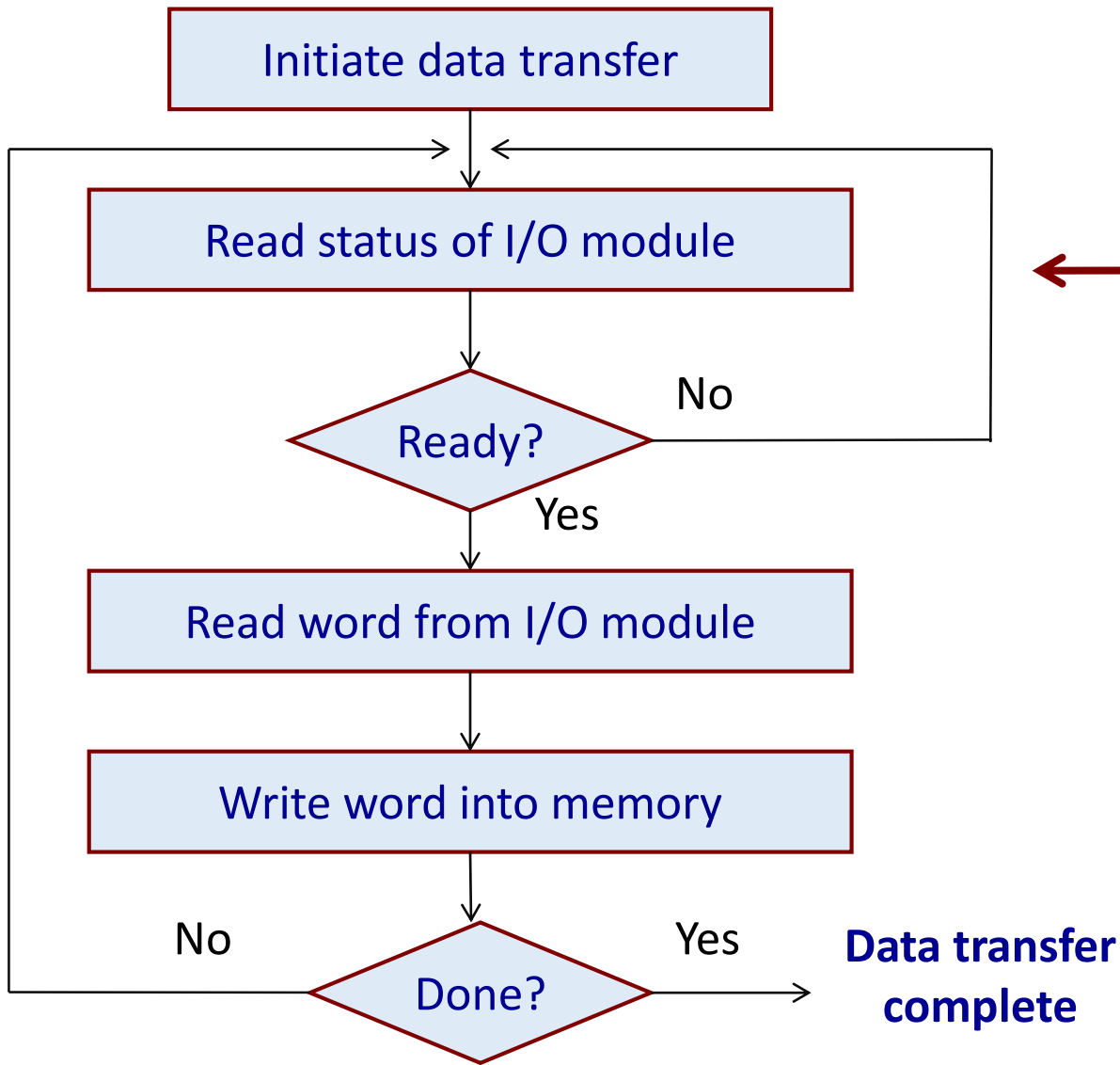
- The I/O device transfers data at a *fixed rate* that is known to the CPU.
- The CPU initiates the I/O operation and transfers successive bytes/words after giving fixed time delays.
- Characteristics:
 - During the time delay, CPU lies idle.
 - Not many I/O devices have strictly synchronous data transfer characteristics.
- A flowchart for synchronous data transfer from an input device is shown on the next slide.



- Error may occur if the input device and the processor get out of synchronization.
- Large number of words cannot be transferred in one go.
- Speed of data transfer depends not only on the speed of I/O device and memory, but also on the execution time of the code.

(b) Asynchronous Data Transfer

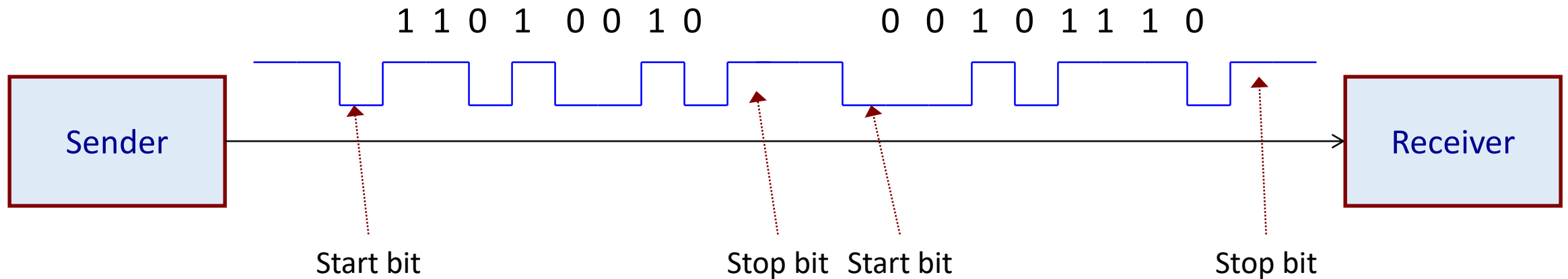
- The CPU does not know when the I/O module will be ready to transfer the next word.
- CPU has to check the status of the I/O module to know when the device is ready to transfer the next word.
 - Called *handshaking*.
- Characteristics:
 - While the CPU is checking whether the I/O module is ready, it cannot do anything else.
 - Wasteful of CPU time for slow devices like keyboard or mouse.



Lot of CPU time is wasted in this loop

- Cannot be used for high-speed devices.
- Speed of data transfer depends not only on the speed of I/O device and memory, but also on the execution time of the code.

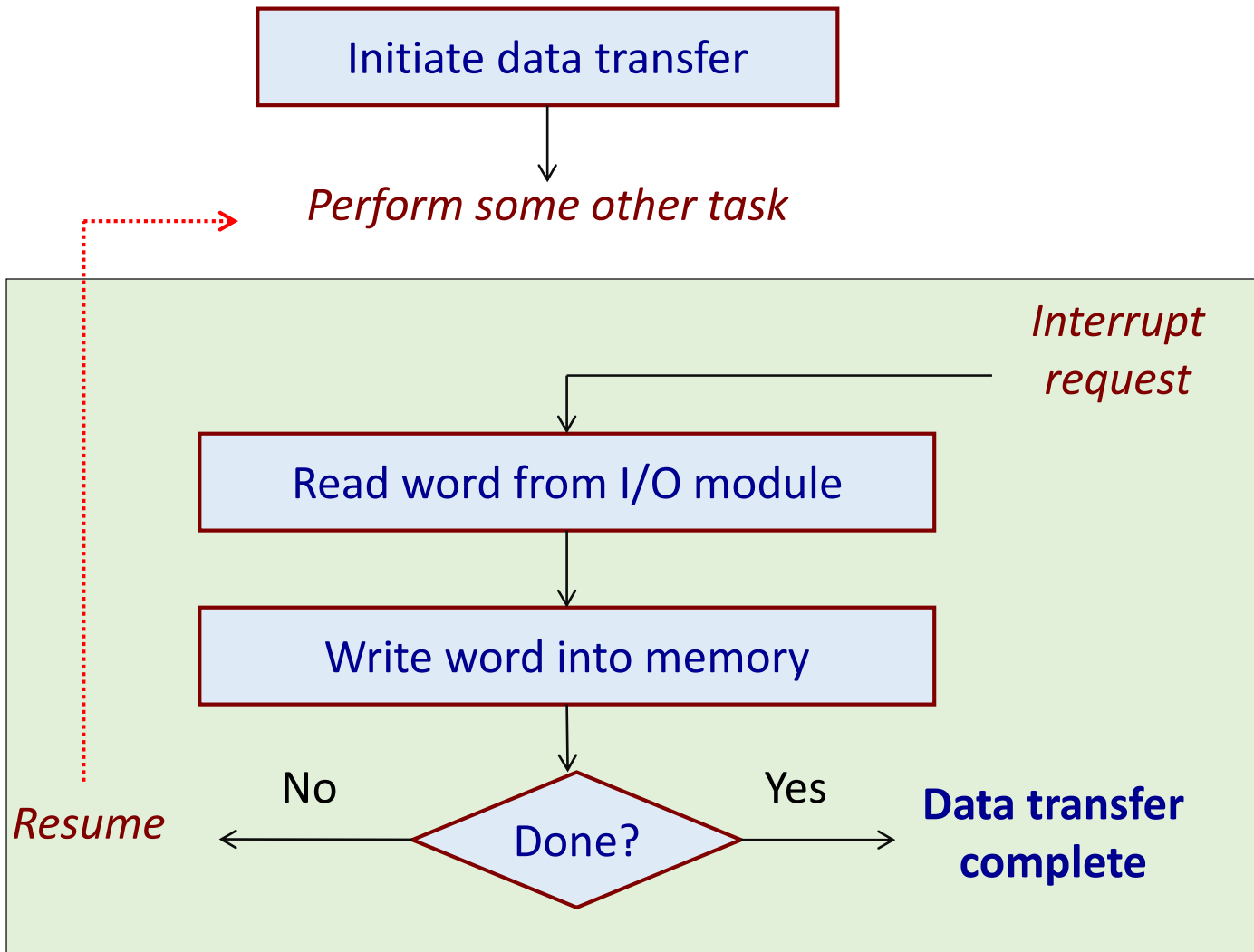
- An example of asynchronous data transfer:
 - Serial data transfer between two devices using start and stop bits.
 - The devices are asynchronous at the level of bytes, but are synchronous at the level of bits within the bytes.



- Just like asynchronous data transfer, receiver waits for the next *START* bit, which indicates the beginning of a new byte transfer.
- After the *START* bit is received, receiver gives (known) bit delays and reads out the 8 bits of the byte.
- The *STOP* bits between the bytes serves to synchronize the data transmission.

(c) Interrupt-Driven Data Transfer

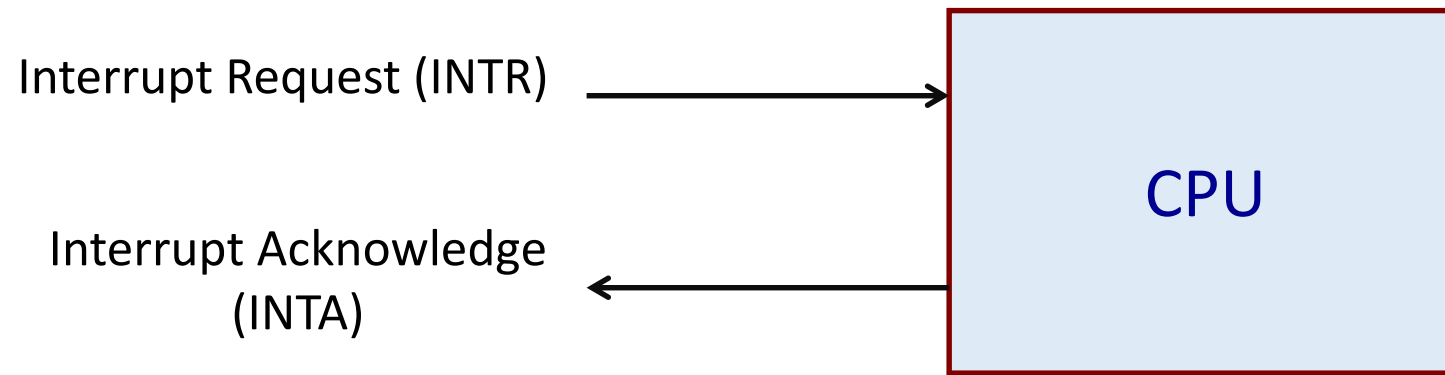
- The CPU initiates the data transfer and proceeds to perform some other task.
- When the I/O module is ready for data transfer, it informs the CPU by activating a signal (called *interrupt request*).
- The CPU suspends the task it was doing, services the request (that is, carries out the data transfer), and returns back to the task it was doing.
- Characteristics:
 - CPU time is not wasted while checking the status of the I/O module.
 - CPU time is required only during data transfer, plus some overheads for transferring and returning control.



*This part of the program that gets activated when interrupt request comes is called **interrupt handler** or **interrupt service subroutine (ISS)**.*

Some Features of Interrupt-Driven Data Transfer

- How is ISS different from a normal subroutine or function?
 - A function is called from well-defined places in the calling program.
 - Only the relevant registers need to be saved on entry to the function, and restored before return.
 - The ISS can get invoked from anywhere in the program that was executing.
 - Depends on when the interrupt request signal arrived.
 - So potentially all the registers that are used in the ISS needs to be saved and restored.



- We shall learn later why the *INTA* signal is required.

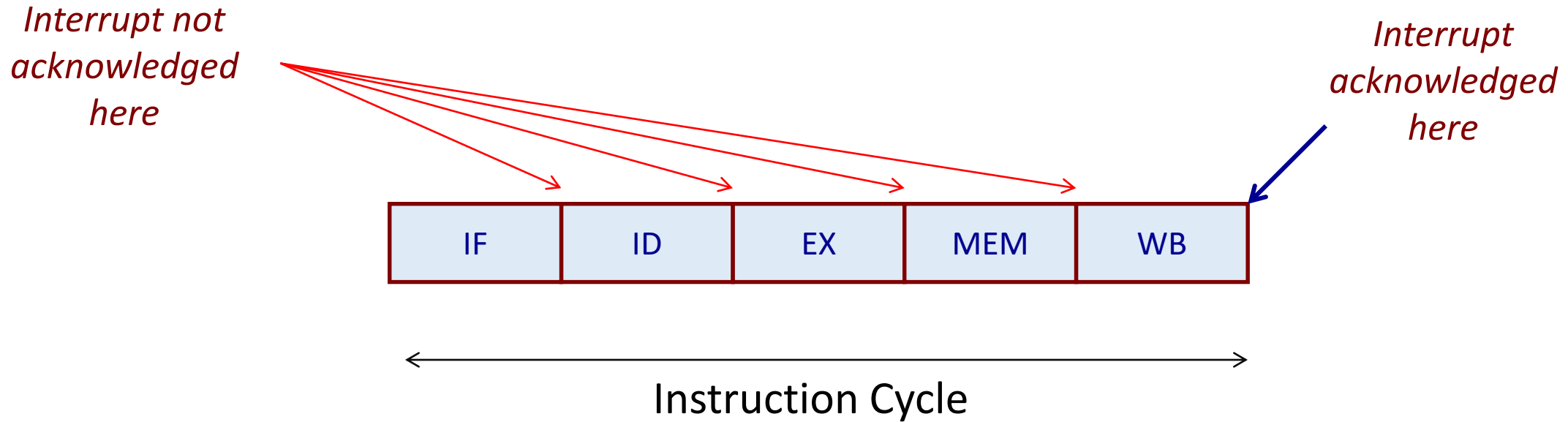
Some Challenges in Interrupts

- For multiple sources of interrupts, how to know the address of the ISR?
- How to handle multiple interrupts?
 - While an interrupt request is being processed, another interrupt request might come.
 - Enabling, disabling and masking of interrupts.
- How to handle simultaneously arriving interrupts?
- Sources of interrupts other than I/O devices.
 - Exceptions, TRAP, etc.

What happens when an interrupt request arrives?

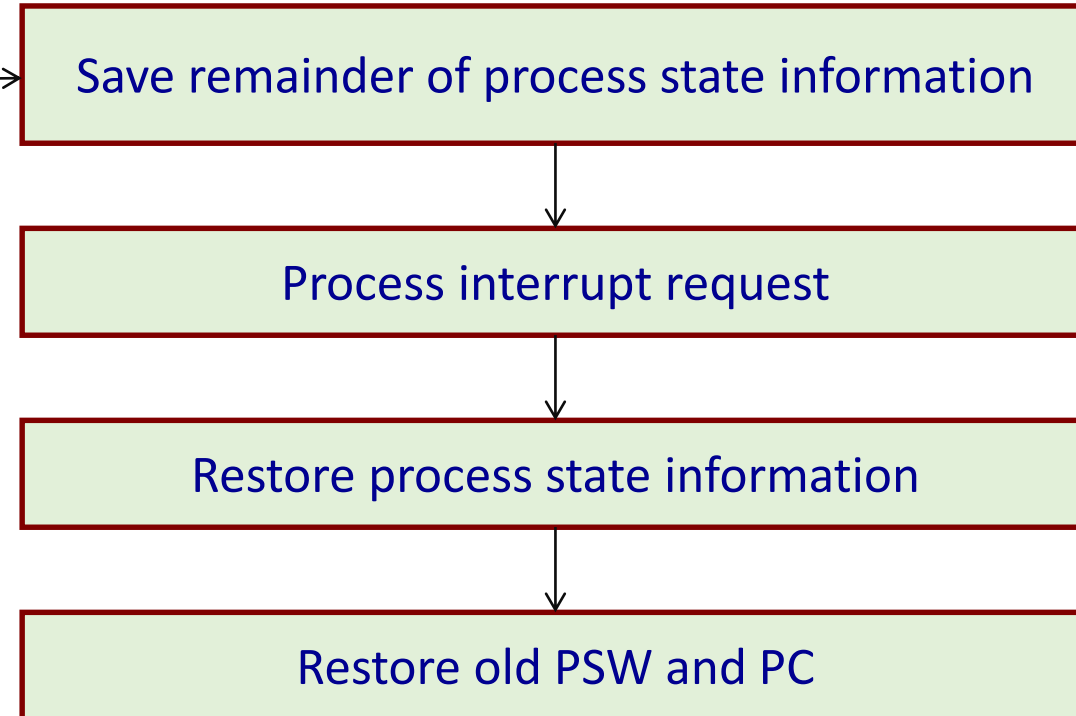
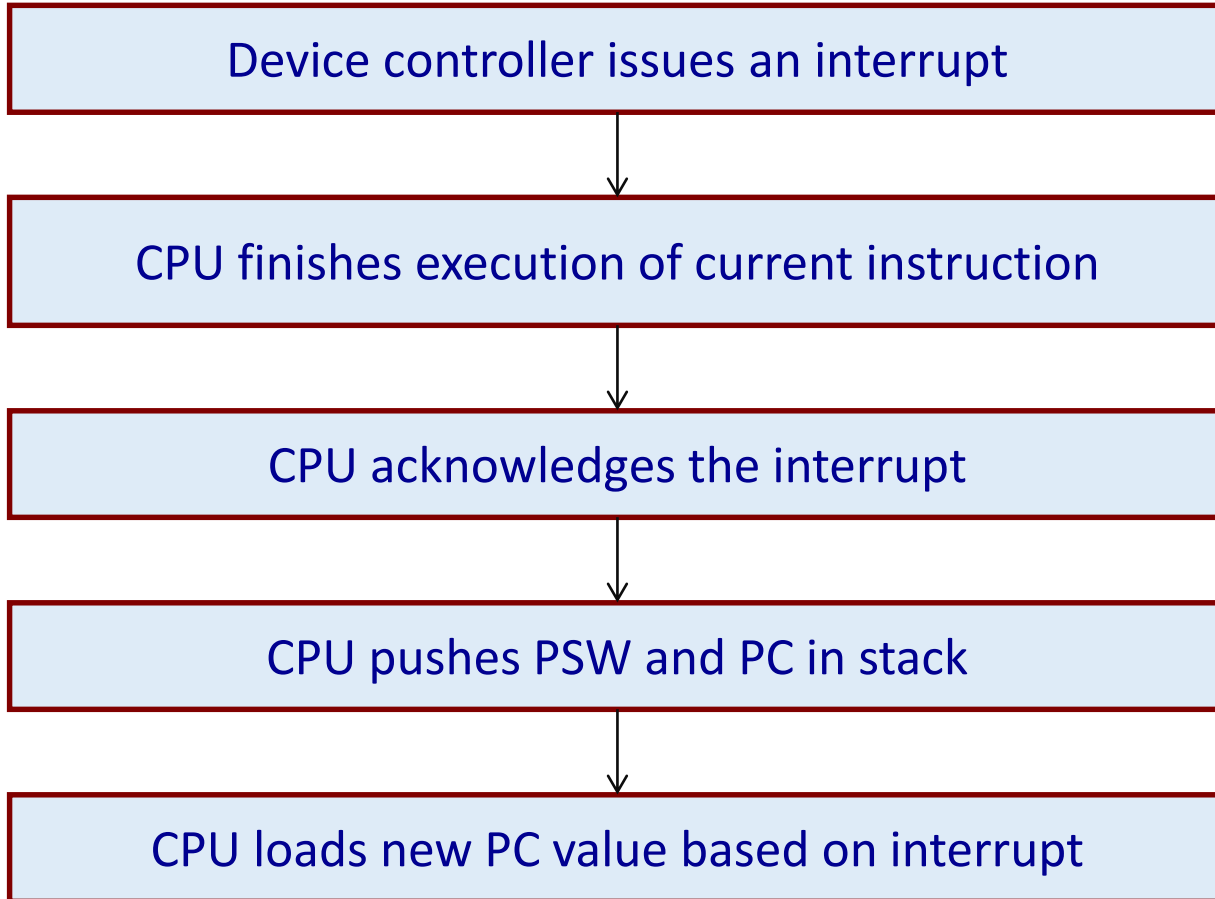
- At the end of the current instruction execution, the PC and *program status word* (PSW) are saved in stack automatically.
 - PSW contains status flags and other processor status information.
- The interrupt is acknowledged, the interrupt vector obtained, based on which control transfers to the appropriate ISS.
 - Different interrupting devices may have different ISS's.
- After handling the interrupt, the ISR executes a special *Return From Interrupt* (RTI) instruction.
 - Restores the PSW and returns control to the saved PC address.
 - Unlike normal RETURN where PSW is not restored.

- An instruction cycle typically consists of several machine cycles.
 - For MIPS32, there are 5 machine cycles.

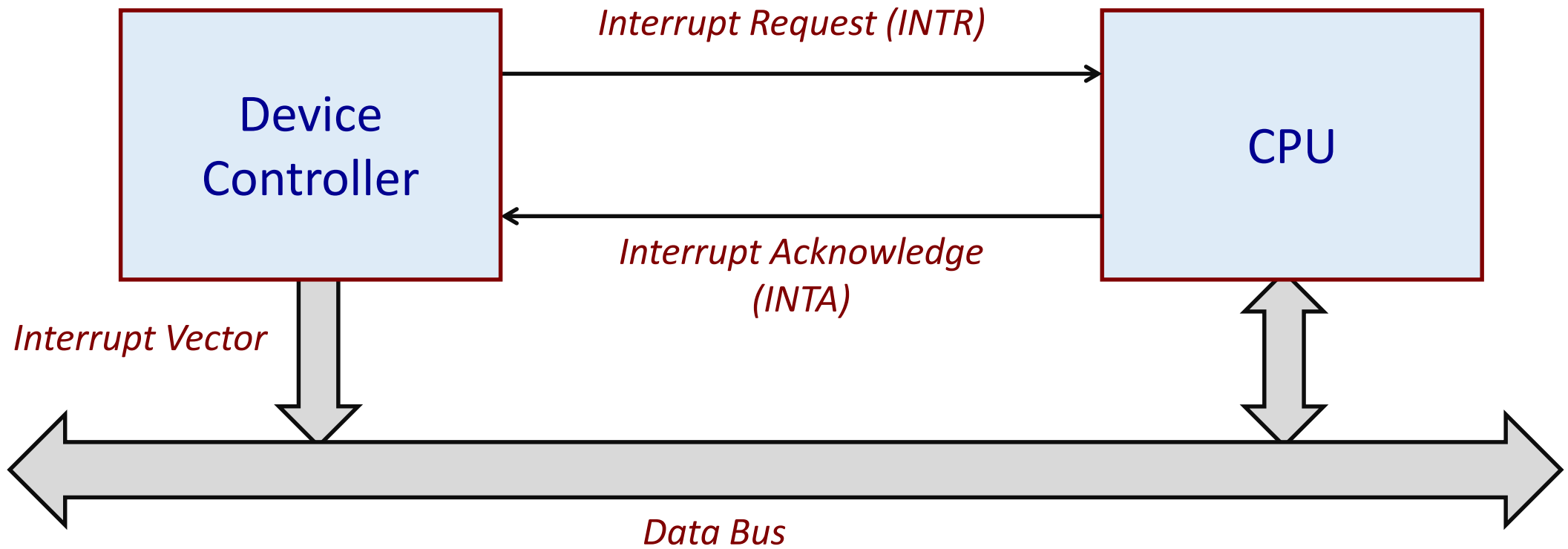


General Interrupt Processing

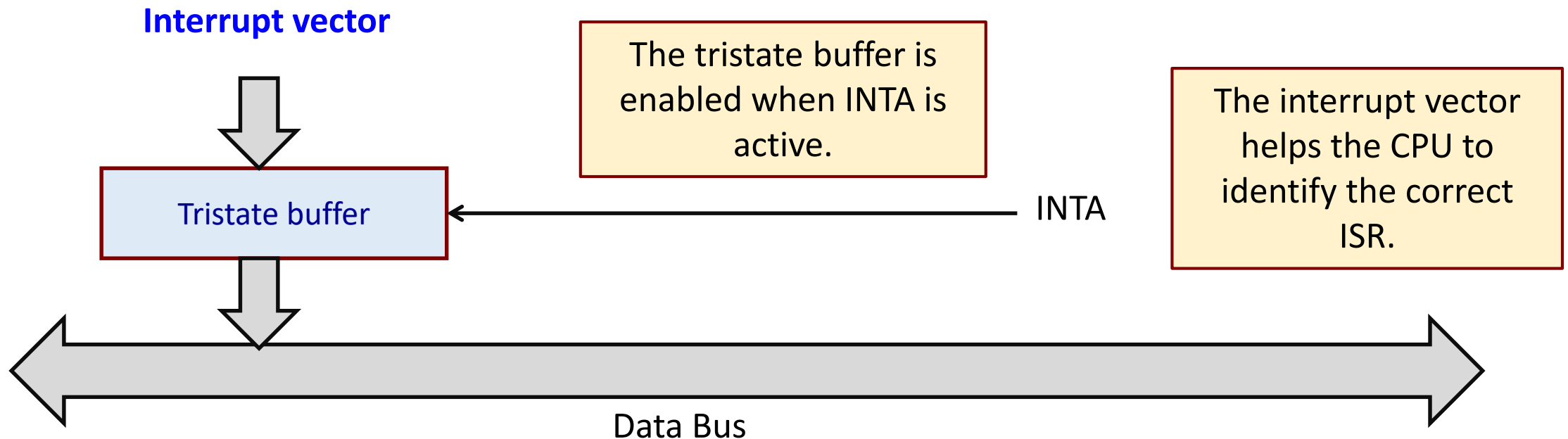
By hardware



By software

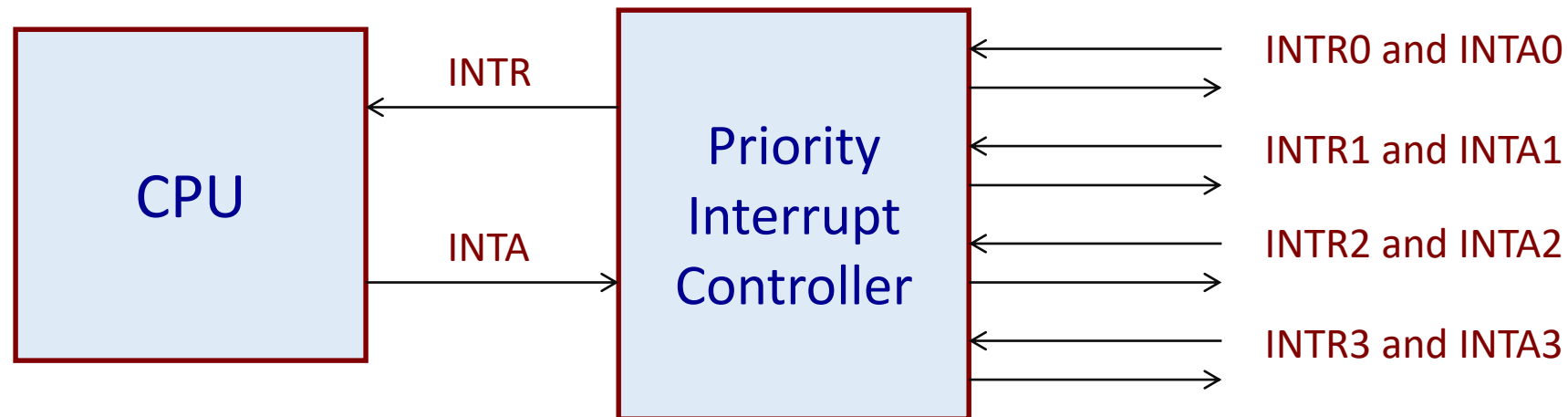


- How is the interrupt vector sent on the data bus in response to INTA?
 - a) Device controller sends *INTR* to the CPU.
 - b) CPU finishes the current instruction and sends back *INTA*.
 - c) Device controller sends *interrupt vector* (or number) over data bus.
 - d) CPU reads the interrupt vector, and identifies the device.



Multiple Devices Interrupting the CPU

- A common solution is to use a *priority interrupt controller*.
 - The interrupt controller interacts with CPU on one side and multiple devices on the other side.
 - For simultaneous interrupt requests, interrupt priority is defined.
 - The interrupt controller is responsible for sending the interrupt vector to CPU.



- How it works?

- The *INTR* line is made active when one or more of the device(s) activate their interrupt request line.

$$INTR = INTR0 + INTR1 + INTR2 + INTR3$$

- When the CPU sends back *INTA*, the interrupt controller sends back the corresponding acknowledge to the interrupting device, and puts the interrupt vector on the data bus.
- The interrupt controller is programmable, where the interrupt vectors for the various interrupts can be programmed (specified).
- For more than one interrupt request simultaneously active, a priority mechanism is used (e.g. *INTR0* is highest priority, followed by *INTR1*, etc.).

How is interrupt nesting handled?

- Consider the scenario:
 - a) A device *D0* has interrupted and the CPU is executing the ISS for *D0*.
 - b) In the mean time, another device *D1* has interrupted.
- Two possible scenarios here:
 - 1) *D1* will interrupt the ISS for *D0*, get processed first, and then the ISS for *D0* will be resumed. → *CREATES PROBLEM FOR MULTI NESTING*
 - 2) Disable the interrupt system automatically whenever an interrupt is acknowledged so that handling of nested interrupts is not required.

- Typical instruction set architectures have the following instructions:
 - EI : Enable interrupt
 - DI : Disable interrupt
- For the second scenario as discussed, the ISS will give an *EI* instruction just before *RTI*.
 - Some ISA combine *EI* and *RTI* in a single instruction.
- The *DI* instruction is sometimes used by the operating system to execute atomic code (e.g. semaphore wait and signal operations).
 - Nobody should interrupt the code while it is being executed.

Cases that make interrupt handling difficult

- For some interrupts, it is not possible to finish the execution of the current instruction.
 - A special *RETURN* instruction is required that would return and restart the interrupted instructions.
- Some examples:
 - a) **Page fault interrupt**: A memory location is being accessed that is not presently available in main memory.
 - b) **Arithmetic exception**: Some error has occurred during some arithmetic operation (e.g. division by zero).

Handling Multiple Devices

- Suppose that a number of devices capable of generating interrupts are connected to the CPU.
- The following questions need to be answered.
 - a) How can the CPU identify the interrupting device?
 - b) How can the CPU obtain the starting address of the appropriate ISS?
 - c) Should interrupt nesting be allowed?
 - d) How should two or more simultaneous interrupt requests be handled?

(a) Device Identification

- Suppose that an external device requests an interrupt by activating an *INTR* line that is common to all the devices. That is,

$$INTR = INTR_1 + INTR_2 + \dots + INTR_n$$

- Each device can have a status bit indicating whether it has interrupted.
 - CPU can *poll* the status bits to find out who has interrupted.
- A better alternative is to use the interrupt vector concept discussed earlier.
 - The interrupting device sends a special identifying code on the data bus upon receiving the interrupt acknowledge.

(b) Find Starting Address of ISS

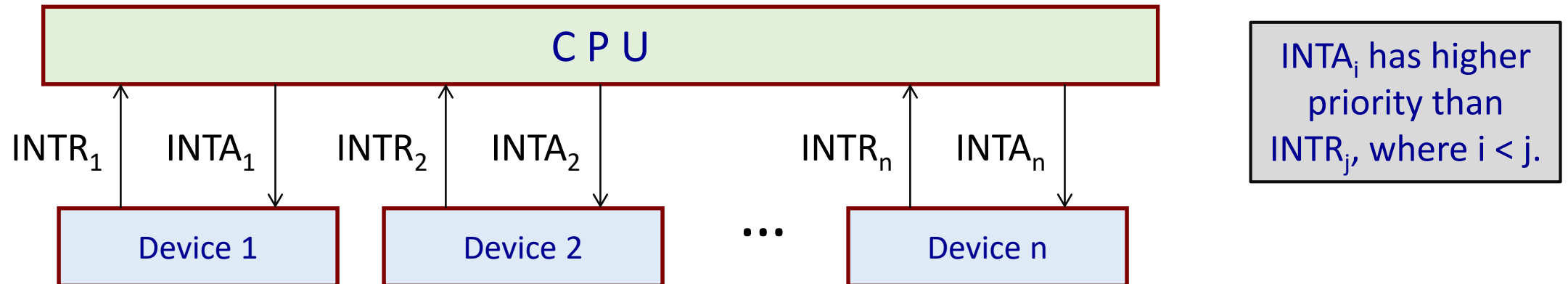
- For a processor with multiple interrupt request inputs, the address of the ISS can be fixed for each individual input.
 - Lacks flexibility.
- If we use the interrupt vector scheme discussed earlier, the device is able to identify itself to the CPU.
 - CPU can then lookup a table where the ISS addresses for all the devices are stored.
 - The interrupt latency is somewhat increased, since we are not immediately jumping to the ISS.

(c) Interrupt Nesting

- A simple approach:
 - Disable all interrupts during the execution of an ISS.
 - This ensures that the interrupt request from one device will not cause more than one interruptions.
 - ISS's are typically short, and the delay they may cause in handling a second interrupt request is often acceptable.
- ***Interrupt priority:***
 - Some interrupting devices may be assigned higher priorities than others.
 - Example: timer interrupt to maintain a real-time clock.
 - Higher priority interrupt may interrupt the ISS of lower priority ones.

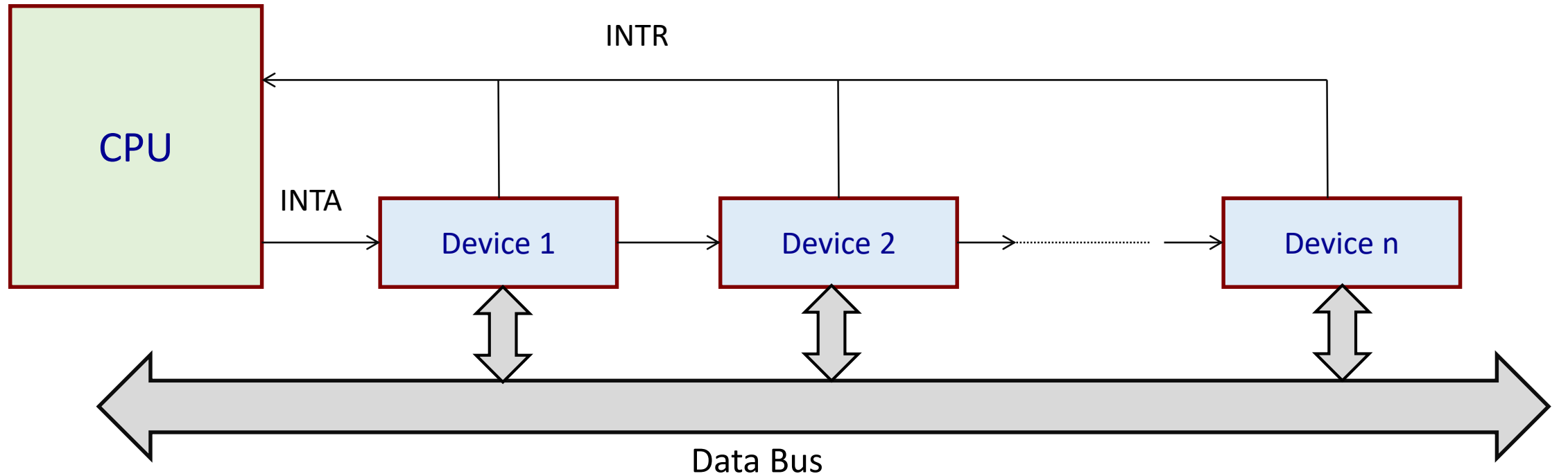
(d) Simultaneous Requests

- Here we consider the problem of simultaneous arrivals of interrupt requests from two or more devices.
 - CPU should have some mechanism by which only one request is serviced while the others are delayed or ignored.
 - If the CPU has multiple interrupt request lines, it can have a *priority scheme* where it accepts the request with the highest priority.

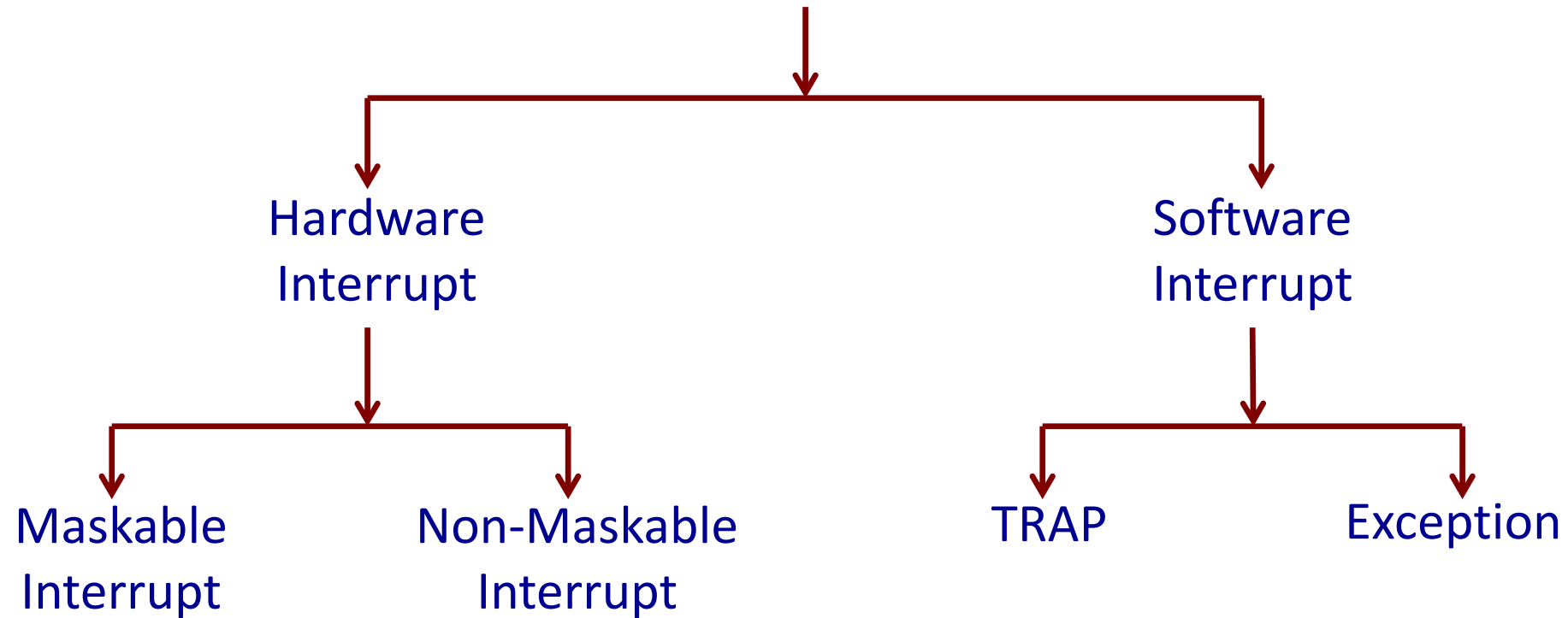


- Another way to assign priority is to use polling using *daisy chaining*.
 - In polling, priority is automatically assigned based on the order in which the devices are polled.
 - In daisy chain connection, the *INTR* line is common to all the devices, but the *INTA* line is connected in a daisy chain fashion allowing it to propagate serially through the devices.
 - A device when it receives *INTA*, passes the signal to the next device only if it had not interrupted. Else, it stops the propagation of *INTA*, and puts the identifying code on the data bus.
 - Thus, the device that is electrically closest to the CPU will have the highest priority.

Daisy Chain Arrangement



Types of Interrupts



- *Hardware Interrupt:*

- The interrupt signal is coming from a device external to the CPU.
- Example: keyboard interrupt, timer interrupt, etc.

- *Maskable Interrupt:*

- Hardware interrupts that can be masked or delayed when a higher priority interrupt request arrives.
- There are processor instructions that can selectively mask and unmask the interrupt request lines of the CPU.

- *Non-Maskable Interrupt:*

- Interrupts that cannot be delayed and should be handled by the CPU immediately.
- Examples: power fail interrupts, real-time system interrupts, etc.

- *Software Interrupt:*

- They are caused due to execution of some instructions.
- Not caused due to external inputs.

- *TRAP:*

- They are special instructions used to request services from the operating system.
- Also called *system calls*.

- *Exception:*

- These are unplanned interrupts generated while executing a program.
- They are generated from within the system.
- Examples: invalid opcode, divide by zero, page fault, invalid memory access, etc.