



# OPERATING SYSTEMS

## Process Management 7

Nitin V Pujari  
Faculty, Computer Science  
Dean - IQAC, PES University

## Course Syllabus - Unit 1

---

### UNIT 1: Introduction and Process Management

Operating-System Structure & Operations, Kernel Data Structures, Computing Environments, Operating-System Services, Operating System Design and Implementation. Process concept: Process in memory, Process State, Process Control Block, Process Creation and Termination, CPU Scheduling and Scheduling Algorithms, IPC - Shared Memory & Message Passing, Pipes - Named and Ordinary. Case Study: Linux/Windows Scheduling Policies.

# OPERATING SYSTEMS

## Course Outline

Class No.	Chapter Title / Reference Literature	Topics to be covered	% of Portions covered	
			Reference chapter	Cumulative
1	1.1-1.2	What Operating Systems Do, Computer-System Organization?	1	21.4
2	1.3,1.4,1.5	Computer-System Architecture, Operating-System Structure & Operations	1	
3	1.10,1.11	Kernel Data Structures, Computing Environments	1	
4	2.1,2.6	Operating-System Services, Operating System Design and Implementation	2	
5	3.1-3.3	Process concept: Process in memory, Process State, Process Control Block, Process Creation and Termination	3	
6	5.1-5.2	CPU Scheduling: Basic Concepts, Scheduling Criteria	5	
7	5.3	Scheduling Algorithms: First-Come, First-Served Scheduling, Shortest-Job-First Scheduling	5	
8	5.3	Scheduling Algorithms: Shortest-Job-First Scheduling (Pre-emptive), Priority Scheduling	5	
9	5.3	Round-Robin Scheduling, Multi-level Queue, Multi-Level Feedback Queue Scheduling	5	
10	5.5,5.6	Multiple-Processor Scheduling, Real-Time CPU Scheduling	5	
11	5.7	Case Study: Linux/Windows Scheduling Policies	5	
12	3.4,3.6.3	IPC - Shared Memory & Message Passing, Pipes – Named and Ordinary	3,6	

## Topics Outline

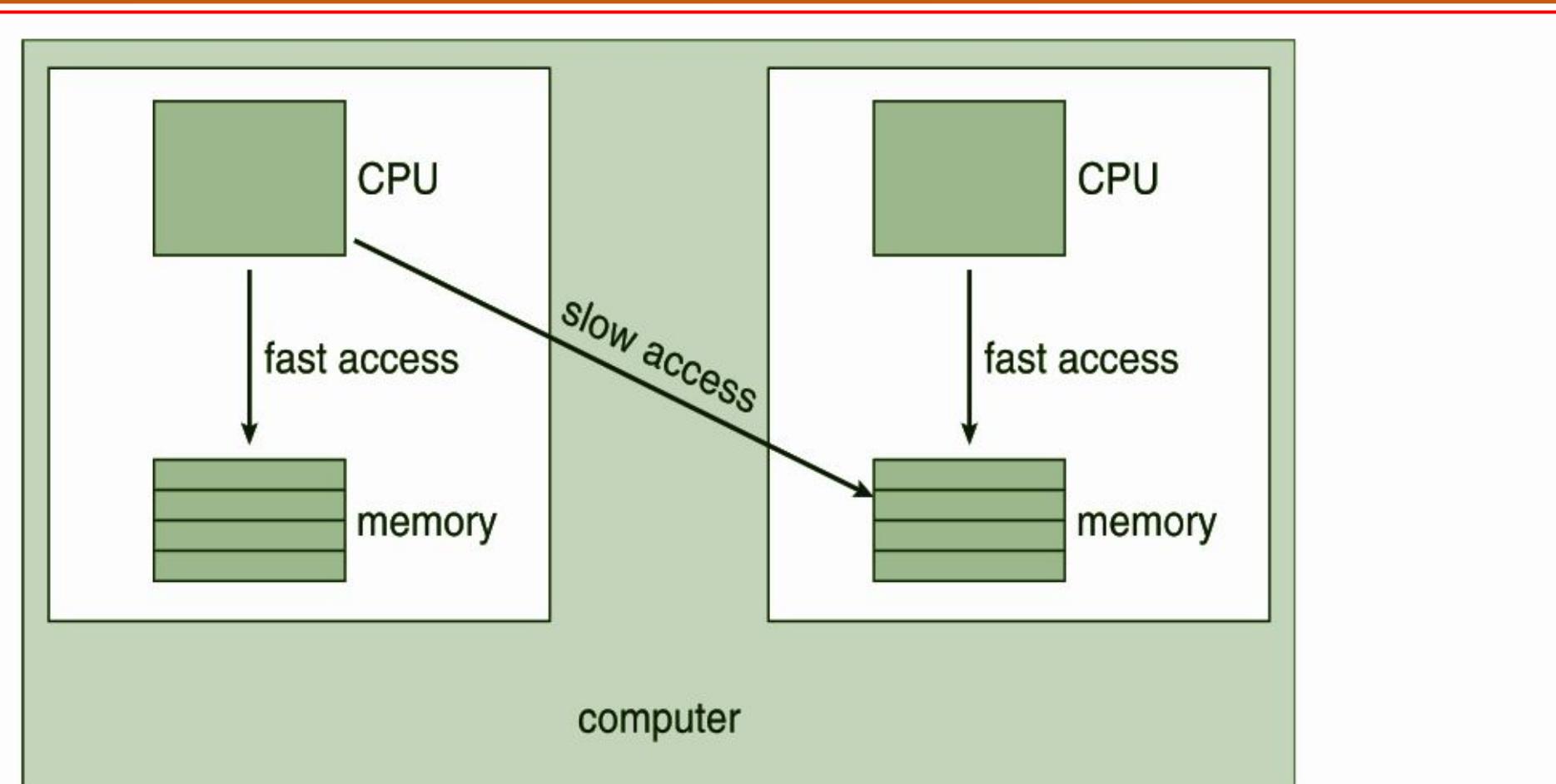
---

- MultiProcessor Scheduling
- NUMA and CPU Scheduling
- MultiProcessor Scheduling - Load Balancing
- Multi Core Processors
- Multithreaded and Multi Core Processors
- Real Time Systems
- Real Time System Requirements
- Real Time System Algorithms
- Real Time CPU Scheduling

# MultiProcessor Scheduling

- CPU scheduling more complex when multiple CPUs are available
- **Homogeneous processors** within a multiprocessor
- **Asymmetric multiprocessing** – only one processor accesses the system data structures, alleviating the need for data sharing
- **Symmetric multiprocessing (SMP)** – each processor is self-scheduling, all processes in common ready queue, or each has its own private queue of ready processes
  - Currently, most common
- **Processor affinity** – process has affinity for processor on which it is currently running
  - **soft affinity, hard affinity, Variations including processor sets**

## Non Uniform Memory access (NUMA and CPU Scheduling)



Note that memory-placement algorithms can also consider affinity

# Multiple-Processor Scheduling – Load Balancing

---

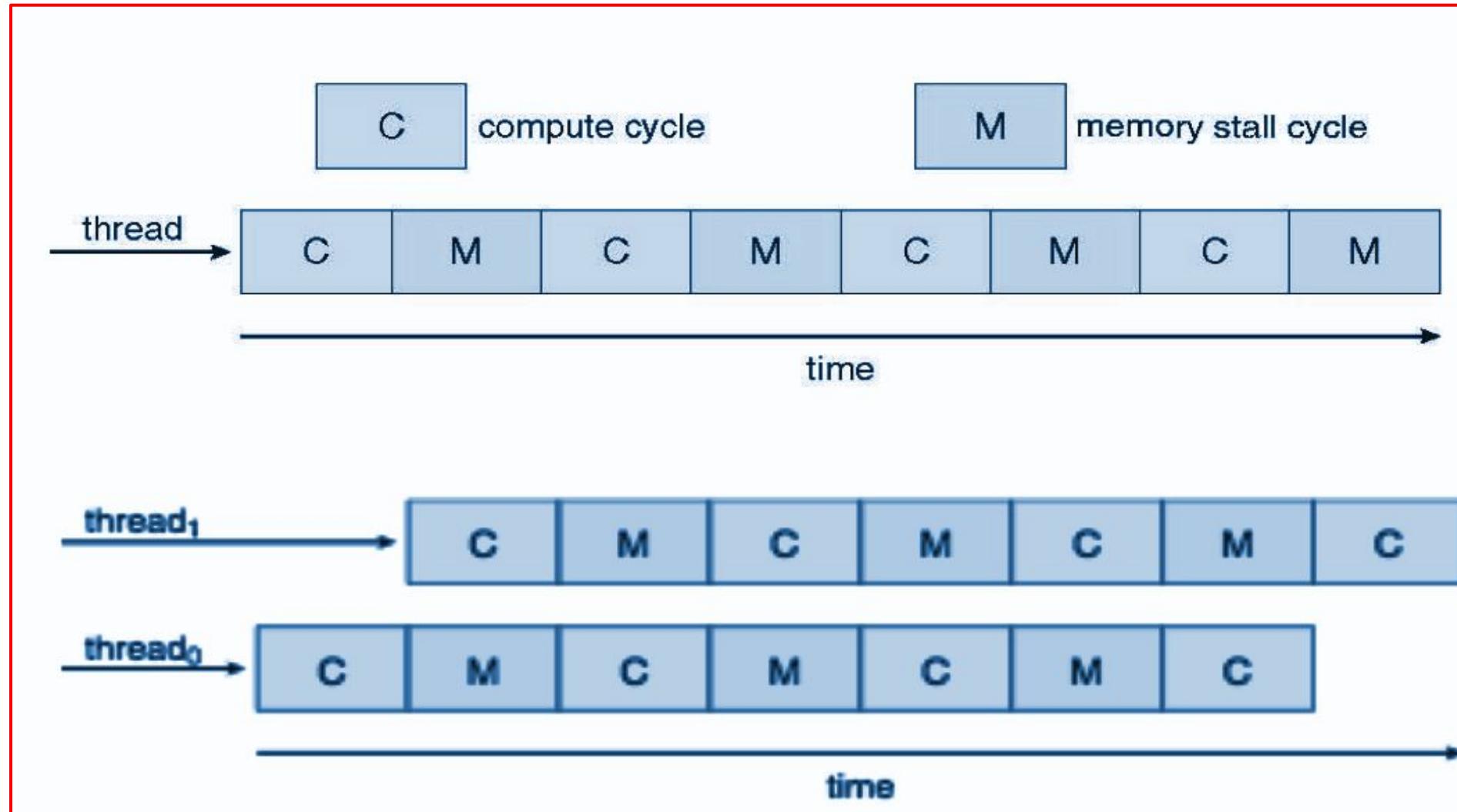
- If SMP, need to keep all CPUs loaded for efficiency
- Load balancing attempts to keep workload evenly distributed
- Push migration – periodic task checks load on each processor, and if found pushes task from overloaded CPU to other CPUs
- Pull migration – idle processors pulls waiting task from busy processor

# Multicore Processors

---

- Recent trend to place multiple processor cores on same physical chip
- Faster and consumes less power
- Multiple threads per core also growing
  - Takes advantage of memory stall to make progress on another thread while memory retrieve happens

## Multithreaded and Multicore System

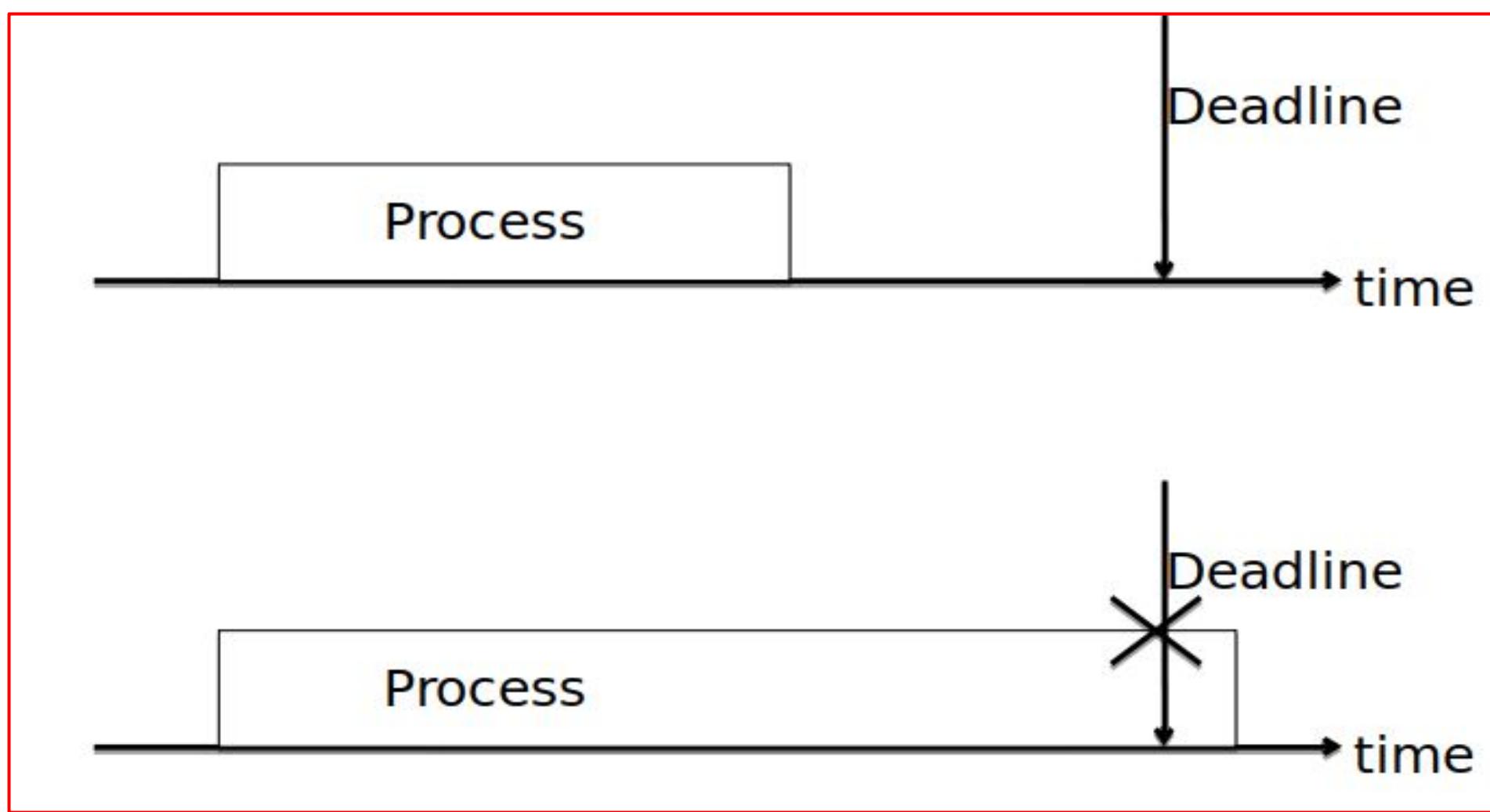


## Real time Systems

---

- A **real-time system** is one whose correctness depends on timing as well as functionality.
- When we discussed more traditional scheduling algorithms, the metrics we looked at were turnaround time (or throughput), fairness, and mean response time. But real-time systems have very different requirements, characterized by different metrics
  - **Timeliness:** How closely does it meet its timing requirements (e.g. ms/day of accumulated tardiness) ?
  - **Predictability:** How much deviation is there in delivered timeliness ?
- A few new concepts
  - **Feasibility:** Whether or not it is possible to meet the requirements of a particular task set ?
  - **Hard real-time :** There are strong requirements that specified tasks be run at specified intervals (or within a specified response time). Failure to meet this requirement (perhaps by as little as microsecond) may result in system failure.
  - **Soft real-time :** We may want to provide very good (e.g. microseconds) response time, the only consequences of missing a deadline are degraded performance or recoverable failures.

## Real Time System Requirements



## Real Time Systems

---

- **Real-time scheduling** is more critical and difficult than traditional time-sharing, and in many ways it is. But real-time systems may have a few characteristics that make scheduling easier:
- We may actually **know** how long each task will take to run. This enables much more **intelligent scheduling**.
- **Starvation** (of low priority tasks) may be acceptable. In aeronautics, a spoiler (sometimes called a lift spoiler or lift dumper) is a device which intentionally reduces the lift component of an airfoil in a controlled way can miss deadline,
- The space shuttle absolutely must sense attitude and acceleration and adjust spoiler positions once per millisecond. But it probably doesn't matter if we update the navigation display once per millisecond or once every ten seconds. Telemetry transmission is probably somewhere in-between. Understanding the relative criticality of each task gives us the freedom to intelligently shed less critical work in times of high demand.
- The workload may be relatively fixed. Normally high utilization implies long queuing delays, as bursty traffic creates long lines. But if the incoming traffic rate is relatively constant, it is possible to simultaneously achieve high utilization and good response time.

## Real Time Systems Algorithms

---

- In the **simplest real-time systems**, where the tasks and their **execution times** are all **known**, there might not even be a scheduler. One task might simply call (or yield to) the next. This model makes a great deal of sense in a system where the tasks form a producer/consumer pipeline (e.g. MPEG frame receipt, protocol decoding, image decompression, display).
- In **more complex real-time system**, with a larger (but still fixed) number of tasks that do not function in a strictly pipeline fashion, it may be possible to do static scheduling. Based on the **list of tasks to be run**, and the expected completion time for each, we can define (at design or build time) a fixed schedule that will ensure timely execution of all tasks.
- But for many real-time systems, the workload changes from moment to moment, based on external events. These require **dynamic scheduling**. For dynamic scheduling algorithms, there are **two** key questions:
  1. how they choose the next (ready) task to run
    - shortest job first
    - static priority ... highest priority ready task
    - soonest start-time deadline first (ASAP)
    - soonest completion-time deadline first (slack time)
  2. how they handle overload (infeasible requirements)
    - best effort
    - periodicity adjustments ... run lower priority tasks less often.
    - work shedding ... stop running lower priority tasks entirely.

## Real Time Systems Algorithms

---

- Preemption may also be a different issue in real-time systems
  - In ordinary time-sharing, preemption is a means of improving mean response time by breaking up the execution of long-running, compute-intensive tasks.
  - A second advantage of preemptive scheduling, particularly important in a general purpose timesharing system, is that it prevents a buggy (infinite loop) program from taking over the CPU. The trade-off, between improved response time and increased overhead (for the added context switches), almost always favors preemptive scheduling. This may not be true for real-time systems

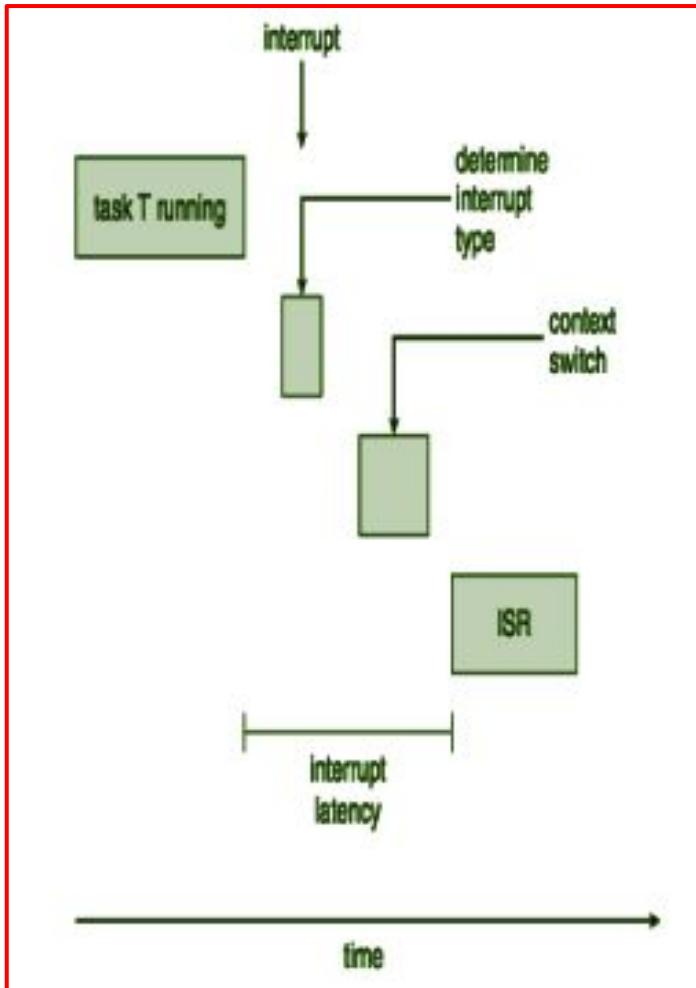
## Real Time Systems Algorithms

---

- For the least demanding real time tasks, a sufficiently **lightly** loaded system might be reasonably successful in meeting its deadlines.
- However, this is achieved simply because the frequency at which the task is run happens to be high enough to meet its real time requirements, not because the scheduler is aware of such requirements.
- A lightly loaded machine running a traditional scheduler can often display a video to a user's satisfaction, not because the scheduler "knows" that a frame must be rendered by a certain deadline, but simply because the machine has enough cycles and a low enough workload to render the frame before the deadline has arrived.

## Real Time CPU Scheduling

- Can present obvious challenges
- **Soft real-time systems** – no guarantee as to when critical real-time process will be scheduled
- **Hard real-time systems** – task must be serviced by its deadline
- Two types of latencies affect performance
  1. Interrupt latency – time from arrival of interrupt to start of routine that services interrupt
  2. Dispatch latency – time for scheduler to take current process off CPU and switch to another

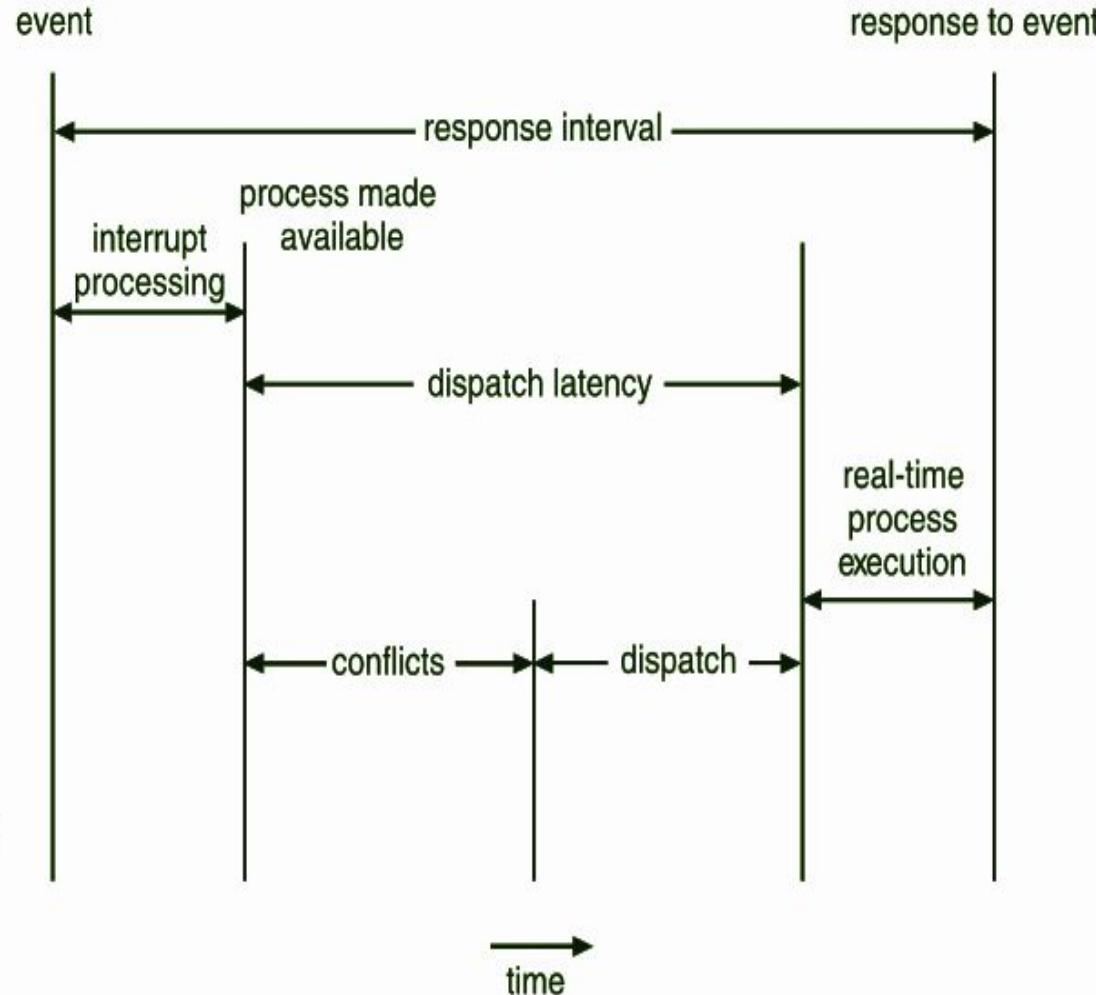


## Real Time CPU Scheduling

- Conflict phase of dispatch latency:

1. Preemption of any process running in kernel mode

2. Release by low-priority process of resources needed by high-priority processes



## Topic Uncovered in this Session

---

- MultiProcessor Scheduling
- NUMA and CPU Scheduling
- MultiProcessor Scheduling - Load Balancing
- Multi Core Processors
- Multithreaded and Multi Core Processors
- Real Time Systems
- Real Time System Requirements
- Real Time System Algorithms
- Real Time CPU Scheduling



**THANK YOU**

**Nitin V Pujari  
Faculty, Computer Science  
Dean - IQAC, PES University**

**nitin.pujari@pes.edu**

**For Course Deliverables by the Anchor Faculty click on [www.pesuacademy.com](http://www.pesuacademy.com)**