CSC 453 Operating Systems

Lecture 5 : Process Scheduling

Concept of Multiprogramming

- Multiprogramming takes advantages of the fact that processes will spend a great deal of their time waiting for I/O operations to finish.
- While process #1 is waiting for I/O, the CPU will execute process #2.

CPU-I/O Burst Cycle

- A program running on a computer has main different "bursts" of activities: bursts of CPU activity and bursts of Input/Output activity.
- Since this is cyclic in nature, it is called the *CPU-I/O Burst Cycle*.











Preemptive vs. Nonpreemptive Scheduling

<u>*Preemptive scheduling*</u> - processes using the CPU can be removed by the system.

<u>Nonpreemptive scheduling</u> - processes using the CPU cannot be removed by the CPU.

<u>Starvation</u> - A situation that arises when a process never gets to the CPU (or to perform an I/O operation, etc.)



- Scheduling Algorithms include:
 - -First-Come-First-Served
 - Shortest Job First
 - Round Robin
 - Priority
 - Guaranteed
 - *Lottery*
 - -Real-Time

First-Come-First-Served

- First-Come-First-Served Algorithm is the simplest CPU scheduling.
- Whichever process requests the CPU first gets it first.
- It is implemented using a standard FIFO single queue.
- Waiting time can be long and it depends heavily on the order in which processes request CPU time:















Predicting The Next CPU Burst Length

• We can substitute for τ_n and expand get the exponential average:

$$\begin{split} \tau_{n+1} &= \alpha \; t_n + (1\text{-}\alpha) \; \alpha \; t_{n-1} + ... \\ &+ (1\text{-}\alpha)^j \; \alpha \; t_{n-j} \text{+}_{...} (1\text{-}\alpha)^{n+1} \; \alpha \; \tau_{0.} \end{split}$$

• More recent terms all have more weight than earlier term in the calculation.



Priority Levels

- There is no general agreement on whether 0 is the highest or lowest priority (priority numbers are assumed to be positive).
 - UNIX uses 0 as the highest priority
 - IBM's MVS uses it as the default (lowest) priority.



- Priorities can be set:
 - <u>internally</u> (by some measurable quantity or quantities such as time limits, memory requirements, number of open files, I/O burstto-CPU burst ratio, etc.)
 - or
 - <u>externally</u> (by system policy, such as process importance, type or availability of funds, sponsoring department, etc.)

Starvation

- *Starvation* is a major problem of priority scheduling algorithms.
- On a busy system, a low-priority process may sit for extremely long periods of time.
- A solution to the problem is *aging*, where we increment the priority (make it a higher priority) for every 1-15 minutes of waiting.

Scenario For Priority Scheduling

Process	Burst Time	<u>Priority</u>
P ₁	10	3
P ₂	1	1
P ₃	2	3
P ₄	1	4
P ₅	5	2





Round-Robin Scheduling and Preemption

- If a process needs less than a time quantum, it releases the CPU voluntarily.
- If a process needs more than a time quantum, it is preempted from the CPU and placed at the back of the ready queue.

Scheduling			
Process	Burst Time		
P ₁	24		
P ₂	3		
P ₃	3		









Lottery Scheduling

- Every process is given, in effect, tickets for a lottery, where the prize is the next time slice (or some other system resource).
- Applied to CPU scheduling, there may be 50 lottery drawings each second, with each winner getting 20 msec of CPU time.
- Important processes can get extra CPU time by being given extra "tickets" for the drawings.
- Cooperating processes can exchange tickets if they wish..







Multiple Processor Scheduling

- Process scheduling on a multiprocessor system is more complex.
- It is easier to schedule *homogeneous* multiprocessor systems than *heterogeneous systems*.
- Identical processors can do *load sharing* with *separate ready queues* or a *common ready queue*.

Symmetric vs. Asymmetric Multiprocessing

- In symmetric multiprocessing, all the processors are considered peers and any one of them can handle any sort of task.
- In asymmetric multiprocessing, there is a hierarchy among the processors and one of them may handle the task of scheduling processes for the others.

What is Real-time Scheduling?

- A real-time system is one in which time plays a crucial role.
- An example is a CD player which must read and then translate the bits into music within a tight time frame.
- Real-time systems can be *hard real time* (where absolute deadlines must always be met) or *soft real time* (where an occasional deadline can be missed).



- Real-time behavior is achieved by by dividing the program into a number of processes, each of which have known behavior.
- Real-time systems react to events which can be *periodic* (happening at regular intervals) or *aperiodic* (not happening at regular intervals).
- If there are m periodic events and event i occurs with a period P_i and requires C_i seconds of CPU time, then the load can only be handled if

 $\Sigma C_i / P_i \ll 1$

Such a system is *schedulable*.

Algorithm Evaluation

- There are several ways in which we can evaluate the scheduling algorithms:
 - Deterministic Modeling
 - Queueing Models
 - Simulation
- Our goal is to see if they help us meet the performance criteria that were discussed earlier.

Deterministic Modeling

- We will assume a predetermined set of data.
- Given that data, we will determine how the scheduling algorithm will perform.
- Deterministic Modeling is easy to understand and implement but it only tells us about the data sets that we use.

Simulations

- We use random numbers to give us a large set of data that should be representative of real-life processing scenarios.
- The distributions can be defined either empirically or mathematically (e.g., a Poisson distribution).
- Such simulations can be expensive and more informative.