



CPU SCHEDULING

Operating Systems
(CS-220)
Fall 2020
FAST NUCES

BCS 5C, BCS 5E

anaum.hamid@nu.edu.pk

CPU SCHEDULING



Basic Concepts



Scheduling Criteria



Scheduling Algorithms



Thread Scheduling



Multiple-Processor Scheduling



Real-Time CPU Scheduling



Operating Systems Examples



Algorithm Evaluation

OBJECTIVES

To introduce CPU scheduling, which is the basis for multi programmed operating systems

To describe various CPU-scheduling algorithms

To discuss evaluation criteria for selecting a CPU-scheduling algorithm for a particular system

To examine the scheduling algorithms of several operating systems

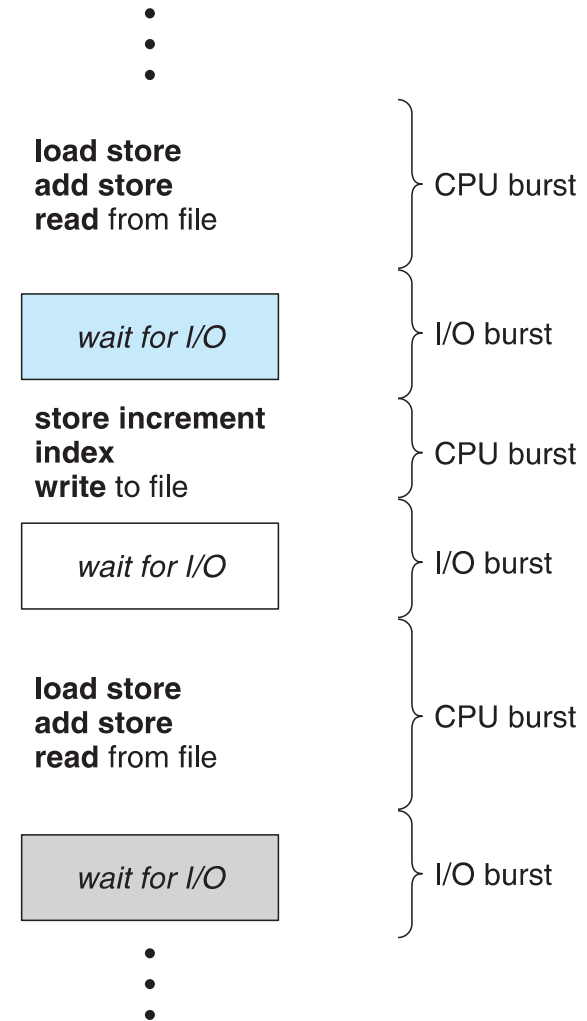
BASIC CONCEPTS

Maximum CPU utilization
obtained with
multiprogramming.

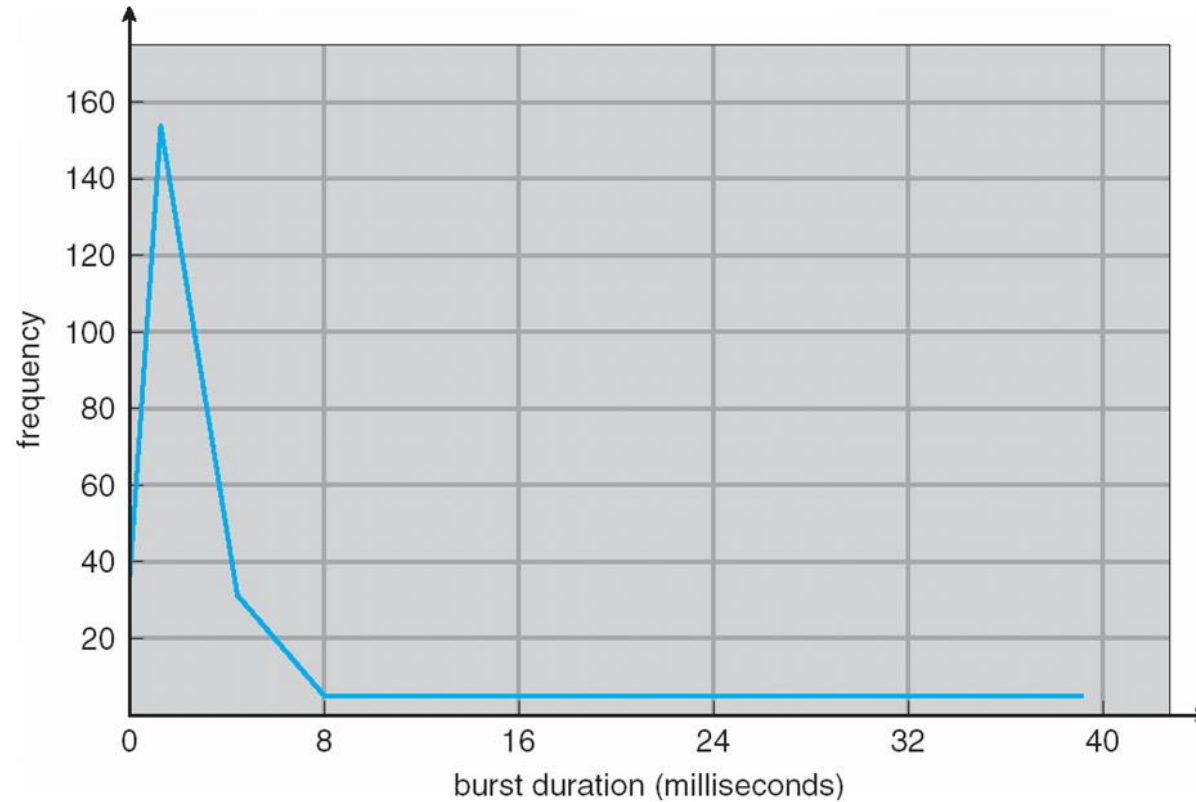
A scheduling system
allows one process to use
the CPU while another is
waiting for I/O, thereby
making full use of
otherwise lost CPU cycles.

CPU—I/O BURST CYCLE

- All processes alternate between two states in a continuing **cycle**; A CPU burst of performing calculations, and An I/O burst, waiting for data transfer in or out of the system.
- Process execution consists of a **cycle** of CPU execution and I/O wait.
- **CPU burst** followed by **I/O burst**.
- CPU burst distribution is of main concern



HISTOGRAM OF CPU-BURST TIMES



CPU bursts vary from process to process, and from program to program, but an extensive study shows frequency patterns similar to that shown in Figure.

CPU SCHEDULER

- **Short-term scheduler** selects from among the processes in ready queue, and allocates the CPU to one of them
 - Queue may be ordered in various ways
- CPU scheduling decisions may take place when a process:
 1. Switches from running to waiting state.
 2. Switches from running to ready state.
 2. Switches from waiting to ready.
 3. Terminates.
- Scheduling under 1 and 4 is **Non-Preemptive** - A new process must be selected.
- All other scheduling is **Preemptive** - There is a choice - To either continue running the current process or select a different one.
 - Consider access to shared data
 - Consider preemption while in kernel mode
 - Consider interrupts occurring during crucial OS activities

PRE-EMPTIVE CAN CAUSE PROBLEMS

- Pre-emptive scheduling can cause problems when two processes share data, because one process may get interrupted in the middle of updating shared data structures.
- Preemption can also be a problem if the kernel is busy implementing a system call (e.g. updating critical kernel data structures) when the preemption occurs. Most modern UNIXes deal with this problem by making the process wait until the system call has either completed or blocked before allowing the preemption. Unfortunately this solution is problematic for real-time systems, as real-time response can no longer be guaranteed.
- Some critical sections of code protect themselves from concurrency problems by disabling interrupts before entering the critical section and re-enabling interrupts on exiting the section.

DISPATCHER

Dispatcher module gives control of the CPU to the process selected by the short-term scheduler; this involves:

- switching context
- switching to user mode
- jumping to the proper location in the user program to restart that program

Dispatch latency – time it takes for the dispatcher to stop one process and start another running

SCHEDULING CRITERIA

CPU utilization – Ideally the CPU would be busy 100% of the time, to waste 0 CPU cycles. On a real system CPU usage should range from 40% (lightly loaded) to 90% (heavily loaded.)

Throughput –Number of processes completed per unit time. May range from 10 / second to 1 / hour depending on the specific processes.

Turnaround time – Amount of time required for a particular process to complete, from submission time to completion. (Wall clock time.)

Waiting time – Amount of time processes spend in the ready queue waiting their turn to get on the CPU

Response time – amount of time it takes from when a request was submitted until the first response is produced, not output (for time-sharing environment)

SCHEDULING ALGORITHM OPTIMIZATION CRITERIA

1. Max CPU utilization
2. Max throughput
3. Min turnaround time
4. Min waiting time
5. Min response time

SCHEDULING ALGORITHM



FIRST- COME, FIRST-SERVED (FCFS) SCHEDULING

<u>Process</u>	<u>Burst Time</u>
P_1	24
P_2	3
P_3	3

Suppose that the processes arrive in the order: P_1 , P_2 , P_3
The Gantt Chart for the schedule is:



Waiting time for $P_1 = 0$; $P_2 = 24$; $P_3 = 27$

Average waiting time: $(0 + 24 + 27)/3 = 17$

FCFS SCHEDULING (CONT.)

Suppose that the processes arrive in the order:

$$P_2, P_3, P_1$$

The Gantt chart for the schedule is:



Waiting time for $P_1 = 6$; $P_2 = 0$; $P_3 = 3$

Average waiting time: $(6 + 0 + 3)/3 = 3$

Much better than previous case

Convoy effect - short process behind long process. All other short processes are left behind a big long process to be done by CPU.

- Consider one CPU-bound and many I/O-bound processes

SHORTEST-JOB-FIRST (SJF) SCHEDULING

Associate with each process the length of its next CPU burst

- Use these lengths to schedule the process with the shortest time

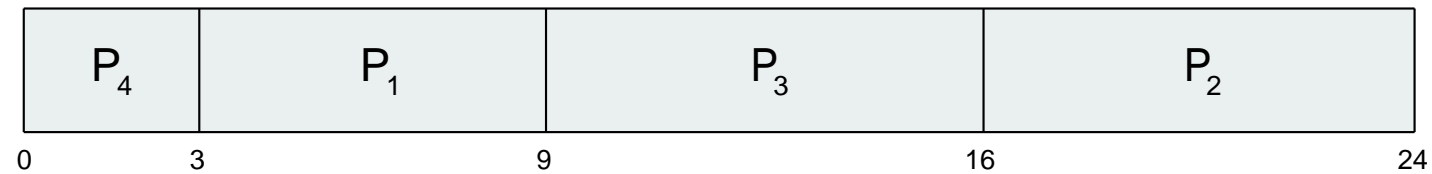
SJF is optimal – gives minimum average waiting time for a given set of processes

- The difficulty is knowing the length of the next CPU request
- Could ask the user

EXAMPLE OF SJF

<u>Process</u>	<u>Burst Time</u>
P_1	6
P_2	8
P_3	7
P_4	3

SJF scheduling chart



$$\text{Average waiting time} = (3 + 16 + 9 + 0) / 4 = 7$$

DETERMINING LENGTH OF NEXT CPU BURST

Can only estimate the length – should be similar to the previous one

- Then pick process with shortest predicted next CPU burst

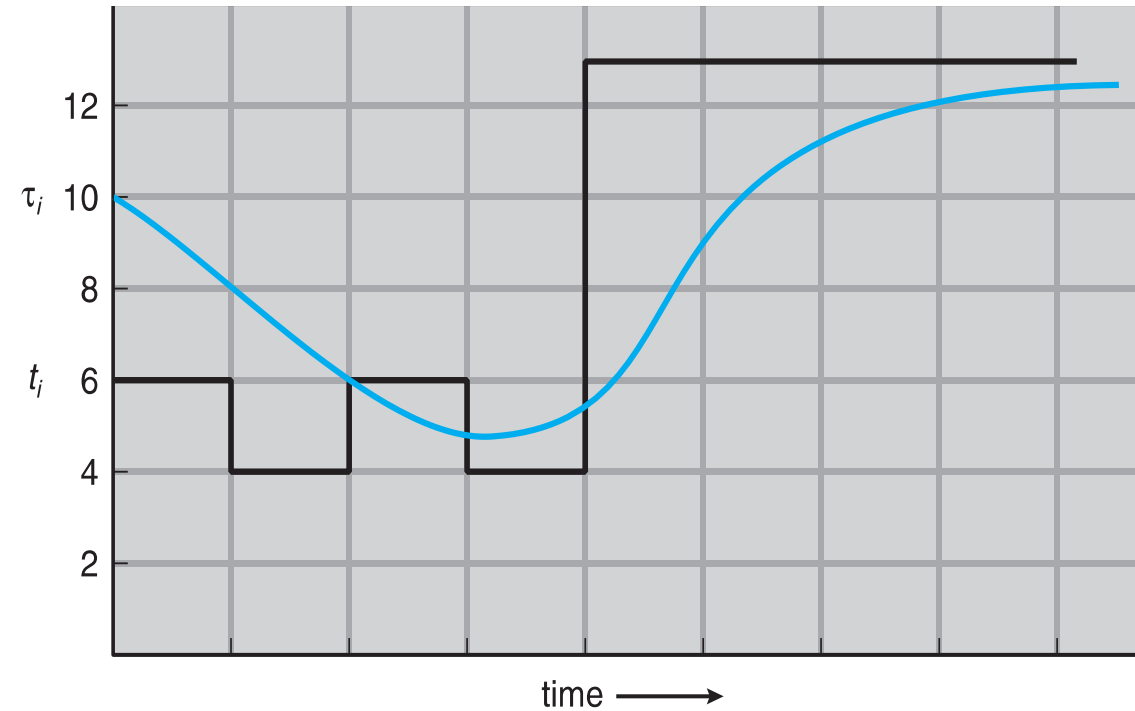
Can be done by using the length of previous CPU bursts, using exponential averaging

1. t_n = actual length of n^{th} CPU burst
2. τ_{n+1} = predicted value for the next CPU burst
3. $\alpha, 0 \leq \alpha \leq 1$
4. Define : $\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n$.

Commonly, α set to $1/2$

Preemptive version called **shortest-remaining-time-first**

PREDICTION OF THE LENGTH OF THE NEXT CPU BURST



CPU burst (t_i)	6	4	6	4	13	13	13	...	
"guess" (τ_i)	10	8	6	6	5	9	11	12	...

EXAMPLES OF EXPONENTIAL AVERAGING

$$\alpha = 0$$

- $\tau_{n+1} = \tau_n$
- Recent history does not count

$$\alpha = 1$$

- $\tau_{n+1} = \alpha t_n$
- Only the actual last CPU burst counts

If we expand the formula, we get:

$$\begin{aligned}\tau_{n+1} = & \alpha t_n + (1 - \alpha)\alpha t_{n-1} + \dots \\ & + (1 - \alpha)^i \alpha t_{n-i} + \dots \\ & + (1 - \alpha)^{n+1} \tau_0\end{aligned}$$

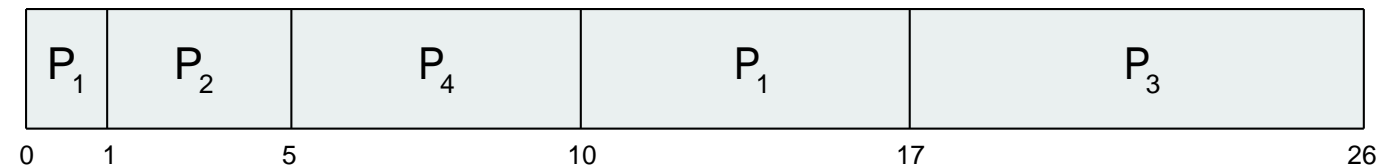
Since both α and $(1 - \alpha)$ are less than or equal to 1, each successive term has less weight than its predecessor

EXAMPLE OF SHORTEST-REMAINING-TIME-FIRST

Now we add the concepts of varying arrival times and preemption to the analysis

<u>Process</u>	<u>Arrival Time</u>	<u>Burst Time</u>
P_1	0	8
P_2	1	4
P_3	2	9
P_4	3	5

Preemptive SJF Gantt Chart



Waiting Time = (Total waiting time - #of MS executing earlier – Arrival Time)

$$WT = (10-1-0)+(1-0-1)+(17-0-2)+(5-0-3) = 26/4 = 6.5$$

Average waiting time = 6.5 msec

PRIORITY SCHEDULING

A priority number (integer) is associated with each process

The CPU is allocated to the process with the highest priority (smallest integer \equiv highest priority)

- Preemptive
- Nonpreemptive

SJF is priority scheduling where priority is the inverse of predicted next CPU burst time

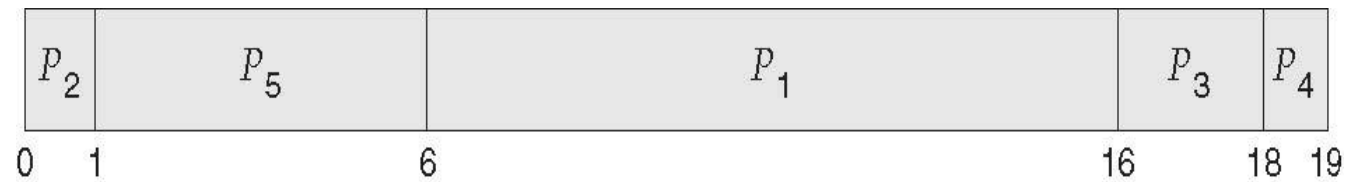
Problem \equiv **Starvation** – low priority processes may never execute

Solution \equiv **Aging** – as time progresses increase the priority of the process

EXAMPLE OF NON-PREEMPTIVE PRIORITY SCHEDULING

<u>Process</u>	<u>Burst Time</u>	<u>Priority</u>
P_1	10	3
P_2	1	1
P_3	2	4
P_4	1	5
P_5	5	2

Priority scheduling Gantt Chart



$$WT = (6+0+16+18+1) = 41/5 = 8.2$$

Average waiting time = 8.2 msec

EXAMPLE OF PREEMPTIVE PRIORITY SCHEDULING (SEE CLASS LECTURE)

ROUND ROBIN (RR)

Specially designed for timesharing/ time slicing systems.

Each process gets a small unit of CPU time (**time quantum** q), usually 10-100 milliseconds. After this time has elapsed, the process is preempted and added to the end of the ready queue.

Timer interrupts every quantum to schedule next process.

Ready queue acts alike circular queue.

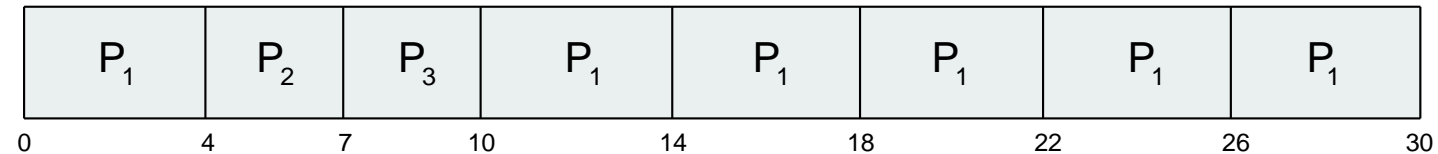
Performance

- q large \Rightarrow FIFO
- q small \Rightarrow q must be large with respect to context switch, otherwise overhead is too high

EXAMPLE OF RR WITH TIME QUANTUM = 4

<u>Process</u>	<u>Burst Time</u>
P_1	24
P_2	3
P_3	3

The Gantt chart is:



Typically, higher average turnaround than SJF, but better **response**

q should be large compared to context switch time

q usually 10ms to 100ms, context switch < 10 usec

Method 1

Turn Around time = Completion time - Arrival time

Waiting time = Turn Around time - Burst time

Process ID	Completion Time	Turnaround Time	Waiting Time
P1	30	$30 - 0 = 30$	$30 - 24 = 6$
P2	7	$7 - 0 = 7$	$7 - 3 = 4$
P3	10	$10 - 0 = 10$	$10 - 3 = 7$

Average Turn Around time

$$\begin{aligned} &= (30 + 7 + 10) / 3 \\ &= 47 / 3 = \underline{15.66 \text{ ms}} \end{aligned}$$

Method 2

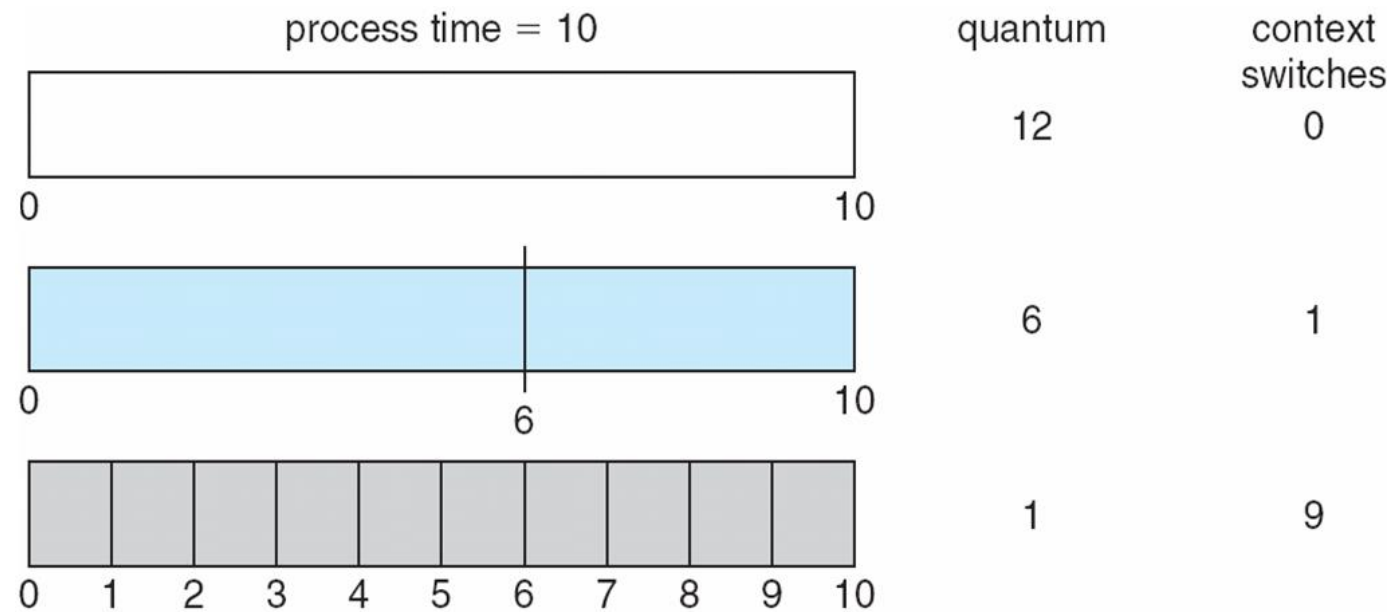
Waiting time = Last Start Time - Arrival Time - (Preemption x Time Quantum)

Process ID	Waiting Time
P1	$26 - 0 - (5 \times 4) = 6$
P2	$4 - 0 - (0 \times 4) = 4$
P3	$7 - 0 - (0 \times 4) = 7$

Average waiting time

$$\begin{aligned} &= (6 + 4 + 7) / 3 \\ &= 17 / 3 = \underline{5.66 \text{ ms}} \end{aligned}$$

TIME QUANTUM AND CONTEXT SWITCH TIME



MULTILEVEL QUEUE

Ready queue is partitioned into separate queues, eg:

- **foreground** (interactive)
- **background** (batch)

Process permanently in each queue

Each queue has its own scheduling algorithm:

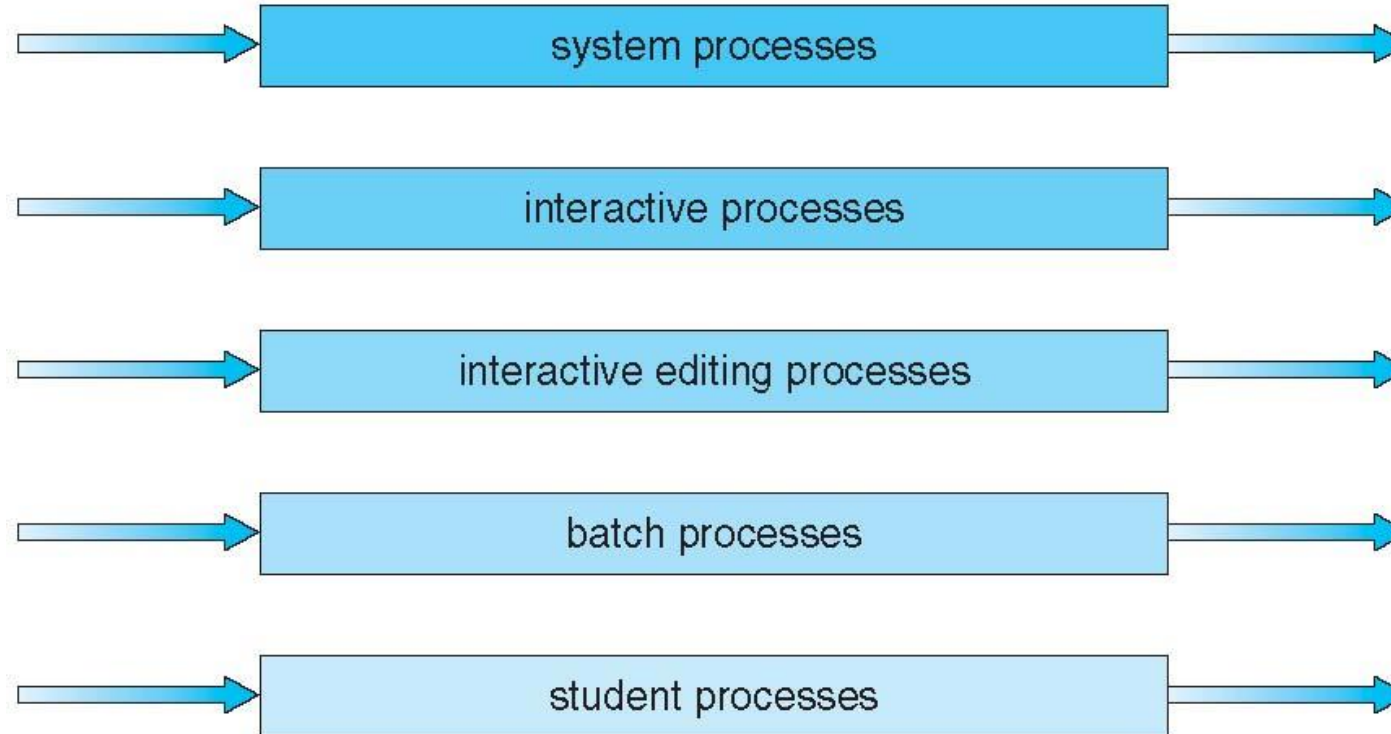
- foreground – RR
- background – FCFS

Scheduling must be done between the queues:

- Fixed priority scheduling; (i.e., serve all from foreground then from background).
Possibility of starvation.
- Time slice – each queue gets a certain amount of CPU time which it can schedule amongst its processes; i.e., 80% to foreground in RR
- 20% to background in FCFS

MULTILEVEL QUEUE SCHEDULING

highest priority



lowest priority

MULTILEVEL FEEDBACK QUEUE

A process can move between the various queues; aging can be implemented this way

Multilevel-feedback-queue scheduler defined by the following parameters:

- ❖ Number of queues
- ❖ Scheduling algorithms for each queue
- ❖ Method used to determine when to upgrade a process
- ❖ Method used to determine when to demote a process.
- ❖ Method used to determine which queue a process will enter when that process needs service

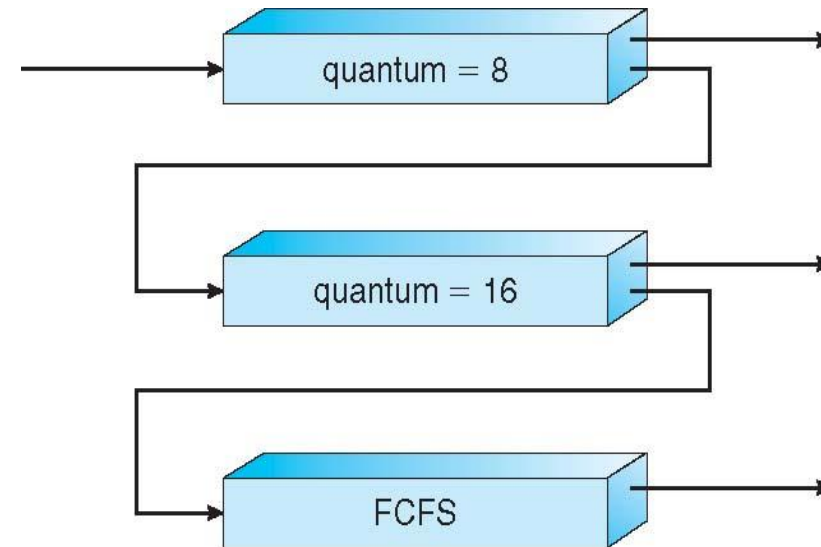
EXAMPLE OF MULTILEVEL FEEDBACK QUEUE

Three queues:

- Q_0 – RR with time quantum 8 milliseconds
- Q_1 – RR time quantum 16 milliseconds
- Q_2 – FCFS

Scheduling

- A new job enters queue Q_0 which is served FCFS
 - When it gains CPU, job receives 8 milliseconds
 - If it does not finish in 8 milliseconds, job is moved to queue Q_1
- At Q_1 job is again served FCFS and receives 16 additional milliseconds
 - If it still does not complete, it is preempted and moved to queue Q_2



MULTIPLE-PROCESSOR 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

Processor affinity – process has affinity for processor on which it is currently running

- **soft affinity**
- **hard affinity**
- Variations including **processor sets**

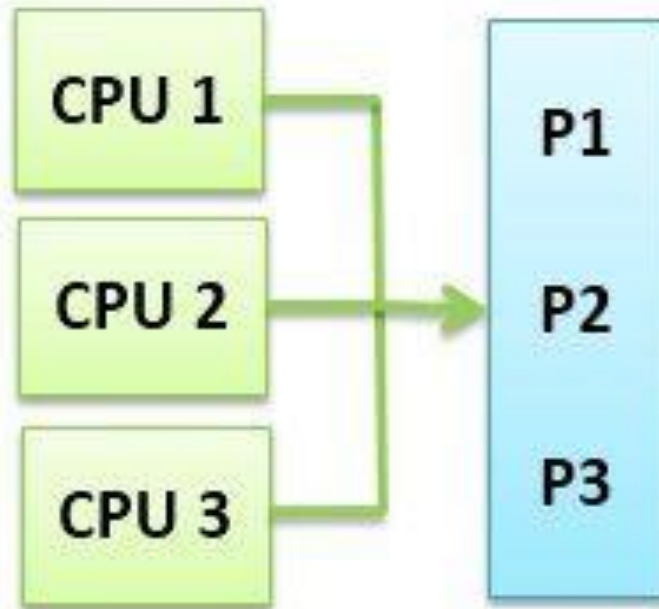
Soft affinity is a mechanism that keeping process on the same CPU as long as possible otherwise, process is migrated to other processor.

Hard affinity is a mechanism that enables processes can run only on a fixed set of one or more processor(s).

SYMMETRIC VS. ASYMMETRIC MULTIPROCESSING ARCHITECTURE

(RECALL FROM CHAPTER 2)

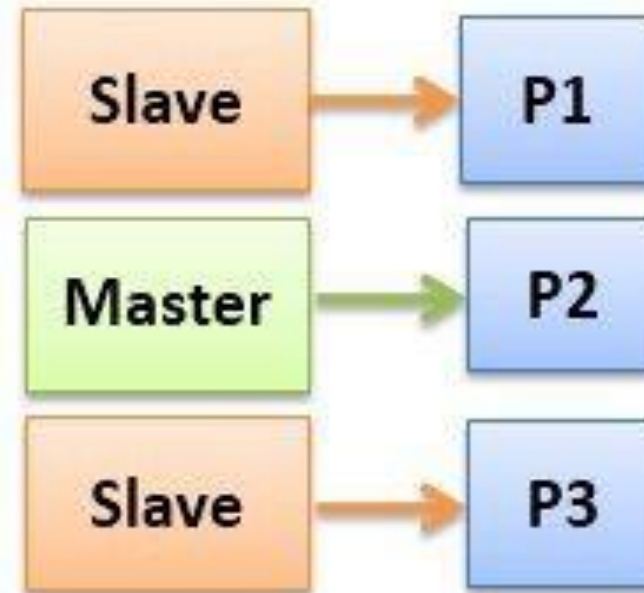
Symmetric Multiprocessing



(Shared Memory)

Vs

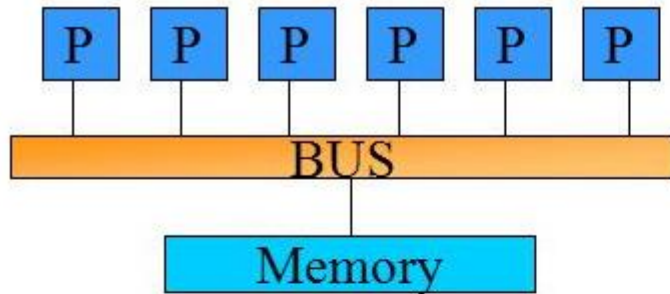
Asymmetric Multiprocessing



(No Shared Memory)

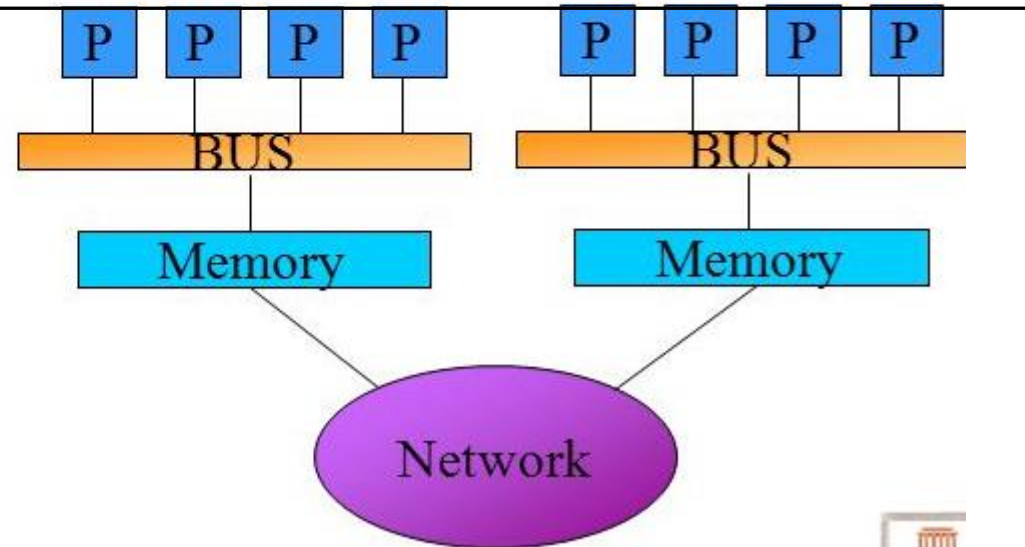
A DUAL-CORE DESIGN (RECALL FROM CHAPTER 2)

UMA and NUMA architecture variations Multi-chip and **multicore**

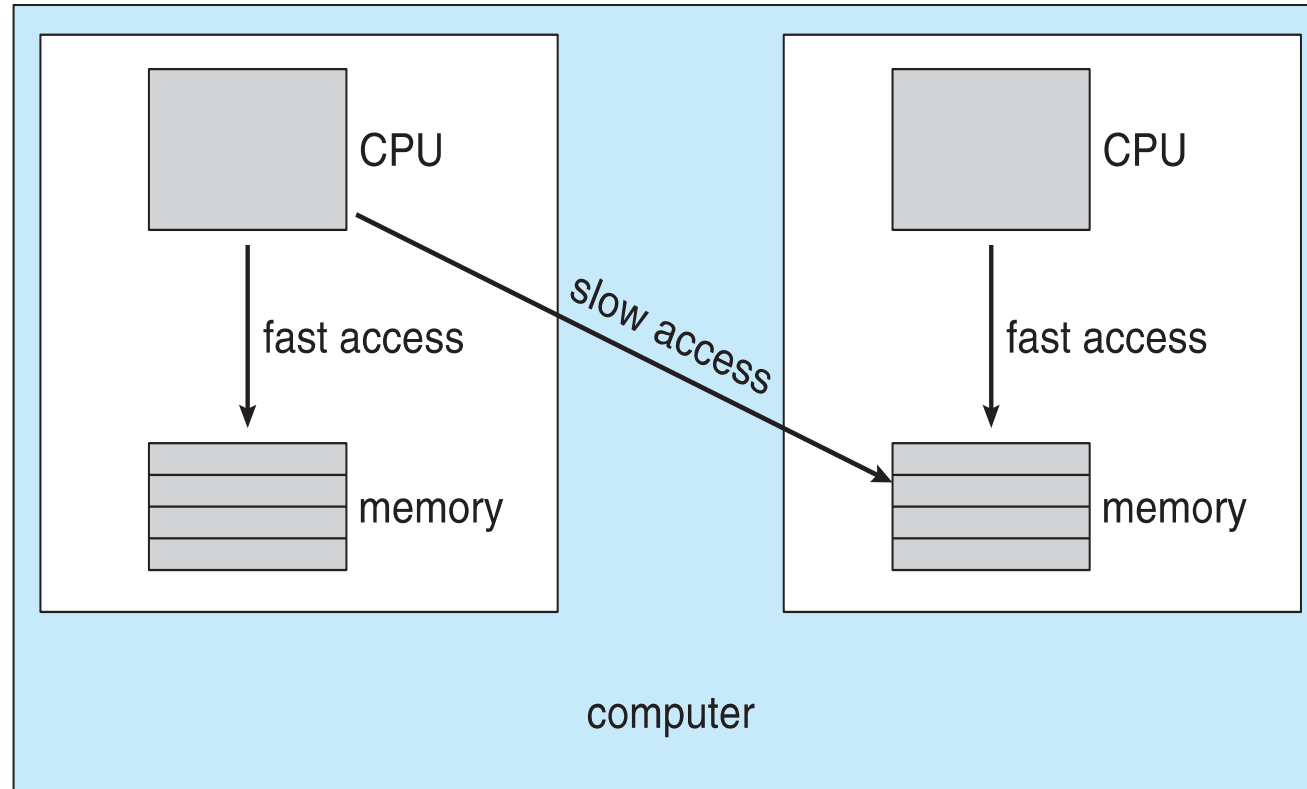


Uniform memory access (UMA):
Each processor has uniform access to memory. Also known as **symmetric multiprocessors**, or SMPs (Sun E10000)

Non-uniform memory access (NUMA): Time for memory access depends on location of data. Local access is faster than non-local access. Easier to scale than SMPs (SGI Origin)



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

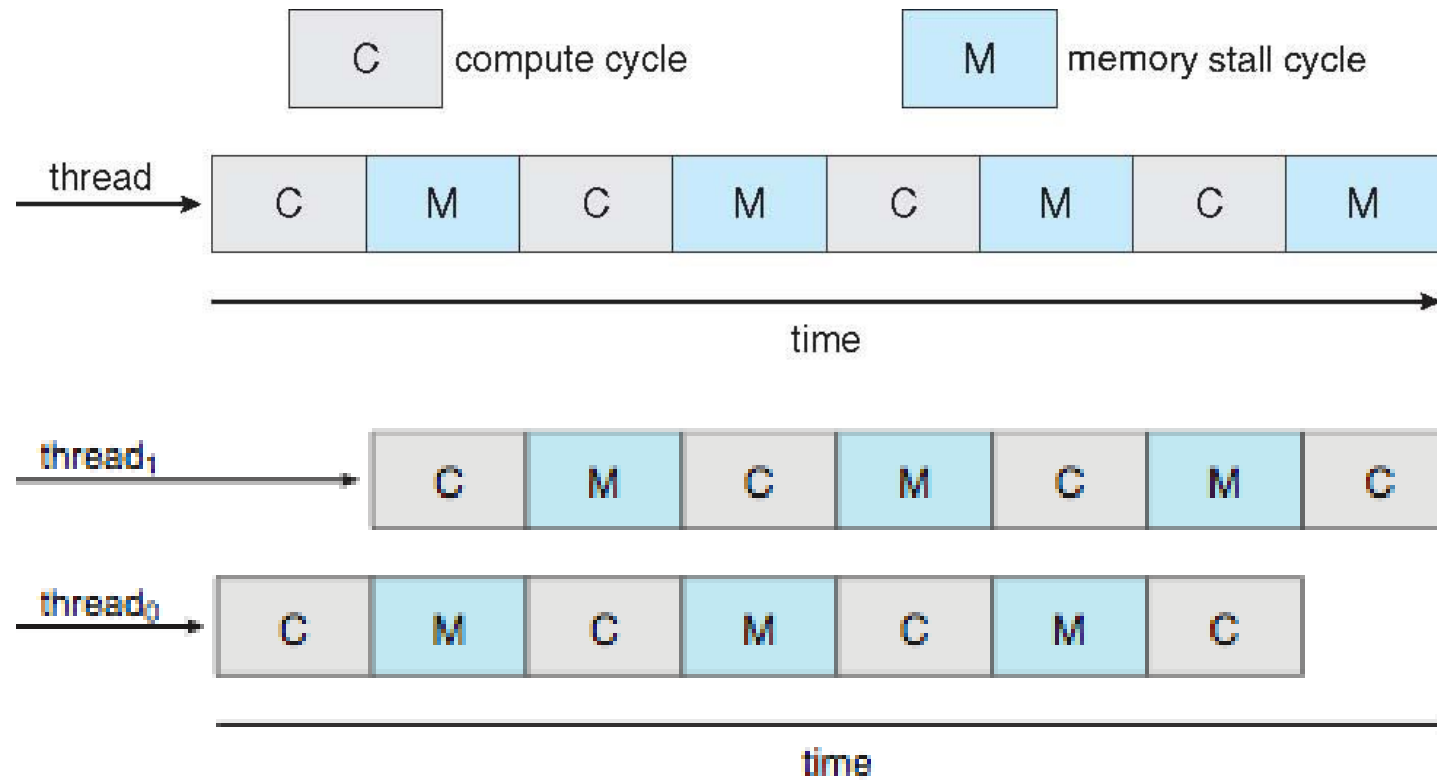
Recent trend to place multiple processor cores on same physical chip

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 MULTICORE SYSTEM



REAL-TIME CPU SCHEDULING

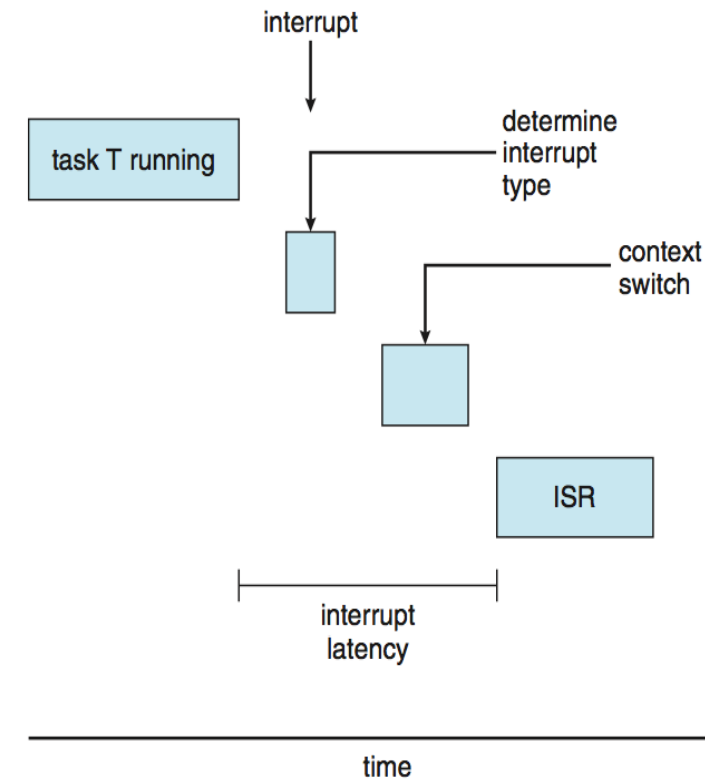
Can present obvious challenges

Soft real-time systems – no guarantee 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.
1. Dispatch latency – time for scheduler to take current process off CPU and switch to another



REAL-TIME CPU SCHEDULING (CONT.)

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

