

# Chapter 3: Processes



# Chapter 3: Processes

---

- Process Concept
- Process Scheduling
- Operations on Processes
- Interprocess Communication
- Communication in Client-Server Systems

# Objectives

---

- To introduce the notion of a process -- a program in execution, which forms the basis of all computation
- To describe the various features of processes, including scheduling, creation and termination, and communication
- To explore interprocess communication using shared memory and message passing
- To describe communication in client-server systems

# Process Concept

---

- An operating system executes a variety of programs:
  - Batch system – **jobs**
  - Time-shared systems – **user programs** or **tasks**
- Textbook uses the terms **job** and **process** almost interchangeably
- **Process** – a program in execution; process execution must progress in sequential fashion
- Multiple parts
  - The program code, also called **text section**
  - Current activity including **program counter**, processor registers
  - **Stack** containing temporary data
    - Function parameters, return addresses, local variables
  - **Data section** containing global variables
  - **Heap** containing memory dynamically allocated during run time

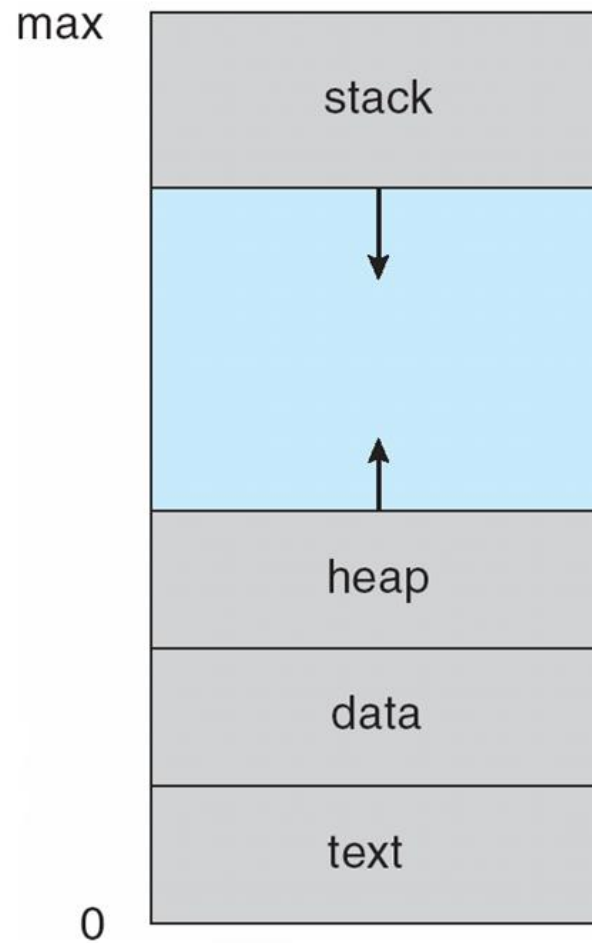
# Process Concept (Cont.)

---

- Program is **passive** entity stored on disk (**executable file**), process is **active**
  - Program becomes process when executable file loaded into memory
- Execution of program started via GUI mouse clicks, command line entry of its name, etc
- One program can be several processes
  - Consider multiple users executing the same program

# Process in Memory

---



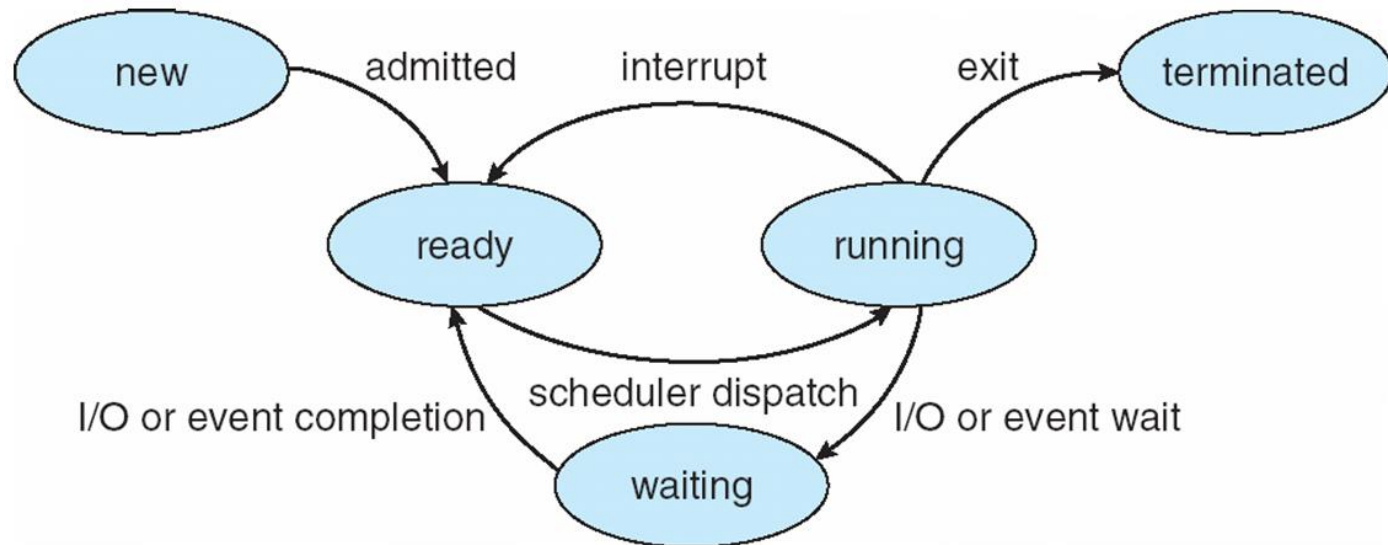
# Process State

---

- As a process executes, it changes **state**
  - **new**: The process is being created
  - **running**: Instructions are being executed
  - **waiting**: The process is waiting for some event to occur
  - **ready**: The process is waiting to be assigned to a processor
  - **terminated**: The process has finished execution

# Diagram of Process State

---

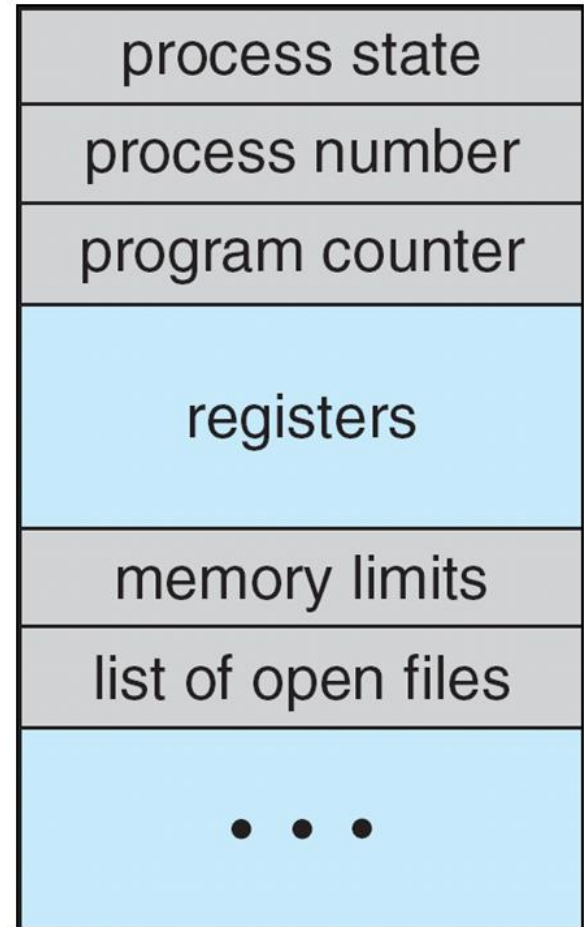




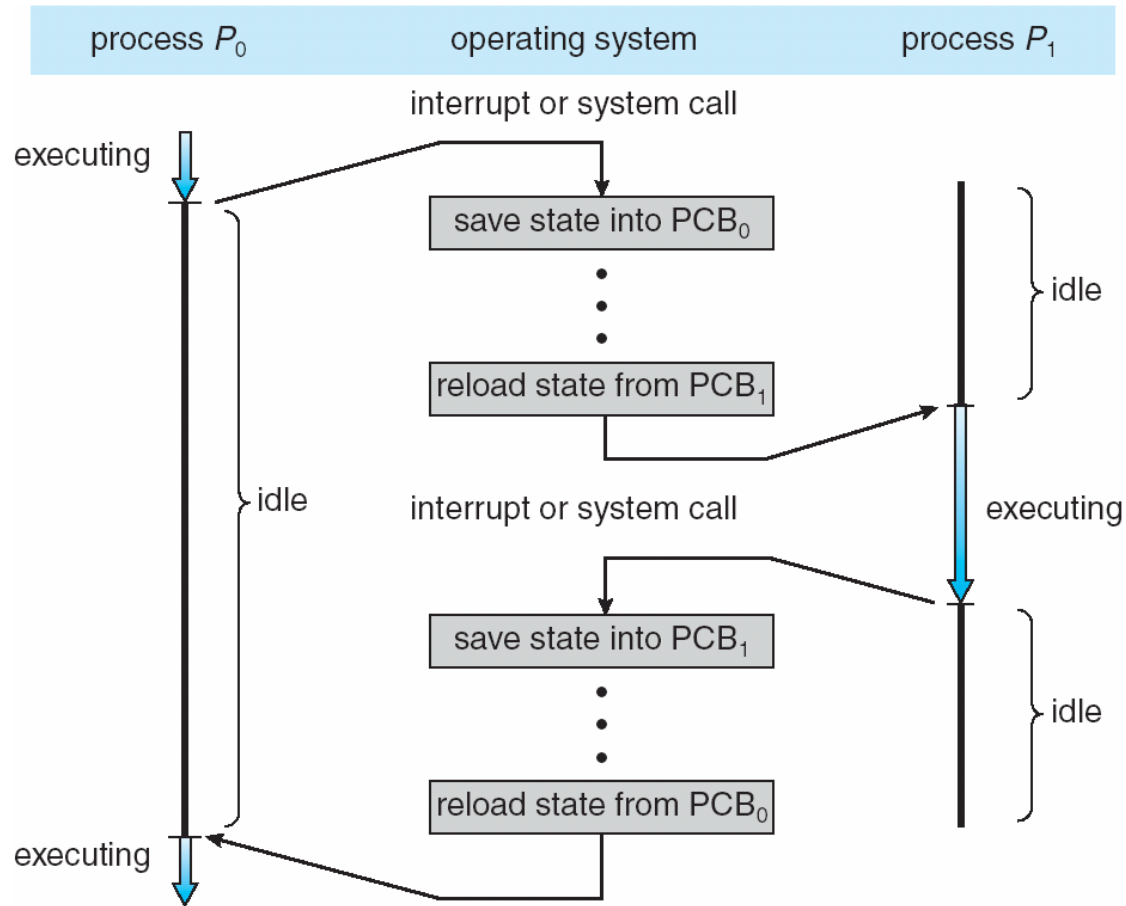
# Process Control Block (PCB)

Information associated with each process  
(also called **task control block**)

- ❑ Process state – running, waiting, etc
- ❑ Program counter – location of instruction to next execute
- ❑ CPU registers – contents of all process-centric registers
- ❑ CPU scheduling information- priorities, scheduling queue pointers
- ❑ Memory-management information – memory allocated to the process
- ❑ Accounting information – CPU used, clock time elapsed since start, time limits
- ❑ I/O status information – I/O devices allocated to process, list of open files



# CPU Switch From Process to Process



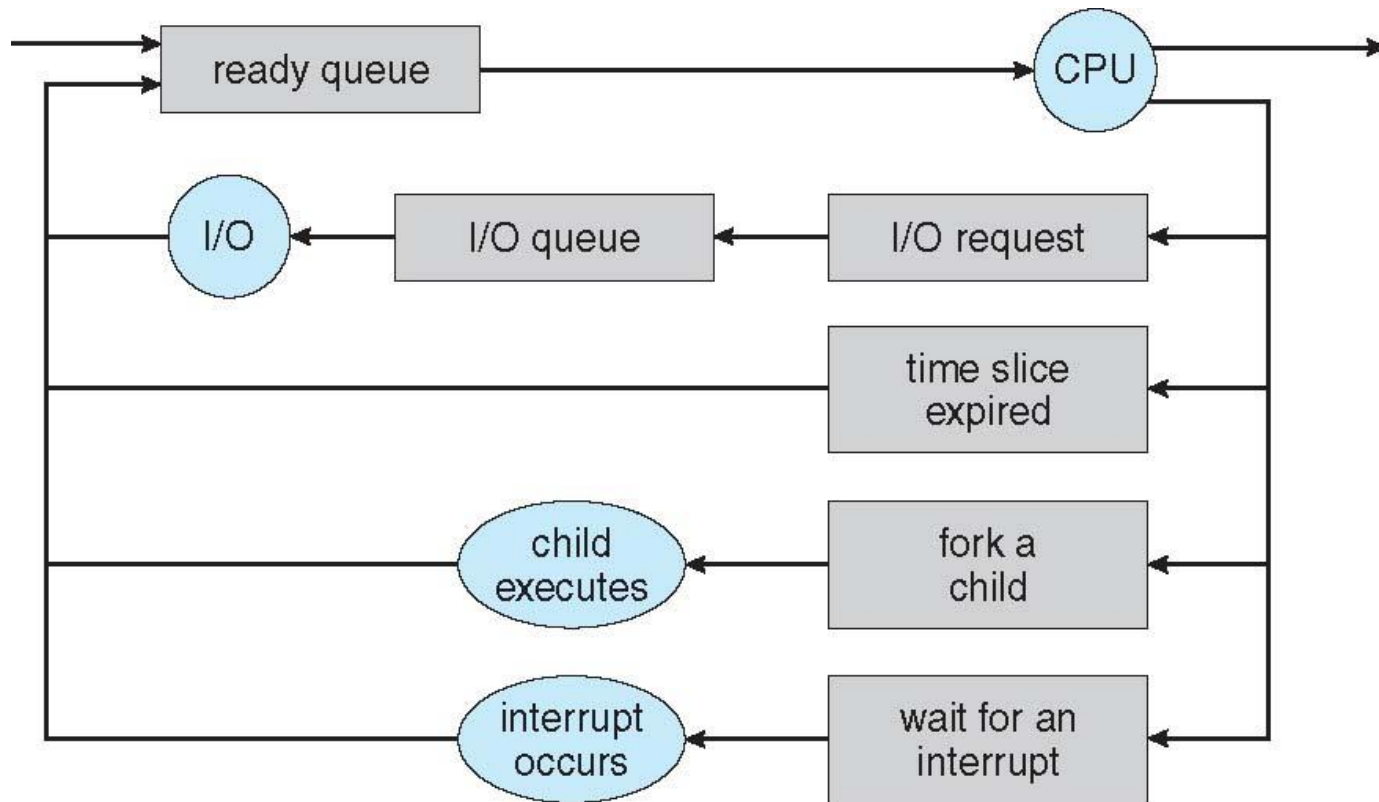
# Process Scheduling

---

- Maximize CPU use, quickly switch processes onto CPU for time sharing
- **Process scheduler** selects among available processes for next execution on CPU
- Maintains **scheduling queues** of processes
  - **Job queue** – set of all processes in the system
  - **Ready queue** – set of all processes residing in main memory, ready and waiting to execute
  - **Device queues** – set of processes waiting for an I/O device
  - Processes migrate among the various queues

# Representation of Process Scheduling

- **Queueing diagram** represents queues, resources, flows



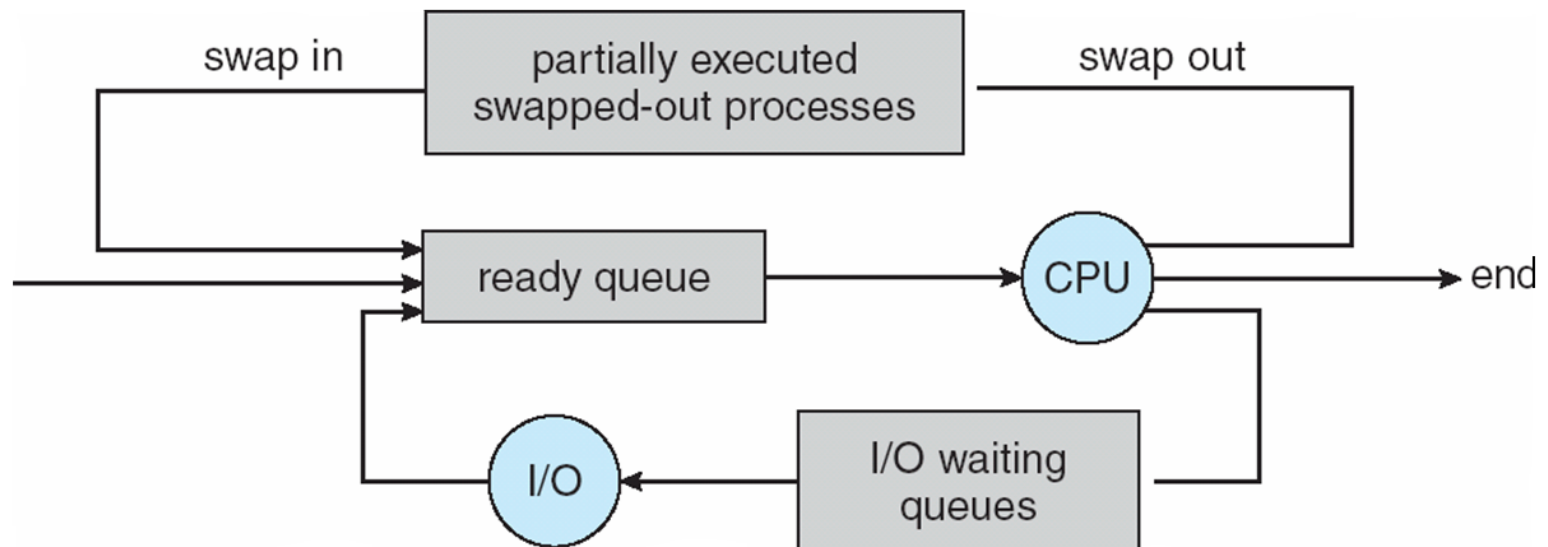
# Schedulers

---

- **Short-term scheduler** (or **CPU scheduler**) – selects which process should be executed next and allocates CPU
  - Sometimes the only scheduler in a system
  - Short-term scheduler is invoked frequently (milliseconds)  $\Rightarrow$  (must be fast)
- **Long-term scheduler** (or **job scheduler**) – selects which processes should be brought into the ready queue
  - Long-term scheduler is invoked infrequently (seconds, minutes)  $\Rightarrow$  (may be slow)
  - The long-term scheduler controls the **degree of multiprogramming**
- Processes can be described as either:
  - **I/O-bound process** – spends more time doing I/O than computations, many short CPU bursts
  - **CPU-bound process** – spends more time doing computations; few very long CPU bursts
- Long-term scheduler strives for good ***process mix***

# Addition of Medium Term Scheduling

- **Medium-term scheduler** can be added if degree of multiple programming needs to decrease
  - Remove process from memory, store on disk, bring back in from disk to continue execution: **swapping**



# Context Switch

---

- When CPU switches to another process, the system must **save the state** of the old process and load the **saved state** for the new process via a **context switch**
- **Context** of a process represented in the PCB
- Context-switch time is overhead; the system does no useful work while switching
  - The more complex the OS and the PCB the longer the context switch
- Time dependent on hardware support
  - Some hardware provides multiple sets of registers per CPU multiple contexts loaded at once

# Operations on Processes

---

- System must provide mechanisms for:
  - process creation,
  - process termination,
  - and so on as detailed next

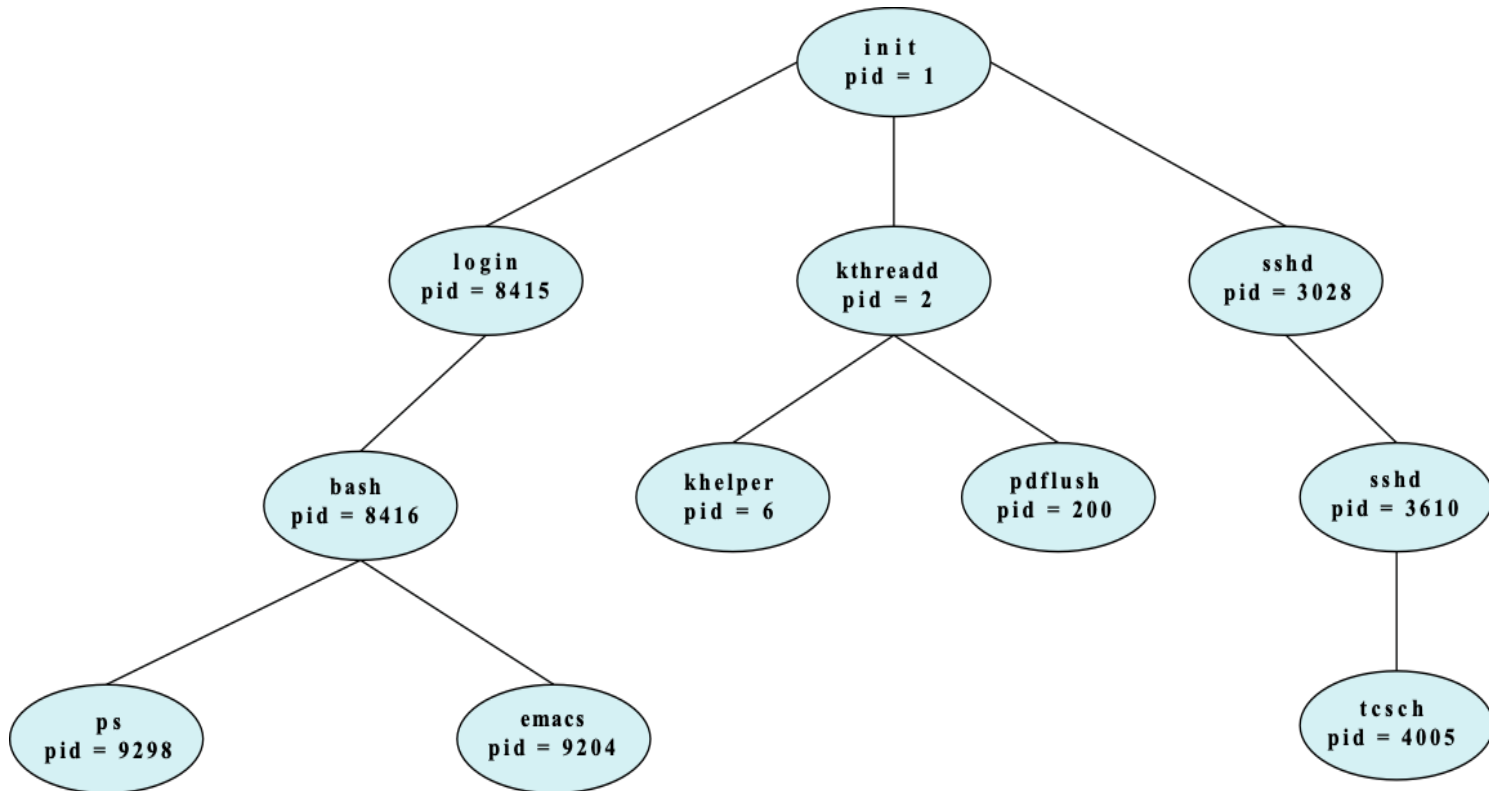


# Process Creation

---

- **Parent** process create **children** processes, which, in turn create other processes, forming a **tree** of processes
- Generally, process identified and managed via a **process identifier** (**pid**)
- Resource sharing options
  - Parent and children share all resources
  - Children share subset of parent's resources
  - Parent and child share no resources
- Execution options
  - Parent and children execute concurrently
  - Parent waits until children terminate

# A Tree of Processes in Linux



# Process Termination

---

- Process executes last statement and then asks the operating system to delete it using the **exit()** system call.
  - Returns status data from child to parent (via **wait()**)
  - Process' resources are deallocated by operating system
- Parent may terminate the execution of children processes using the **abort()** system call. Some reasons for doing so:
  - Child has exceeded allocated resources
  - Task assigned to child is no longer required
  - The parent is exiting and the operating systems does not allow a child to continue if its parent terminates

# Process Termination

---

- Some operating systems do not allow child to exist if its parent has terminated. If a process terminates, then all its children must also be terminated.
  - **cascading termination.** All children, grandchildren, etc. are terminated.
  - The termination is initiated by the operating system.
- The parent process may wait for termination of a child process by using the `wait()` system call. The call returns status information and the pid of the terminated process

```
pid = wait(&status);
```
- If no parent waiting (did not invoke `wait()`) process is a **zombie**
- If parent terminated without invoking `wait`, process is an **orphan**

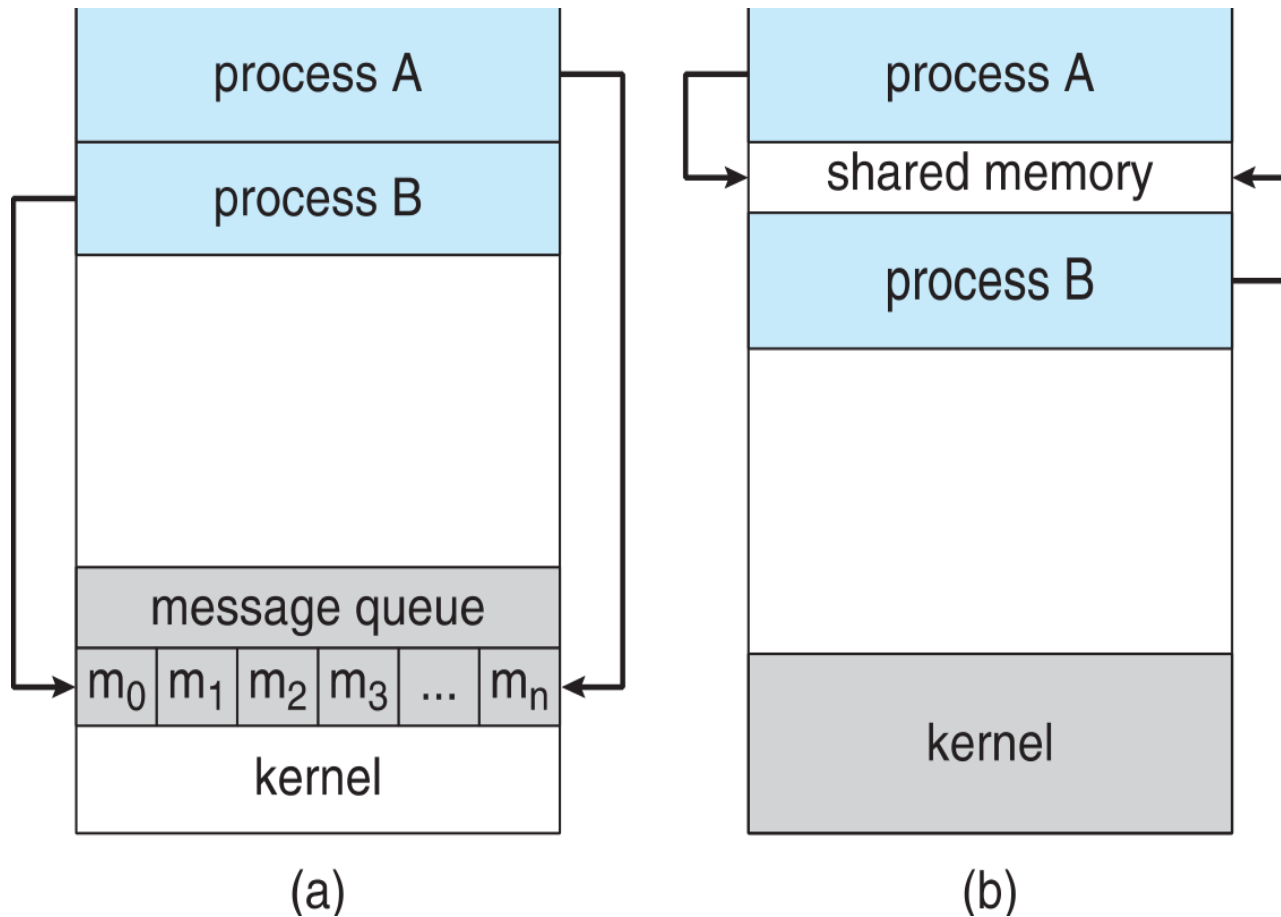
# Interprocess Communication

---

- Processes within a system may be *independent* or *cooperating*
- Cooperating process can affect or be affected by other processes, including sharing data
- Reasons for cooperating processes:
  - Information sharing
  - Computation speedup
  - Modularity
  - Convenience
- Cooperating processes need **interprocess communication (IPC)**
- Two models of IPC
  - **Shared memory**
  - **Message passing**

# Communications Models

(a) Message passing. (b) shared memory.



# Cooperating Processes

---

- **Independent** process cannot affect or be affected by the execution of another process
- **Cooperating** process can affect or be affected by the execution of another process
- Advantages of process cooperation
  - Information sharing
  - Computation speed-up
  - Modularity
  - Convenience

# Producer-Consumer Problem

---

- Paradigm for cooperating processes, *producer* process produces information that is consumed by a *consumer* process
  - **unbounded-buffer** places no practical limit on the size of the buffer
  - **bounded-buffer** assumes that there is a fixed buffer size



# Interprocess Communication – Shared Memory

---

- An area of memory shared among the processes that wish to communicate
- The communication is under the control of the users processes not the operating system.
- Major issues is to provide mechanism that will allow the user processes to synchronize their actions when they access shared memory.

# Interprocess Communication – Message Passing

---

- Mechanism for processes to communicate and to synchronize their actions
- Message system – processes communicate with each other without resorting to shared variables
- IPC facility provides two operations:
  - **send**(*message*)
  - **receive**(*message*)
- The *message* size is either fixed or variable

# Message Passing (Cont.)

---

- If processes  $P$  and  $Q$  wish to communicate, they need to:
  - Establish a **communication link** between them
  - Exchange messages via send/receive
- Implementation issues:
  - How are links established?
  - Can a link be associated with more than two processes?
  - How many links can there be between every pair of communicating processes?
  - What is the capacity of a link?
  - Is the size of a message that the link can accommodate fixed or variable?
  - Is a link unidirectional or bi-directional?

# Message Passing (Cont.)

---

- Implementation of communication link
  - Physical:
    - Shared memory
    - Hardware bus
    - Network
  - Logical:
    - Direct or indirect
    - Synchronous or asynchronous
    - Automatic or explicit buffering

# Direct Communication

---

- Processes must name each other explicitly:
  - **send** (*P*, *message*) – send a message to process P
  - **receive**(*Q*, *message*) – receive a message from process Q
- Properties of communication link
  - Links are established automatically
  - A link is associated with exactly one pair of communicating processes
  - Between each pair there exists exactly one link
  - The link may be unidirectional, but is usually bi-directional

# Indirect Communication

---

- Messages are directed and received from mailboxes (also referred to as ports)
  - Each mailbox has a unique id
  - Processes can communicate only if they share a mailbox
- Properties of communication link
  - Link established only if processes share a common mailbox
  - A link may be associated with many processes
  - Each pair of processes may share several communication links
  - Link may be unidirectional or bi-directional

# Indirect Communication

---

- Operations
  - create a new mailbox (port)
  - send and receive messages through mailbox
  - destroy a mailbox
- Primitives are defined as:
  - send**(*A*, *message*) – send a message to mailbox *A*
  - receive**(*A*, *message*) – receive a message from mailbox *A*

# Indirect Communication

---

- Mailbox sharing
  - $P_1$ ,  $P_2$ , and  $P_3$  share mailbox A
  - $P_1$  sends;  $P_2$  and  $P_3$  receive
  - Who gets the message?
- Solutions
  - Allow a link to be associated with at most two processes
  - Allow only one process at a time to execute a receive operation
  - Allow the system to select arbitrarily the receiver.  
Sender is notified who the receiver was.



# Synchronization

---

- Message passing may be either blocking or non-blocking
- **Blocking** is considered **synchronous**
  - **Blocking send** -- the sender is blocked until the message is received
  - **Blocking receive** -- the receiver is blocked until a message is available
- **Non-blocking** is considered **asynchronous**
  - **Non-blocking send** -- the sender sends the message and continue
  - **Non-blocking receive** -- the receiver receives:
    - A valid message, or
    - Null message
- Different combinations possible
  - If both send and receive are blocking, we have a **rendezvous**

# Buffering

---

- Queue of messages attached to the link.
- implemented in one of three ways
  1. Zero capacity – no messages are queued on a link.  
Sender must wait for receiver (rendezvous)
  2. Bounded capacity – finite length of  $n$  messages  
Sender must wait if link full
  3. Unbounded capacity – infinite length  
Sender never waits

# Examples of IPC Systems - POSIX

---

- POSIX Shared Memory

- Process first creates shared memory segment

- ```
shm_fd = shm_open(name, O_CREAT | O_RDWR, 0666);
```

- Also used to open an existing segment to share it

- Set the size of the object

- ```
ftruncate(shm_fd, 4096);
```

- Now the process could write to the shared memory

- ```
    sprintf(shared_memory, "Writing to shared  
memory");
```

# IPC POSIX Producer

---

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <fcntl.h>
#include <sys/shm.h>
#include <sys/stat.h>

int main()
{
    /* the size (in bytes) of shared memory object */
    const int SIZE = 4096;
    /* name of the shared memory object */
    const char *name = "OS";
    /* strings written to shared memory */
    const char *message_0 = "Hello";
    const char *message_1 = "World!";

    /* shared memory file descriptor */
    int shm_fd;
    /* pointer to shared memory object */
    void *ptr;

    /* create the shared memory object */
    shm_fd = shm_open(name, O_CREAT | O_RDWR, 0666);

    /* configure the size of the shared memory object */
    ftruncate(shm_fd, SIZE);

    /* memory map the shared memory object */
    ptr = mmap(0, SIZE, PROT_WRITE, MAP_SHARED, shm_fd, 0);

    /* write to the shared memory object */
    sprintf(ptr,"%s",message_0);
    ptr += strlen(message_0);
    sprintf(ptr,"%s",message_1);
    ptr += strlen(message_1);

    return 0;
}
```

# IPC POSIX Consumer

---

```
#include <stdio.h>
#include <stdlib.h>
#include <fcntl.h>
#include <sys/shm.h>
#include <sys/stat.h>

int main()
{
    /* the size (in bytes) of shared memory object */
    const int SIZE = 4096;
    /* name of the shared memory object */
    const char *name = "OS";
    /* shared memory file descriptor */
    int shm_fd;
    /* pointer to shared memory object */
    void *ptr;

    /* open the shared memory object */
    shm_fd = shm_open(name, O_RDONLY, 0666);

    /* memory map the shared memory object */
    ptr = mmap(0, SIZE, PROT_READ, MAP_SHARED, shm_fd, 0);

    /* read from the shared memory object */
    printf("%s", (char *)ptr);

    /* remove the shared memory object */
    shm_unlink(name);

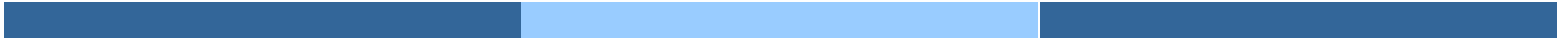
    return 0;
}
```

# Examples of IPC Systems - Mach

---

- Mach communication is message based
  - Even system calls are messages
  - Each task gets two mailboxes at creation- Kernel and Notify
  - Only three system calls needed for message transfer  
`msg_send()` , `msg_receive()` , `msg_rpc()`
  - Mailboxes needed for communication, created via  
`port_allocate()`
- Send and receive are flexible, for example four options if mailbox full:
  - Wait indefinitely
  - Wait at most n milliseconds
  - Return immediately
  - Temporarily cache a message

# End of Chapter 3



# Chapter 6: CPU Scheduling





# Chapter 6: CPU Scheduling

---

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms

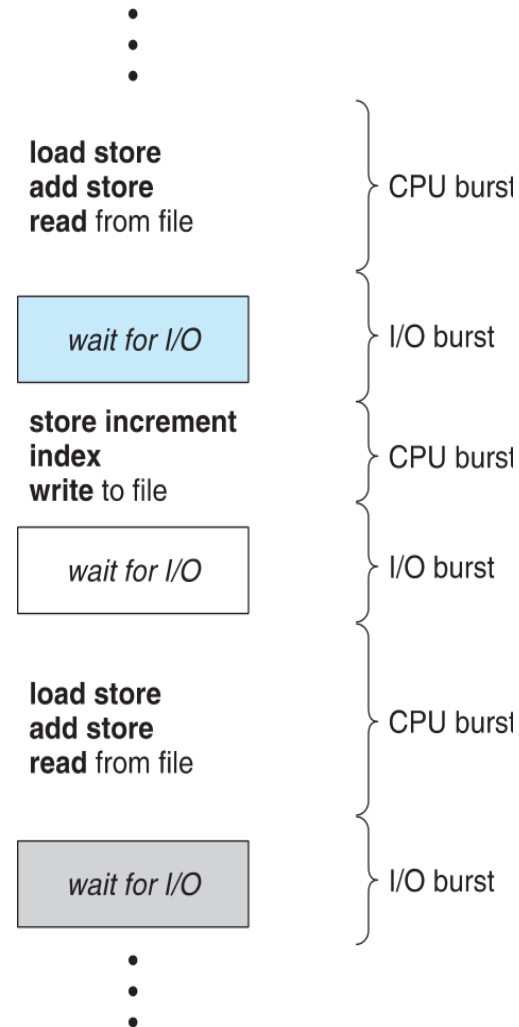
# Objectives

---

- To introduce CPU scheduling, which is the basis for multiprogrammed operating systems
- To describe various CPU-scheduling algorithms
- To discuss evaluation criteria for selecting a CPU-scheduling algorithm for a particular system

# Basic Concepts

- Maximum CPU utilization obtained with multiprogramming
- CPU–I/O Burst Cycle – 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



# 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
  3. Switches from waiting to ready
  4. Terminates
- Scheduling under 1 and 4 is **nonpreemptive**
- All other scheduling is **preemptive**
  - Consider access to shared data
  - Consider preemption while in kernel mode
  - Consider interrupts occurring during crucial OS activities

# 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** – keep the CPU as busy as possible
- **Throughput** – # of processes that complete their execution per time unit
- **Turnaround time** – amount of time to execute a particular process
- **Waiting time** – amount of time a process has been waiting in the ready queue
- **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 Criteria

---

- » **Arrival Time:** Time at which the process arrives in the ready queue.

**Completion Time:** Time at which process completes its execution.

**Burst Time:** Time required by a process for CPU execution.

**Turn Around Time:** Time Difference between completion time and arrival time.

$\text{Turn Around Time} = \text{Completion Time} - \text{Arrival Time}$

- » **Waiting Time(W.T):** Time Difference between turn around time and burst time.

$\text{Waiting Time} = \text{Turn Around Time} - \text{Burst Time}$

# Scheduling Algorithm Optimization Criteria

---

- Max CPU utilization
- Max throughput
- Min turnaround time
- Min waiting time
- Min response time

$$TAT = CT - AT$$

$$TAT = WT + BT$$

$$WT = ST - AT$$

$$WT = TAT - BT$$



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

# Shortest-Job-First (SJF) Scheduling

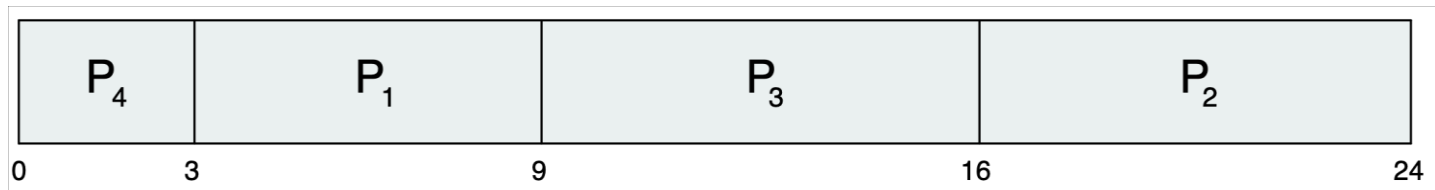
---

- 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



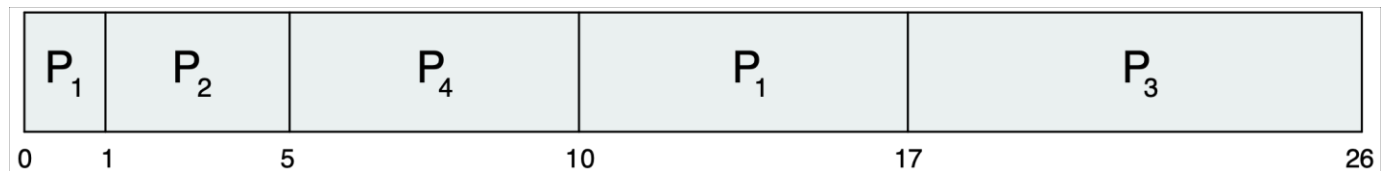
□ Average waiting time =  $(3 + 16 + 9 + 0) / 4 = 7$

# 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



- Average waiting time =  $[(10-1)+(1-1)+(17-2)+5-3]/4 = 26/4 = 6.5$  msec

# Priority Scheduling

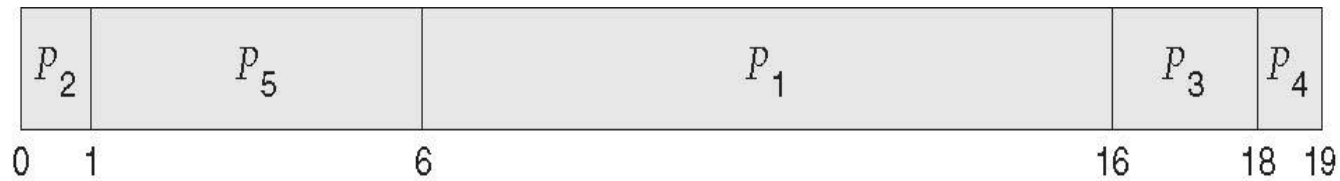
---

- 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 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



□ Average waiting time = 8.2 msec

# Round Robin (RR)

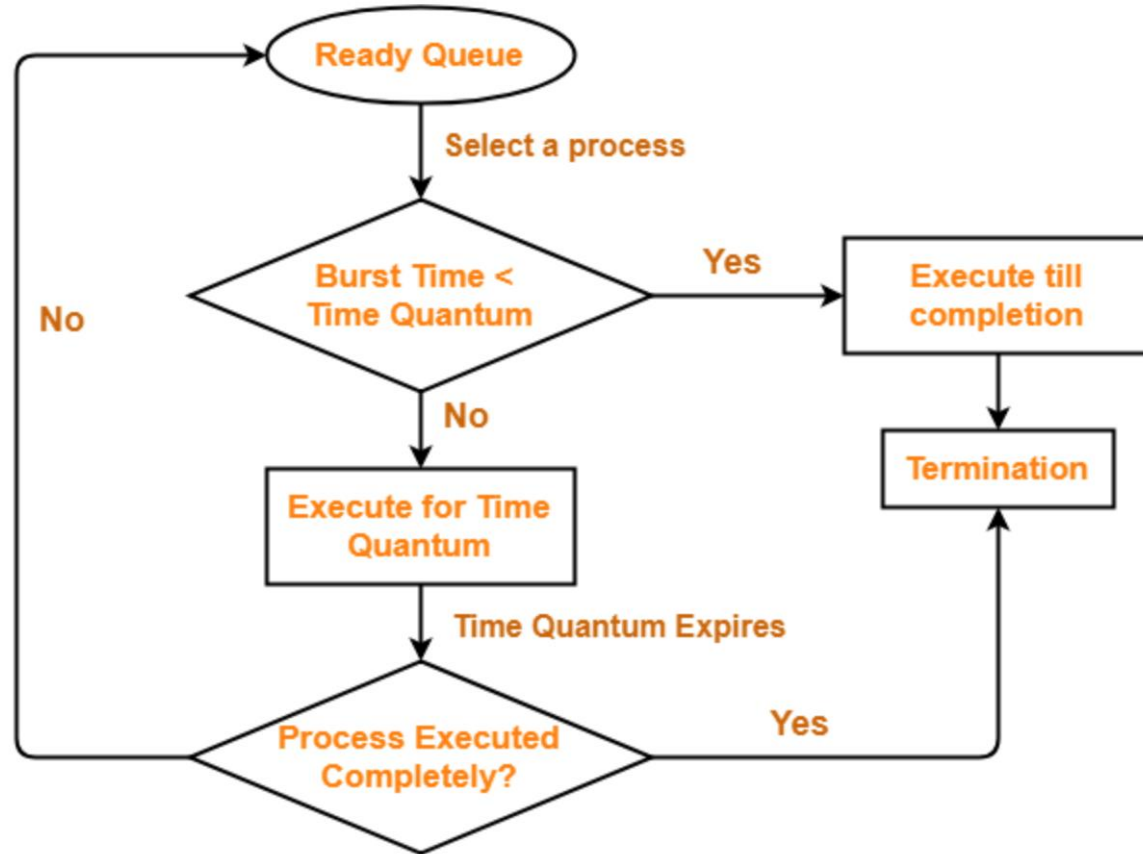
---

- 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.
- If there are  $n$  processes in the ready queue and the time quantum is  $q$ , then each process gets  $1/n$  of the CPU time in chunks of at most  $q$  time units at once. No process waits more than  $(n-1)q$  time units.
- Timer interrupts every quantum to schedule next process
- Performance
  - $q$  large  $\Rightarrow$  FIFO
  - $q$  small  $\Rightarrow q$  must be large with respect to context switch, otherwise overhead is too high



# Round Robin (RR)

Round Robin Scheduling is FCFS Scheduling with preemptive mode

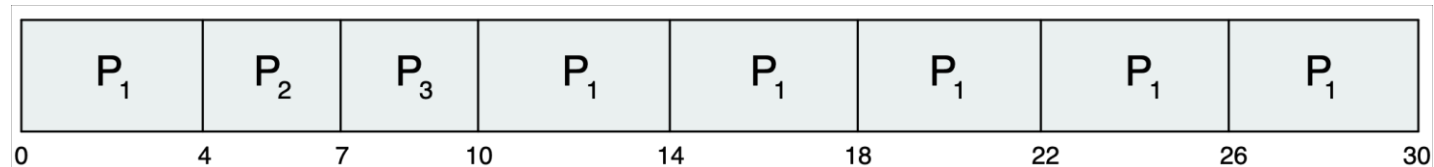


**Round Robin Scheduling**

# 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

---

Consider the set of 5 processes whose arrival time and burst time are given below-

| Process Id | Arrival time | Burst time |
|------------|--------------|------------|
| P1         | 0            | 5          |
| P2         | 1            | 3          |
| P3         | 2            | 1          |
| P4         | 3            | 2          |
| P5         | 4            | 3          |

If the CPU scheduling policy is Round Robin with time quantum = 2 unit, calculate the average waiting time and average turn around time.

Consider the set of 5 processes whose arrival time and burst time are given below-

| Process Id | Arrival time | Burst time |
|------------|--------------|------------|
| P1         | 0            | 5          |
| P2         | 1            | 3          |
| P3         | 2            | 1          |
| P4         | 3            | 2          |
| P5         | 4            | 3          |

If the CPU scheduling policy is Round Robin with time quantum = 2 unit, calculate the average waiting time and average turn around time.

### Ready Queue-

P5, P1, P2, P5, P4, P1, P3, P2, P1



Round Robin Scheduling | Problem-01 | Gantt Chart

### Gantt Chart

### Ready Queue-

P5, P1, P2, P5, P4, P1, P3, P2, P1

- Turn Around time = Exit time – Arrival time
- Waiting time = Turn Around time – Burst time



Round Robin Scheduling | Problem-01 | Gantt Chart

**Gantt Chart**

| Process Id | Exit time | Turn Around time | Waiting time |
|------------|-----------|------------------|--------------|
| P1         | 13        | $13 - 0 = 13$    | $13 - 5 = 8$ |
| P2         | 12        | $12 - 1 = 11$    | $11 - 3 = 8$ |
| P3         | 5         | $5 - 2 = 3$      | $3 - 1 = 2$  |
| P4         | 9         | $9 - 3 = 6$      | $6 - 2 = 4$  |
| P5         | 14        | $14 - 4 = 10$    | $10 - 3 = 7$ |

Now,

- Average Turn Around time =  $(13 + 11 + 3 + 6 + 10) / 5 = 43 / 5 = 8.6$  unit
- Average waiting time =  $(8 + 8 + 2 + 4 + 7) / 5 = 29 / 5 = 5.8$  unit

---

Consider the set of 6 processes whose arrival time and burst time are given below-

| Process Id | Arrival time | Burst time |
|------------|--------------|------------|
| P1         | 0            | 4          |
| P2         | 1            | 5          |
| P3         | 2            | 2          |
| P4         | 3            | 1          |
| P5         | 4            | 6          |
| P6         | 6            | 3          |

If the CPU scheduling policy is Round Robin with time quantum = 2, calculate the average waiting time and average turn around time.

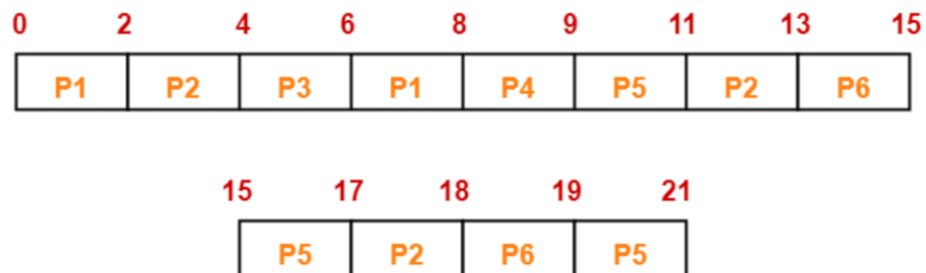
Consider the set of 6 processes whose arrival time and burst time are given below-

| Process Id | Arrival time | Burst time |
|------------|--------------|------------|
| P1         | 0            | 4          |
| P2         | 1            | 5          |
| P3         | 2            | 2          |
| P4         | 3            | 1          |
| P5         | 4            | 6          |
| P6         | 6            | 3          |

If the CPU scheduling policy is Round Robin with time quantum = 2, calculate the average waiting time and average turn around time.

**Ready Queue-**

P5, P6, P2, P5, P6, P2, P5, P4, P1, P3, P2, P1



**Gantt Chart**

Round Robin Scheduling | Problem-02 | Gantt Chart

---

Now, we know-

- Turn Around time = Exit time – Arrival time
- Waiting time = Turn Around time – Burst time

| Process Id | Exit time | Turn Around time | Waiting time  |
|------------|-----------|------------------|---------------|
| P1         | 8         | $8 - 0 = 8$      | $8 - 4 = 4$   |
| P2         | 18        | $18 - 1 = 17$    | $17 - 5 = 12$ |
| P3         | 6         | $6 - 2 = 4$      | $4 - 2 = 2$   |
| P4         | 9         | $9 - 3 = 6$      | $6 - 1 = 5$   |
| P5         | 21        | $21 - 4 = 17$    | $17 - 6 = 11$ |
| P6         | 19        | $19 - 6 = 13$    | $13 - 3 = 10$ |

Now,

- Average Turn Around time =  $(8 + 17 + 4 + 6 + 17 + 13) / 6 = 65 / 6 = 10.84$  unit
- Average waiting time =  $(4 + 12 + 2 + 5 + 11 + 10) / 6 = 44 / 6 = 7.33$  unit



---

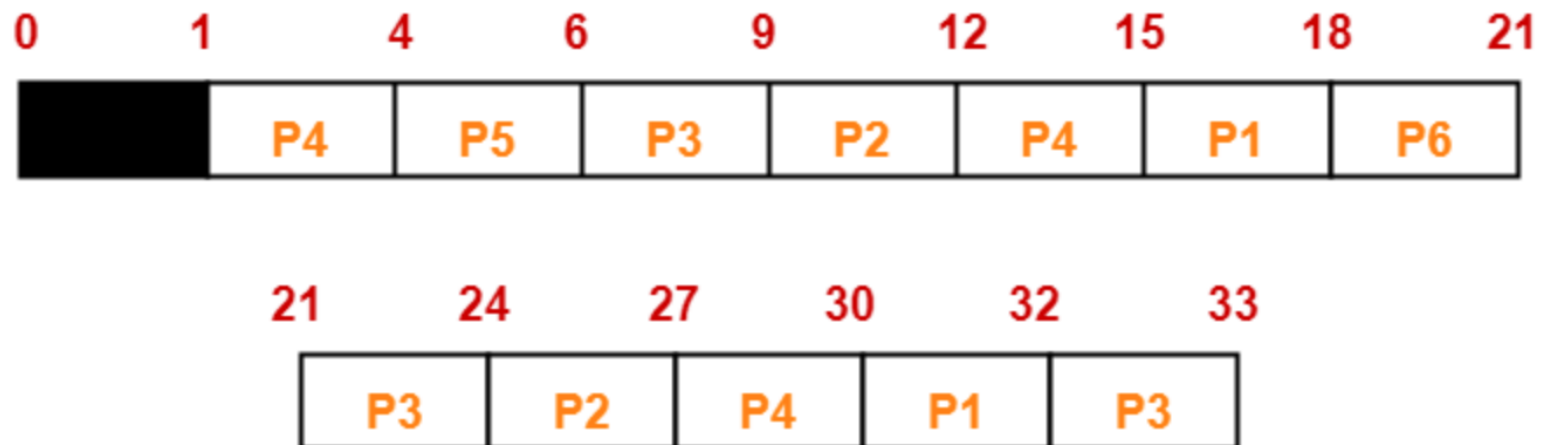
Consider the set of 6 processes whose arrival time and burst time are given below-

| Process Id | Arrival time | Burst time |
|------------|--------------|------------|
| P1         | 5            | 5          |
| P2         | 4            | 6          |
| P3         | 3            | 7          |
| P4         | 1            | 9          |
| P5         | 2            | 2          |
| P6         | 6            | 3          |

If the CPU scheduling policy is Round Robin with time quantum = 3, calculate the average waiting time and average turn around time.

### Ready Queue-

P3, P1, P4, P2, P3, P6, P1, P4, P2, P3, P5, P4



**Gantt Chart**

---

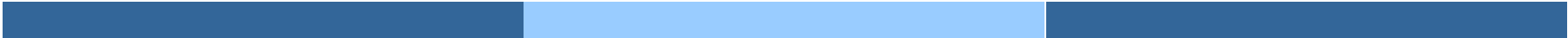
Now, we know-

- Turn Around time = Exit time – Arrival time
- Waiting time = Turn Around time – Burst time

| Process Id | Exit time | Turn Around time | Waiting time  |
|------------|-----------|------------------|---------------|
| P1         | 32        | $32 - 5 = 27$    | $27 - 5 = 22$ |
| P2         | 27        | $27 - 4 = 23$    | $23 - 6 = 17$ |
| P3         | 33        | $33 - 3 = 30$    | $30 - 7 = 23$ |
| P4         | 30        | $30 - 1 = 29$    | $29 - 9 = 20$ |
| P5         | 6         | $6 - 2 = 4$      | $4 - 2 = 2$   |
| P6         | 21        | $21 - 6 = 15$    | $15 - 3 = 12$ |

Now,

- Average Turn Around time =  $(27 + 23 + 30 + 29 + 4 + 15) / 6 = 128 / 6 = 21.33$  unit
- Average waiting time =  $(22 + 17 + 23 + 20 + 2 + 12) / 6 = 96 / 6 = 16$  unit



<https://www.gatevidyalay.com/round-robin-round-robin-scheduling-examples/>

## Shortest Remaining Time First(SRTF)

| Process        | Burst Time(CPU) | Arrival Time(ms) |
|----------------|-----------------|------------------|
| P <sub>1</sub> | 3               | 0                |
| P <sub>2</sub> | 6               | 2                |
| P <sub>3</sub> | 4               | 4                |
| P <sub>4</sub> | 5               | 6                |
| P <sub>5</sub> | 2               | 8                |



### Average Waiting Time and Turnaround Time

#### Average Waiting Time

First of all, we have to find the waiting time for each process.

Waiting Time of process

$$P_1 = 0\text{ms}$$

$$P_2 = (3 - 2) + (10 - 4) = 7\text{ms}$$

$$P_3 = (4 - 4) = 0\text{ms}$$

$$P_4 = (15 - 6) = 9\text{ms}$$

$$P_5 = (8 - 8) = 0\text{ms}$$

$$\text{Therefore, Average Waiting Time} = (0 + 7 + 0 + 9 + 0) / 5 = 3.2\text{ms}$$

## Shortest Remaining Time First(SRTF)

| Process        | Burst Time(CPU) | Arrival Time(ms) |
|----------------|-----------------|------------------|
| P <sub>1</sub> | 3               | 0                |
| P <sub>2</sub> | 6               | 2                |
| P <sub>3</sub> | 4               | 4                |
| P <sub>4</sub> | 5               | 6                |
| P <sub>5</sub> | 2               | 8                |

|                |                |                |                |                |                |    |
|----------------|----------------|----------------|----------------|----------------|----------------|----|
| P <sub>1</sub> | P <sub>2</sub> | P <sub>3</sub> | P <sub>5</sub> | P <sub>2</sub> | P <sub>4</sub> |    |
| 0              | 3              | 4              | 8              | 10             | 15             | 20 |

*SRTF Scheduling*

### Average Turnaround Time

First of all, we have to find the turnaround time of each process.

Turnaround Time of process

$$P_1 = (0 + 3) = 3\text{ms}$$

$$P_2 = (7 + 6) = 13\text{ms}$$

$$P_3 = (0 + 4) = 4\text{ms}$$

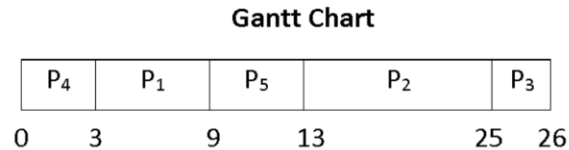
$$P_4 = (9 + 5) = 14\text{ms}$$

$$P_5 = (0 + 2) = 2\text{ms}$$

$$\text{Therefore, Average Turnaround Time} = (3 + 13 + 4 + 14 + 2) / 5 = 7.2\text{ms}$$

## Priority Scheduling Example

| Process        | CPU Burst Time | Priority |
|----------------|----------------|----------|
| P <sub>1</sub> | 6              | 2        |
| P <sub>2</sub> | 12             | 4        |
| P <sub>3</sub> | 1              | 5        |
| P <sub>4</sub> | 3              | 1        |
| P <sub>5</sub> | 4              | 3        |



*Priority Scheduling Example with Gantt Chart*

## Average Waiting Time and Turnaround Time

### Average Waiting Time

First of all, we have to find out the waiting time of each process.

Waiting Time of process

P<sub>1</sub> = 3ms

P<sub>2</sub> = 13ms

P<sub>3</sub> = 25ms

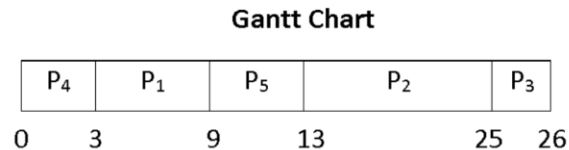
P<sub>4</sub> = 0ms

P<sub>5</sub> = 9ms

Therefore, Average Waiting Time =  $(3 + 13 + 25 + 0 + 9) / 5 = 10\text{ms}$

## Priority Scheduling Example

| Process        | CPU Burst Time | Priority |
|----------------|----------------|----------|
| P <sub>1</sub> | 6              | 2        |
| P <sub>2</sub> | 12             | 4        |
| P <sub>3</sub> | 1              | 5        |
| P <sub>4</sub> | 3              | 1        |
| P <sub>5</sub> | 4              | 3        |



*Priority Scheduling Example with Gantt Chart*

## Average Turnaround Time

First finding Turnaround Time of each process.

Turnaround Time of process

$$P_1 = (3 + 6) = 9\text{ms}$$

$$P_2 = (13 + 12) = 25\text{ms}$$

$$P_3 = (25 + 1) = 26\text{ms}$$

$$P_4 = (0 + 3) = 3\text{ms}$$

$$P_5 = (9 + 4) = 13\text{ms}$$

$$\text{Therefore, Average Turnaround Time} = (9 + 25 + 26 + 3 + 13) / 5 = 15.2\text{ms}$$



---

| Process Id | Arrival time | Burst time | Priority |
|------------|--------------|------------|----------|
| P1         | 0            | 4          | 2        |
| P2         | 1            | 3          | 3        |
| P3         | 2            | 1          | 4        |
| P4         | 3            | 5          | 5        |
| P5         | 4            | 2          | 5        |

If the CPU scheduling policy is priority non-preemptive, calculate the average waiting time and average turn around time. (Higher number represents higher priority)

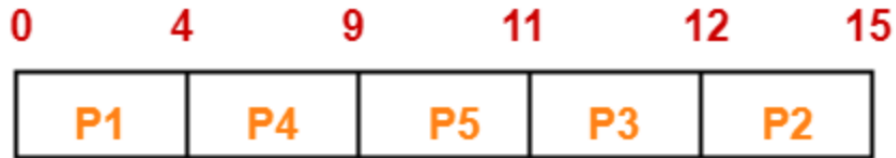
---

| Process Id | Arrival time | Burst time | Priority |
|------------|--------------|------------|----------|
| P1         | 0            | 4          | 2        |
| P2         | 1            | 3          | 3        |
| P3         | 2            | 1          | 4        |
| P4         | 3            | 5          | 5        |
| P5         | 4            | 2          | 5        |

If the CPU scheduling policy is priority non-preemptive, calculate the average waiting time and average turn around time. (Higher number represents higher priority)



**Gantt Chart**



**Gantt Chart**

| Process Id | Exit time | Turn Around time | Waiting time  |
|------------|-----------|------------------|---------------|
| P1         | 4         | $4 - 0 = 4$      | $4 - 4 = 0$   |
| P2         | 15        | $15 - 1 = 14$    | $14 - 3 = 11$ |
| P3         | 12        | $12 - 2 = 10$    | $10 - 1 = 9$  |
| P4         | 9         | $9 - 3 = 6$      | $6 - 5 = 1$   |
| P5         | 11        | $11 - 4 = 7$     | $7 - 2 = 5$   |

Now,

- Average Turn Around time =  $(4 + 14 + 10 + 6 + 7) / 5 = 41 / 5 = 8.2$  unit
- Average waiting time =  $(0 + 11 + 9 + 1 + 5) / 5 = 26 / 5 = 5.2$  unit

---

Consider the set of 5 processes whose arrival time and burst time are given below-

| Process Id | Arrival time | Burst time | Priority |
|------------|--------------|------------|----------|
| P1         | 0            | 4          | 2        |
| P2         | 1            | 3          | 3        |
| P3         | 2            | 1          | 4        |
| P4         | 3            | 5          | 5        |
| P5         | 4            | 2          | 5        |

If the CPU scheduling policy is priority preemptive, calculate the average waiting time and average turn around time. (Higher number represents higher priority)

Consider the set of 5 processes whose arrival time and burst time are given below-

| Process Id | Arrival time | Burst time | Priority |
|------------|--------------|------------|----------|
| P1         | 0            | 4          | 2        |
| P2         | 1            | 3          | 3        |
| P3         | 2            | 1          | 4        |
| P4         | 3            | 5          | 5        |
| P5         | 4            | 2          | 5        |



**Gantt Chart**



**Gantt Chart**

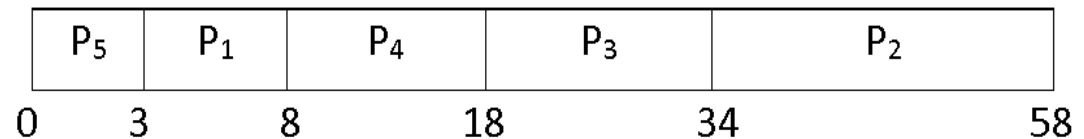
Now, we know-

- Turn Around time = Exit time – Arrival time
- Waiting time = Turn Around time – Burst time

| Process Id | Exit time | Turn Around time | Waiting time  |
|------------|-----------|------------------|---------------|
| P1         | 15        | $15 - 0 = 15$    | $15 - 4 = 11$ |
| P2         | 12        | $12 - 1 = 11$    | $11 - 3 = 8$  |
| P3         | 3         | $3 - 2 = 1$      | $1 - 1 = 0$   |
| P4         | 8         | $8 - 3 = 5$      | $5 - 5 = 0$   |
| P5         | 10        | $10 - 4 = 6$     | $6 - 2 = 4$   |

- Average Turn Around time =  $(15 + 11 + 1 + 5 + 6) / 5 = 38 / 5 = 7.6$  unit
- Average waiting time =  $(11 + 8 + 0 + 0 + 4) / 5 = 23 / 5 = 4.6$  unit

| Process        | Burst Time(ms) |
|----------------|----------------|
| P <sub>1</sub> | 5              |
| P <sub>2</sub> | 24             |
| P <sub>3</sub> | 16             |
| P <sub>4</sub> | 10             |
| P <sub>5</sub> | 3              |



## Average Waiting Time and Turnaround Time

### Average Waiting Time

We will apply the same formula to find average waiting time in this problem. Here arrival time is common to all processes(i.e., zero).

Waiting Time for

$$P_1 = 3 - 0 = 3\text{ms}$$

$$P_2 = 34 - 0 = 34\text{ms}$$

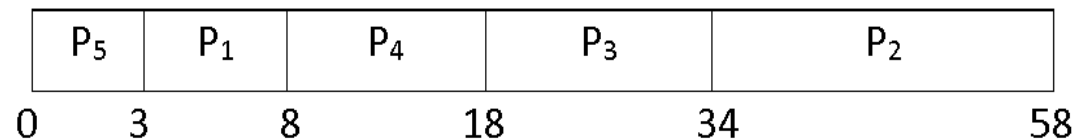
$$P_3 = 18 - 0 = 18\text{ms}$$

$$P_4 = 8 - 0 = 8\text{ms}$$

$$P_5 = 0\text{ms}$$

$$\text{Now, Average Waiting Time} = (3 + 34 + 18 + 8 + 0) / 5 = 12.6\text{ms}$$

| Process        | Burst Time(ms) |
|----------------|----------------|
| P <sub>1</sub> | 5              |
| P <sub>2</sub> | 24             |
| P <sub>3</sub> | 16             |
| P <sub>4</sub> | 10             |
| P <sub>5</sub> | 3              |



### Average Turnaround Time

According to the SJF Gantt chart and the turnaround time formulae,  
Turnaround Time of

$$P_1 = 3 + 5 = 8\text{ms}$$

$$P_2 = 34 + 24 = 58\text{ms}$$

$$P_3 = 18 + 16 = 34\text{ms}$$

$$P_4 = 8 + 10 = 18\text{ms}$$

$$P_5 = 0 + 3 = 3\text{ms}$$

$$\text{Therefore, Average Turnaround Time} = (8 + 58 + 34 + 18 + 3) / 5 = 24.2\text{ms}$$

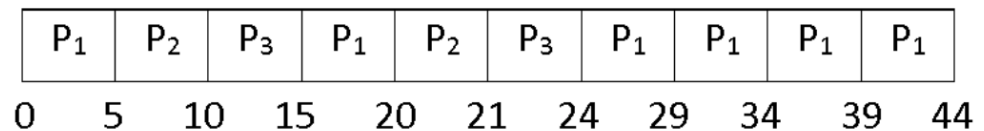


## Round Robin Scheduling Example

Here is the Round Robin scheduling example with gantt chart. Time Quantum is 5ms.

| Process        | CPU Burst Time |
|----------------|----------------|
| P <sub>1</sub> | 30             |
| P <sub>2</sub> | 6              |
| P <sub>3</sub> | 8              |

**Gantt Chart**



Round Robin Scheduling Example  
*Round Robin Scheduling*

### Average Waiting Time and Turnaround Time

#### Average Waiting Time

For finding Average Waiting Time, we have to find out the waiting time of each process.

Waiting Time of

$$P_1 = 0 + (15 - 5) + (24 - 20) = 14\text{ms}$$

$$P_2 = 5 + (20 - 10) = 15\text{ms}$$

$$P_3 = 10 + (21 - 15) = 16\text{ms}$$

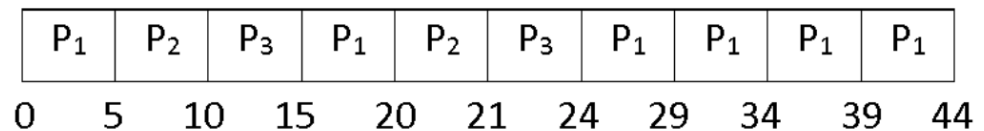
$$\text{Therefore, Average Waiting Time} = (14 + 15 + 16) / 3 = 15\text{ms}$$

## Round Robin Scheduling Example

Here is the Round Robin scheduling example with gantt chart. Time Quantum is 5ms.

| Process        | CPU Burst Time |
|----------------|----------------|
| P <sub>1</sub> | 30             |
| P <sub>2</sub> | 6              |
| P <sub>3</sub> | 8              |

**Gantt Chart**



Round Robin Scheduling Example  
*Round Robin Scheduling*

### Average Turnaround Time

Same concept for finding the Turnaround Time.

Turnaround Time of

$$P_1 = 14 + 30 = 44\text{ms}$$

$$P_2 = 15 + 6 = 21\text{ms}$$

$$P_3 = 16 + 8 = 24\text{ms}$$

Therefore, Average Turnaround Time =  $(44 + 21 + 24) / 3 = 29.66\text{ms}$

# Multilevel Queue

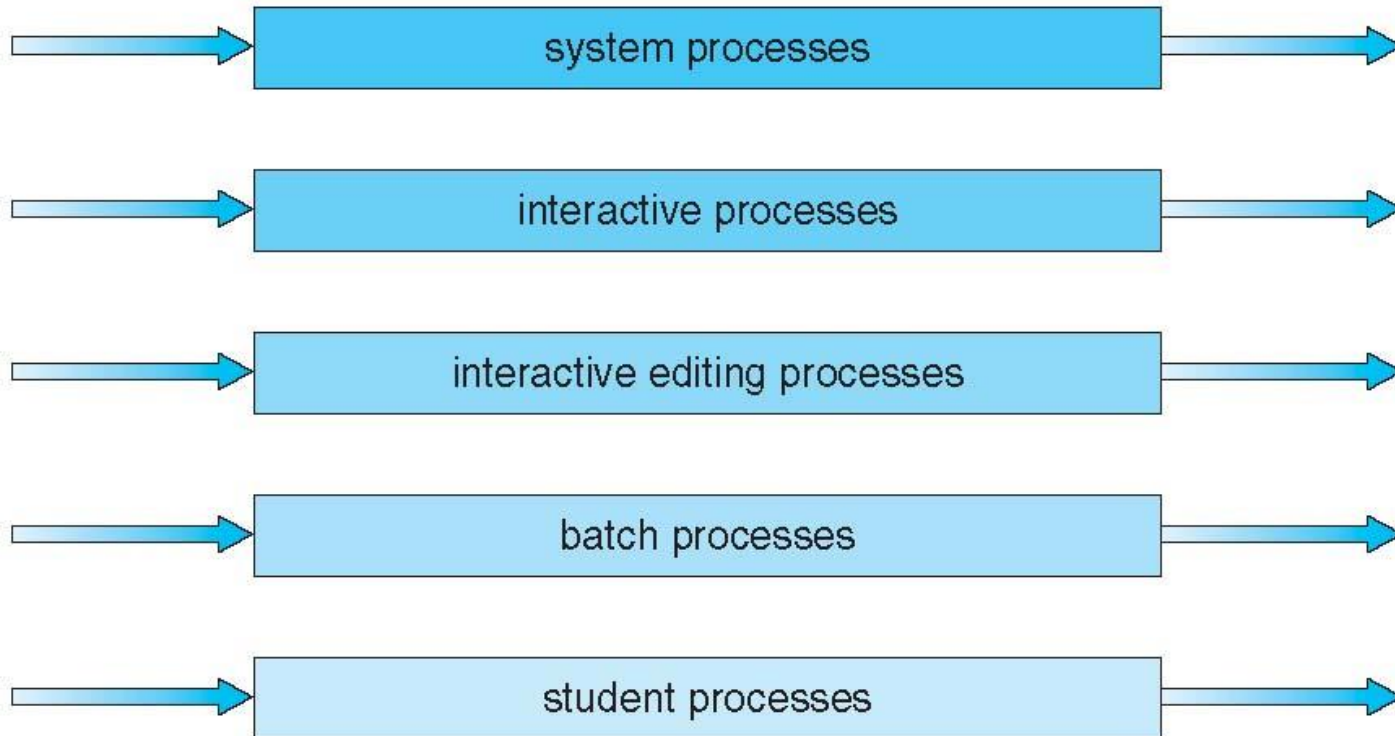
---

- Ready queue is partitioned into separate queues, eg:
  - **foreground** (interactive)
  - **background** (batch)
- Process permanently in a given 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

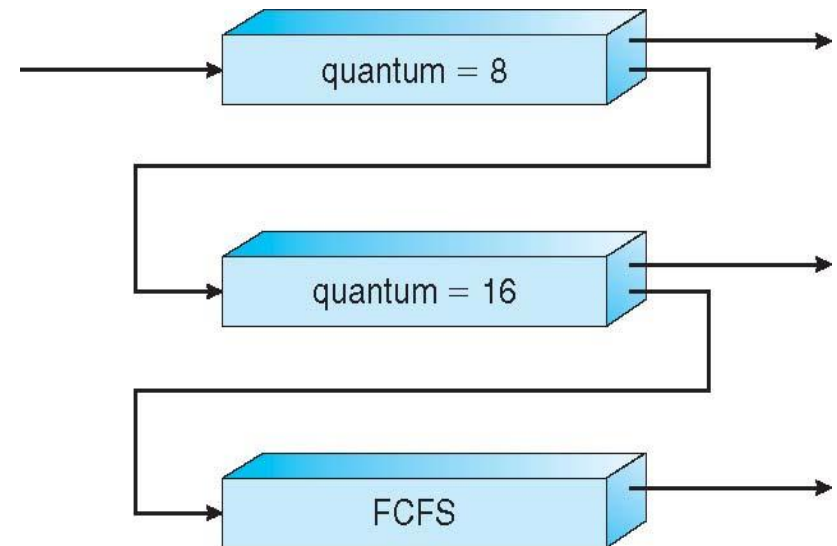
# Example of Multilevel Feedback Queue

## □ 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  $Q_2$



# Thread Scheduling

---

- Distinction between user-level and kernel-level threads
- When threads supported, threads scheduled, not processes
- Many-to-one and many-to-many models, thread library schedules user-level threads to run on LWP
  - Known as **process-contention scope (PCS)** since scheduling competition is within the process
  - Typically done via priority set by programmer
- Kernel thread scheduled onto available CPU is **system-contention scope (SCS)** – competition among all threads in system

# Pthread Scheduling

---

- API allows specifying either PCS or SCS during thread creation
  - PTHREAD\_SCOPE\_PROCESS schedules threads using PCS scheduling
  - PTHREAD\_SCOPE\_SYSTEM schedules threads using SCS scheduling
- Can be limited by OS – Linux and Mac OS X only allow PTHREAD\_SCOPE\_SYSTEM



# Pthread Scheduling API

---

```
#include <pthread.h>
#include <stdio.h>
#define NUM_THREADS 5
int main(int argc, char *argv[]) {
    int i, scope;
    pthread_t tid[NUM_THREADS];
    pthread_attr_t attr;
    /* get the default attributes */
    pthread_attr_init(&attr);
    /* first inquire on the current scope */
    if (pthread_attr_getscope(&attr, &scope) != 0)
        fprintf(stderr, "Unable to get scheduling scope\n");
    else {
        if (scope == PTHREAD_SCOPE_PROCESS)
            printf("PTHREAD_SCOPE_PROCESS");
        else if (scope == PTHREAD_SCOPE_SYSTEM)
            printf("PTHREAD_SCOPE_SYSTEM");
        else
            fprintf(stderr, "Illegal scope value.\n");
    }
}
```

# Pthread Scheduling API

---

```
/* set the scheduling algorithm to PCS or SCS */
pthread_attr_setscope(&attr, PTHREAD_SCOPE_SYSTEM);
/* create the threads */
for (i = 0; i < NUM_THREADS; i++)
    pthread_create(&tid[i], &attr, runner, NULL);
/* now join on each thread */
for (i = 0; i < NUM_THREADS; i++)
    pthread_join(tid[i], NULL);
}

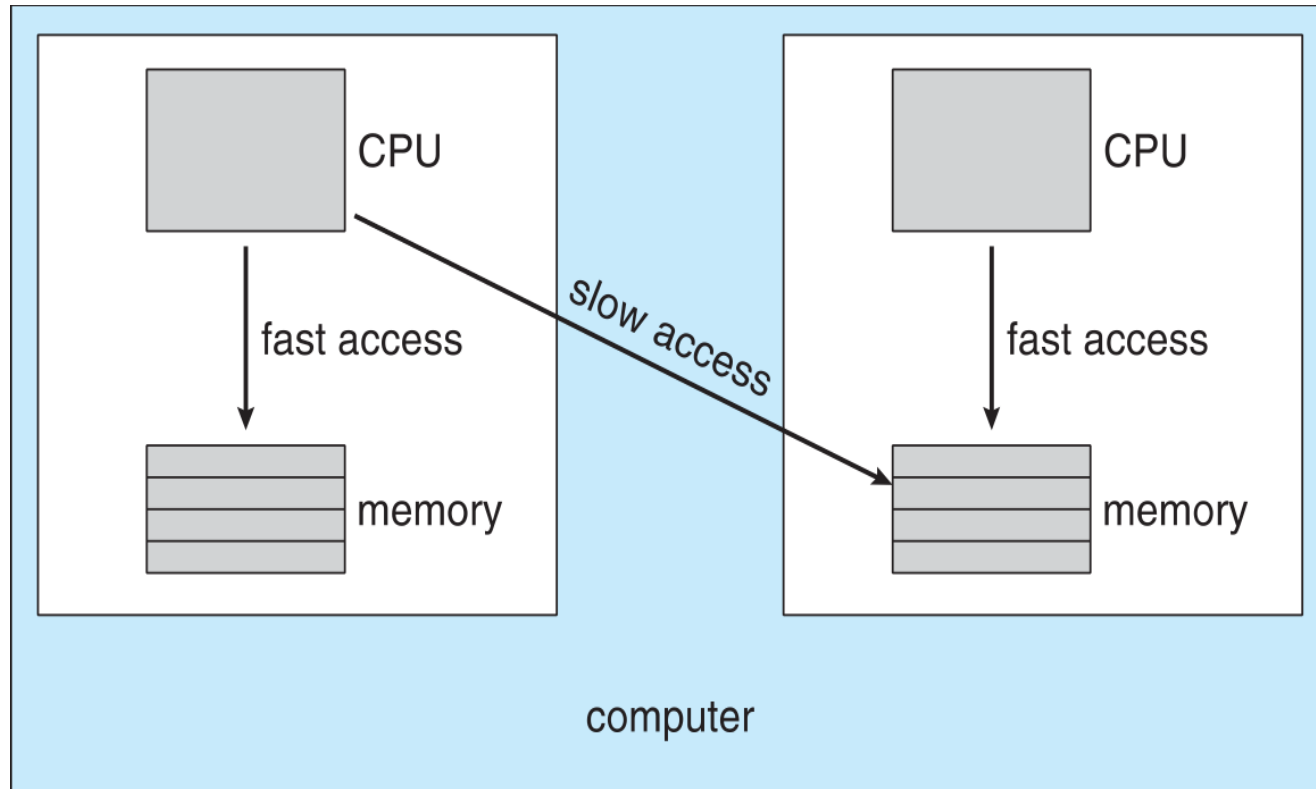
/* Each thread will begin control in this function */
void *runner(void *param)
{
    /* do some work ... */
    pthread_exit(0);
}
```

# 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** – process has affinity for processor on which it is currently running
  - **soft affinity**
  - **hard affinity**
  - Variations including **processor sets**

# NUMA and CPU Scheduling



Note that memory-placement algorithms can also consider affinity

# Multiple-Processor Scheduling – Load Balancing

---

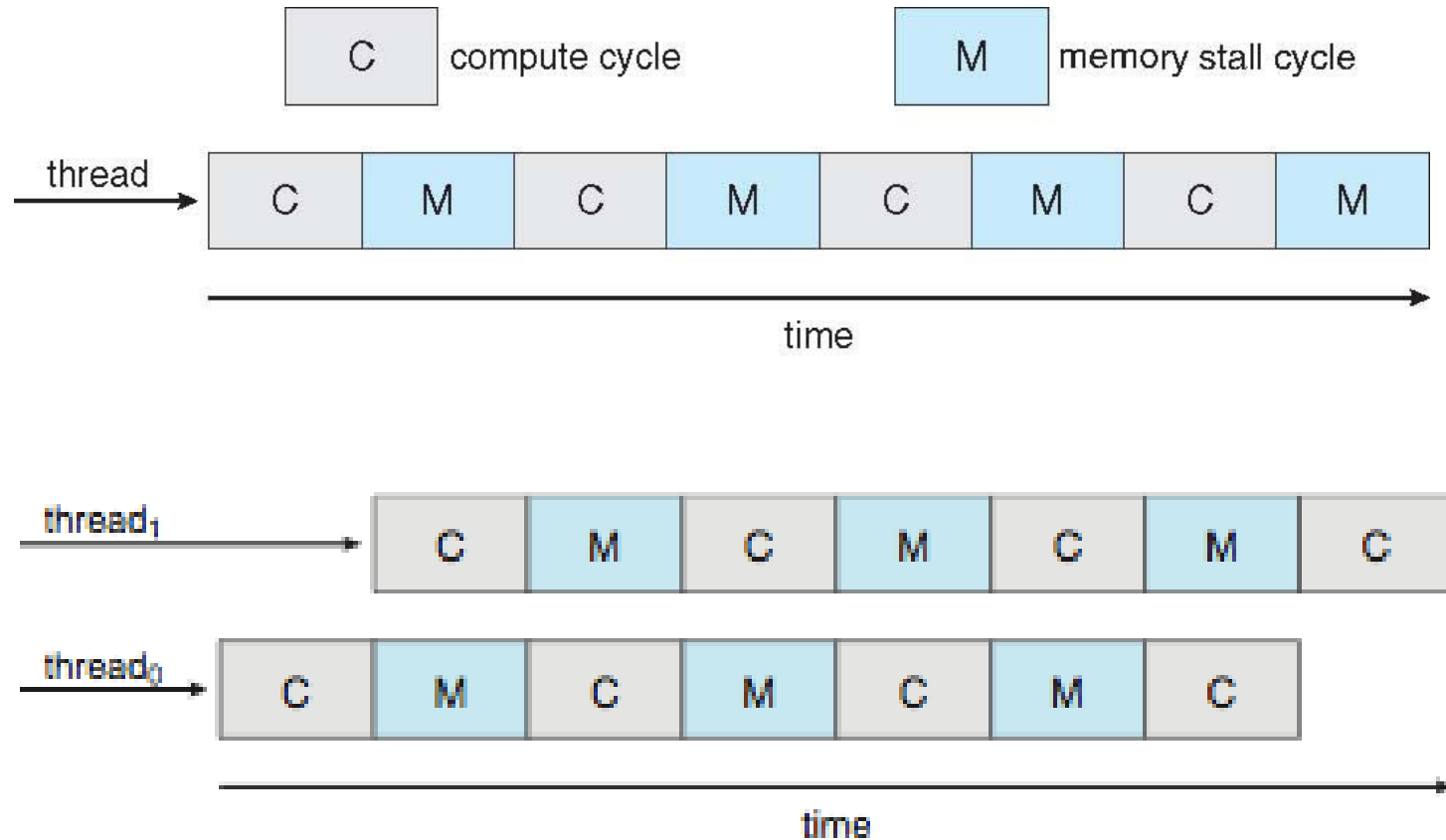
- If SMP, need to keep all CPUs loaded for efficiency
- **Load balancing** attempts to keep workload evenly distributed
- **Push migration** – periodic task checks load on each processor, and if found pushes task from overloaded CPU to other CPUs
- **Pull migration** – idle processors pulls waiting task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

# Multicore Processors

---

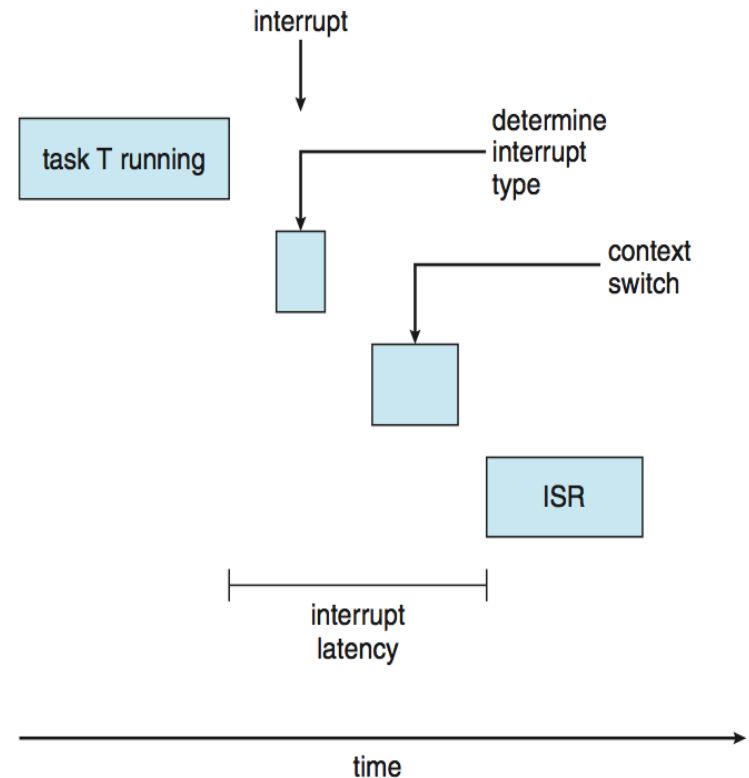
- Recent trend to place multiple processor cores on same physical chip
- Faster and consumes less power
- Multiple threads per core also growing
  - Takes advantage of memory stall to make progress on another thread while memory retrieve happens

# Multithreaded Multicore System



# Real-Time CPU Scheduling

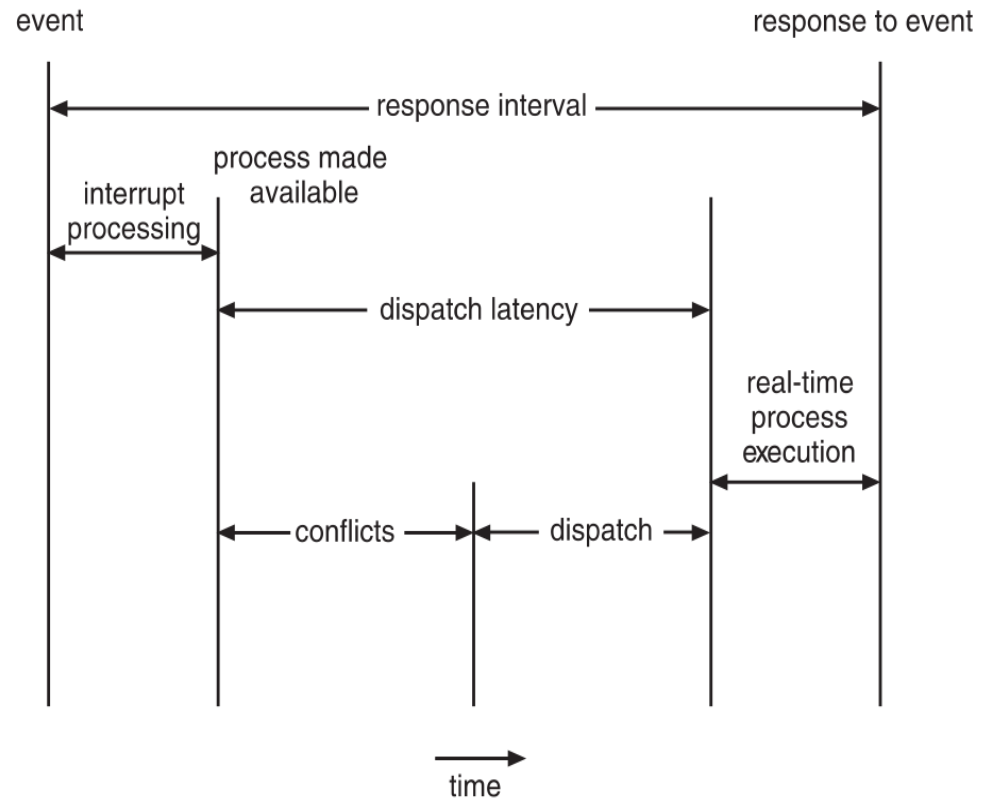
- Can present obvious challenges
- **Soft real-time systems** – no guarantee as to when critical real-time process will be scheduled
- **Hard real-time systems** – task must be serviced by its deadline
- Two types of latencies affect performance
  1. Interrupt latency – time from arrival of interrupt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





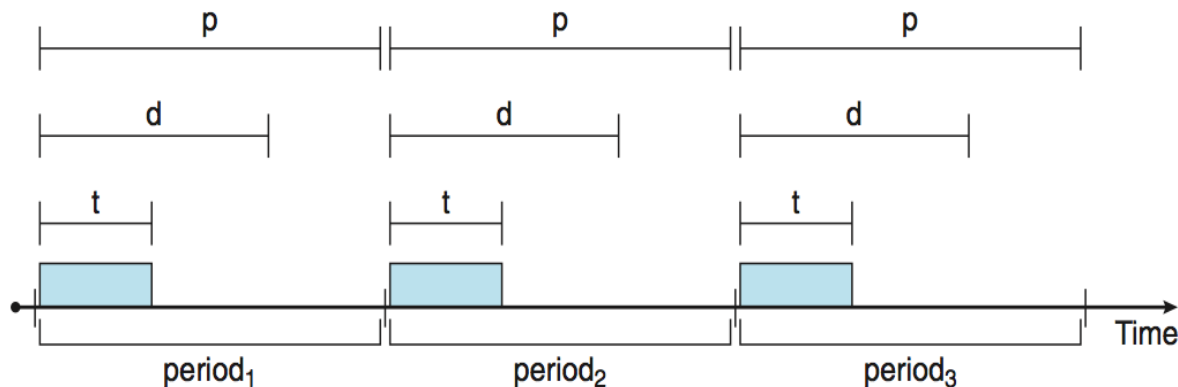
# Real-Time CPU Scheduling (Cont.)

- 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



# Priority-based Scheduling

- For real-time scheduling, scheduler must support preemptive, priority-based scheduling
  - But only guarantees soft real-time
- For hard real-time must also provide ability to meet deadlines
- Processes have new characteristics: **periodic** ones require CPU at constant intervals
  - Has processing time  $t$ , deadline  $d$ , period  $p$
  - $0 \leq t \leq d \leq p$
  - **Rate** of periodic task is  $1/p$



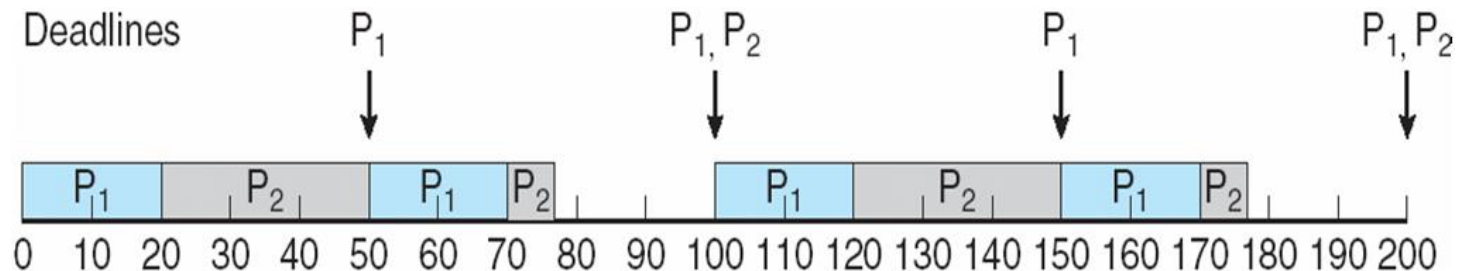
# Virtualization and Scheduling

---

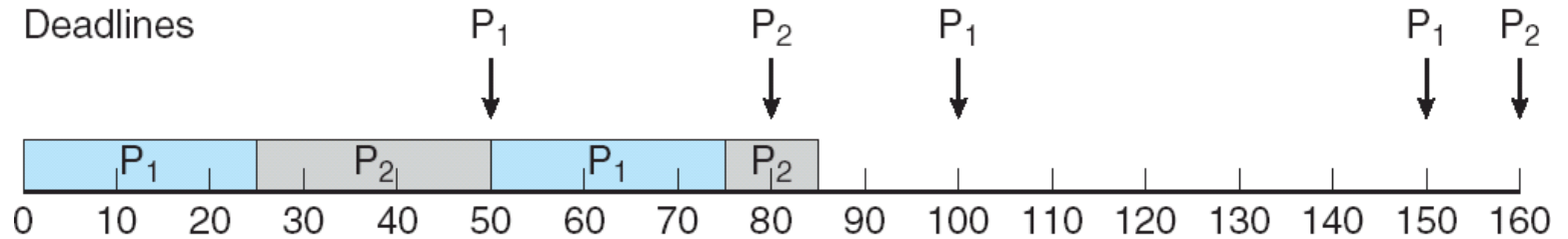
- Virtualization software schedules multiple guests onto CPU(s)
- Each guest doing its own scheduling
  - Not knowing it doesn't own the CPUs
  - Can result in poor response time
  - Can effect time-of-day clocks in guests
- Can undo good scheduling algorithm efforts of guests

# Rate Monotonic Scheduling

- A priority is assigned based on the inverse of its period
- Shorter periods = higher priority;
- Longer periods = lower priority
- $P_1$  is assigned a higher priority than  $P_2$ .



# Missed Deadlines with Rate Monotonic Scheduling

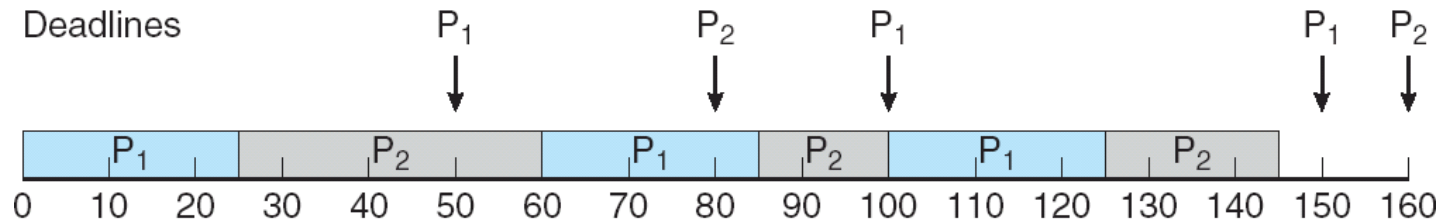


# Earliest Deadline First Scheduling (EDF)

- Priorities are assigned according to deadlines:

the earlier the deadline, the higher the priority;

the later the deadline, the lower the priority



# Proportional Share Scheduling

---

- $T$  shares are allocated among all processes in the system
- An application receives  $N$  shares where  $N < T$
- This ensures each application will receive  $N / T$  of the total processor time

# POSIX Real-Time Scheduling

---

- The POSIX.1b standard
- API provides functions for managing real-time threads
- Defines two scheduling classes for real-time threads:
  1. SCHED\_FIFO - threads are scheduled using a FCFS strategy with a FIFO queue. There is no time-slicing for threads of equal priority
  2. SCHED\_RR - similar to SCHED\_FIFO except time-slicing occurs for threads of equal priority
- Defines two functions for getting and setting scheduling policy:
  - n `pthread_attr_getsched_policy(pthread_attr_t *attr, int *policy)`
  - n `pthread_attr_setsched_policy(pthread_attr_t *attr, int policy)`



# POSIX Real-Time Scheduling API

---

```
#include <pthread.h>
#include <stdio.h>
#define NUM_THREADS 5
int main(int argc, char *argv[])
{
    int i, policy;
    pthread_t_tid[NUM_THREADS];
    pthread_attr_t attr;
    /* get the default attributes */
    pthread_attr_init(&attr);
    /* get the current scheduling policy */
    if (pthread_attr_getschedpolicy(&attr, &policy) != 0)
        fprintf(stderr, "Unable to get policy.\n");
    else {
        if (policy == SCHED_OTHER) printf("SCHED_OTHER\n");
        else if (policy == SCHED_RR) printf("SCHED_RR\n");
        else if (policy == SCHED_FIFO) printf("SCHED_FIFO\n");
    }
}
```

# POSIX Real-Time Scheduling API (Cont.)

---

```
/* set the scheduling policy - FIFO, RR, or OTHER */
if (pthread_attr_setschedpolicy(&attr, SCHED_FIFO) != 0)
    fprintf(stderr, "Unable to set policy.\n");
/* create the threads */
for (i = 0; i < NUM_THREADS; i++)
    pthread_create(&tid[i], &attr, runner, NULL);
/* now join on each thread */
for (i = 0; i < NUM_THREADS; i++)
    pthread_join(tid[i], NULL);
}

/* Each thread will begin control in this function */
void *runner(void *param)
{
    /* do some work ... */
    pthread_exit(0);
}
```

# Operating System Examples

---

- Linux scheduling
- Windows scheduling
- Solaris scheduling

# Linux Scheduling Through Version 2.5

---

- Prior to kernel version 2.5, ran variation of standard UNIX scheduling algorithm
- Version 2.5 moved to constant order  $O(1)$  scheduling time
  - Preemptive, priority based
  - Two priority ranges: time-sharing and real-time
  - **Real-time** range from 0 to 99 and **nice** value from 100 to 140
  - Map into global priority with numerically lower values indicating higher priority
  - Higher priority gets larger  $q$
  - Task run-able as long as time left in time slice (**active**)
  - If no time left (**expired**), not run-able until all other tasks use their slices
  - All run-able tasks tracked in per-CPU **runqueue** data structure
    - Two priority arrays (active, expired)
    - Tasks indexed by priority
    - When no more active, arrays are exchanged
  - Worked well, but poor response times for interactive processes

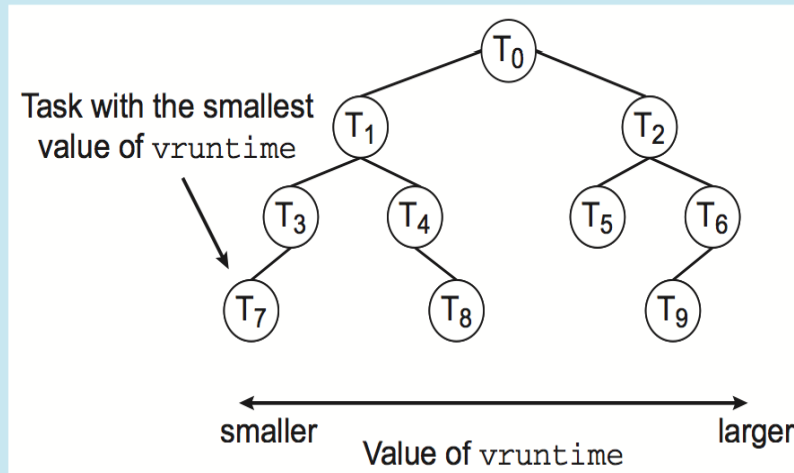
# Linux Scheduling in Version 2.6.23 +

---

- **Completely Fair Scheduler** (CFS)
- **Scheduling classes**
  - Each has specific priority
  - Scheduler picks highest priority task in highest scheduling class
  - Rather than quantum based on fixed time allotments, based on proportion of CPU time
  - 2 scheduling classes included, others can be added
    1. default
    2. real-time
- Quantum calculated based on **nice value** from -20 to +19
  - Lower value is higher priority
  - Calculates **target latency** – interval of time during which task should run at least once
  - Target latency can increase if say number of active tasks increases
- CFS scheduler maintains per task **virtual run time** in variable **vruntime**
  - Associated with decay factor based on priority of task – lower priority is higher decay rate
  - Normal default priority yields virtual run time = actual run time
- To decide next task to run, scheduler picks task with lowest virtual run time

# CFS Performance

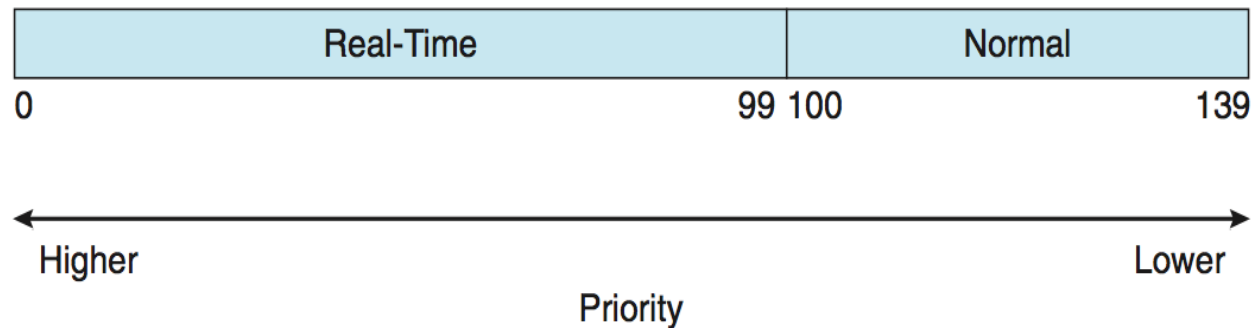
The Linux CFS scheduler provides an efficient algorithm for selecting which task to run next. Each runnable task is placed in a red-black tree—a balanced binary search tree whose key is based on the value of `vruntime`. This tree is shown below:



When a task becomes runnable, it is added to the tree. If a task on the tree is not runnable (for example, if it is blocked while waiting for I/O), it is removed. Generally speaking, tasks that have been given less processing time (smaller values of `vruntime`) are toward the left side of the tree, and tasks that have been given more processing time are on the right side. According to the properties of a binary search tree, the leftmost node has the smallest key value, which for the sake of the CFS scheduler means that it is the task with the highest priority. Because the red-black tree is balanced, navigating it to discover the leftmost node will require  $O(\lg N)$  operations (where  $N$  is the number of nodes in the tree). However, for efficiency reasons, the Linux scheduler caches this value in the variable `rb_leftmost`, and thus determining which task to run next requires only retrieving the cached value.

# Linux Scheduling (Cont.)

- Real-time scheduling according to POSIX.1b
  - Real-time tasks have static priorities
- Real-time plus normal map into global priority scheme
- Nice value of -20 maps to global priority 100
- Nice value of +19 maps to priority 139



# Windows Scheduling

---

- ❑ Windows uses priority-based preemptive scheduling
- ❑ Highest-priority thread runs next
- ❑ **Dispatcher** is scheduler
- ❑ Thread runs until (1) blocks, (2) uses time slice, (3) preempted by higher-priority thread
- ❑ Real-time threads can preempt non-real-time
- ❑ 32-level priority scheme
- ❑ **Variable class** is 1-15, **real-time class** is 16-31
- ❑ Priority 0 is memory-management thread
- ❑ Queue for each priority
- ❑ If no run-able thread, runs **idle thread**



# Windows Priority Classes

---

- Win32 API identifies several priority classes to which a process can belong
  - REALTIME\_PRIORITY\_CLASS, HIGH\_PRIORITY\_CLASS, ABOVE\_NORMAL\_PRIORITY\_CLASS, NORMAL\_PRIORITY\_CLASS, BELOW\_NORMAL\_PRIORITY\_CLASS, IDLE\_PRIORITY\_CLASS
  - All are variable except REALTIME
- A thread within a given priority class has a relative priority
  - TIME\_CRITICAL, HIGHEST, ABOVE\_NORMAL, NORMAL, BELOW\_NORMAL, LOWEST, IDLE
- Priority class and relative priority combine to give numeric priority
- Base priority is NORMAL within the class
- If quantum expires, priority lowered, but never below base

# Windows Priority Classes (Cont.)

---

- If wait occurs, priority boosted depending on what was waited for
- Foreground window given 3x priority boost
- Windows 7 added **user-mode scheduling (UMS)**
  - Applications create and manage threads independent of kernel
  - For large number of threads, much more efficient
  - UMS schedulers come from programming language libraries like C++ **Concurrent Runtime** (ConcRT) framework

# Windows Priorities

---

|               | real-time | high | above normal | normal | below normal | idle priority |
|---------------|-----------|------|--------------|--------|--------------|---------------|
| time-critical | 31        | 15   | 15           | 15     | 15           | 15            |
| highest       | 26        | 15   | 12           | 10     | 8            | 6             |
| above normal  | 25        | 14   | 11           | 9      | 7            | 5             |
| normal        | 24        | 13   | 10           | 8      | 6            | 4             |
| below normal  | 23        | 12   | 9            | 7      | 5            | 3             |
| lowest        | 22        | 11   | 8            | 6      | 4            | 2             |
| idle          | 16        | 1    | 1            | 1      | 1            | 1             |

# Solaris

---

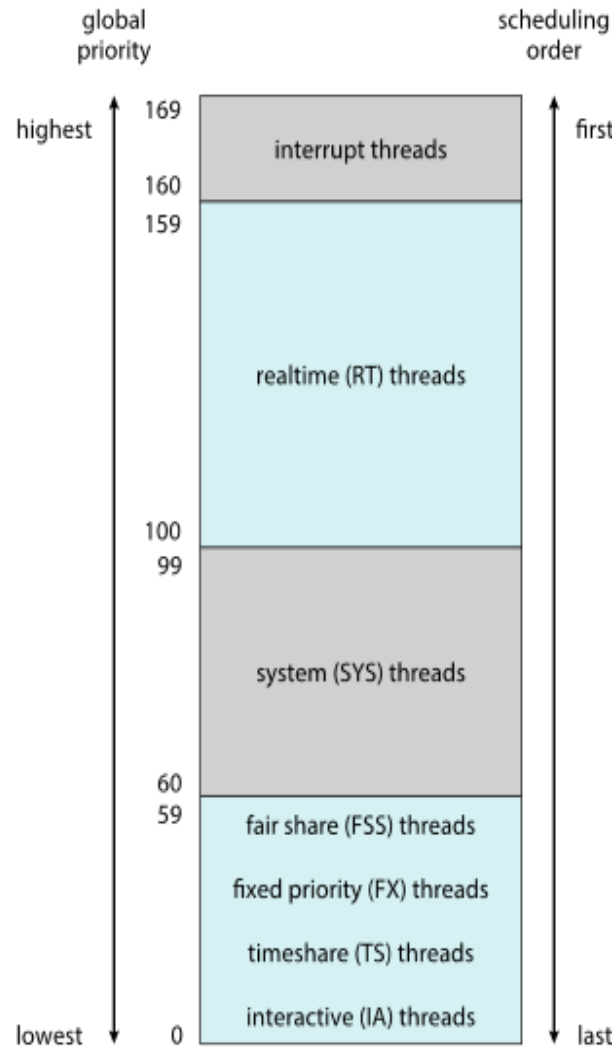
- Priority-based scheduling
- Six classes available
  - Time sharing (default) (TS)
  - Interactive (IA)
  - Real time (RT)
  - System (SYS)
  - Fair Share (FSS)
  - Fixed priority (FP)
- Given thread can be in one class at a time
- Each class has its own scheduling algorithm
- Time sharing is multi-level feedback queue
  - Loadable table configurable by sysadmin

# Solaris Dispatch Table

---

| priority | time quantum | time quantum expired | return from sleep |
|----------|--------------|----------------------|-------------------|
| 0        | 200          | 0                    | 50                |
| 5        | 200          | 0                    | 50                |
| 10       | 160          | 0                    | 51                |
| 15       | 160          | 5                    | 51                |
| 20       | 120          | 10                   | 52                |
| 25       | 120          | 15                   | 52                |
| 30       | 80           | 20                   | 53                |
| 35       | 80           | 25                   | 54                |
| 40       | 40           | 30                   | 55                |
| 45       | 40           | 35                   | 56                |
| 50       | 40           | 40                   | 58                |
| 55       | 40           | 45                   | 58                |
| 59       | 20           | 49                   | 59                |

# Solaris Scheduling



# Solaris Scheduling (Cont.)

---

- Scheduler converts class-specific priorities into a per-thread global priority
  - Thread with highest priority runs next
  - Runs until (1) blocks, (2) uses time slice, (3) preempted by higher-priority thread
  - Multiple threads at same priority selected via RR

# Algorithm Evaluation

---

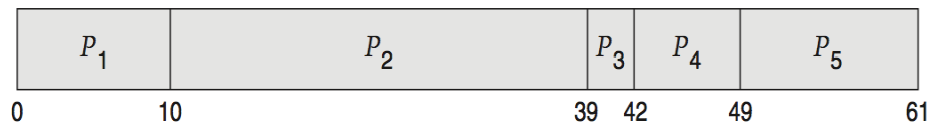
- How to select CPU-scheduling algorithm for an OS?
- Determine criteria, then evaluate algorithms
- **Deterministic modeling**
  - Type of **analytic evaluation**
  - Takes a particular predetermined workload and defines the performance of each algorithm for that workload
- Consider 5 processes arriving at time 0:

| <u>Process</u> | <u>Burst Time</u> |
|----------------|-------------------|
| $P_1$          | 10                |
| $P_2$          | 29                |
| $P_3$          | 3                 |
| $P_4$          | 7                 |
| $P_5$          | 12                |

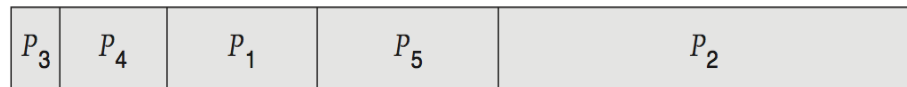


# Deterministic Evaluation

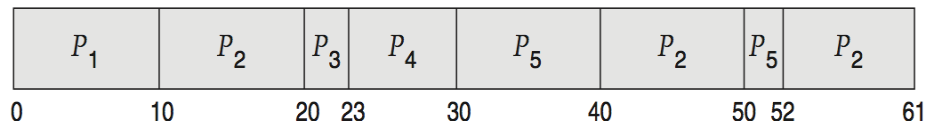
- For each algorithm, calculate minimum average waiting time
- Simple and fast, but requires exact numbers for input, applies only to those inputs
  - FCS is 28ms:



- Not preemptive scheduling



- RR



# Queueing Models

---

- Describes the arrival of processes, and CPU and I/O bursts probabilistically
  - Commonly exponential, and described by mean
  - Computes average throughput, utilization, waiting time, etc
- Computer system described as network of servers, each with queue of waiting processes
  - Knowing arrival rates and service rates
  - Computes utilization, average queue length, average wait time, etc

# Little's Formula

---

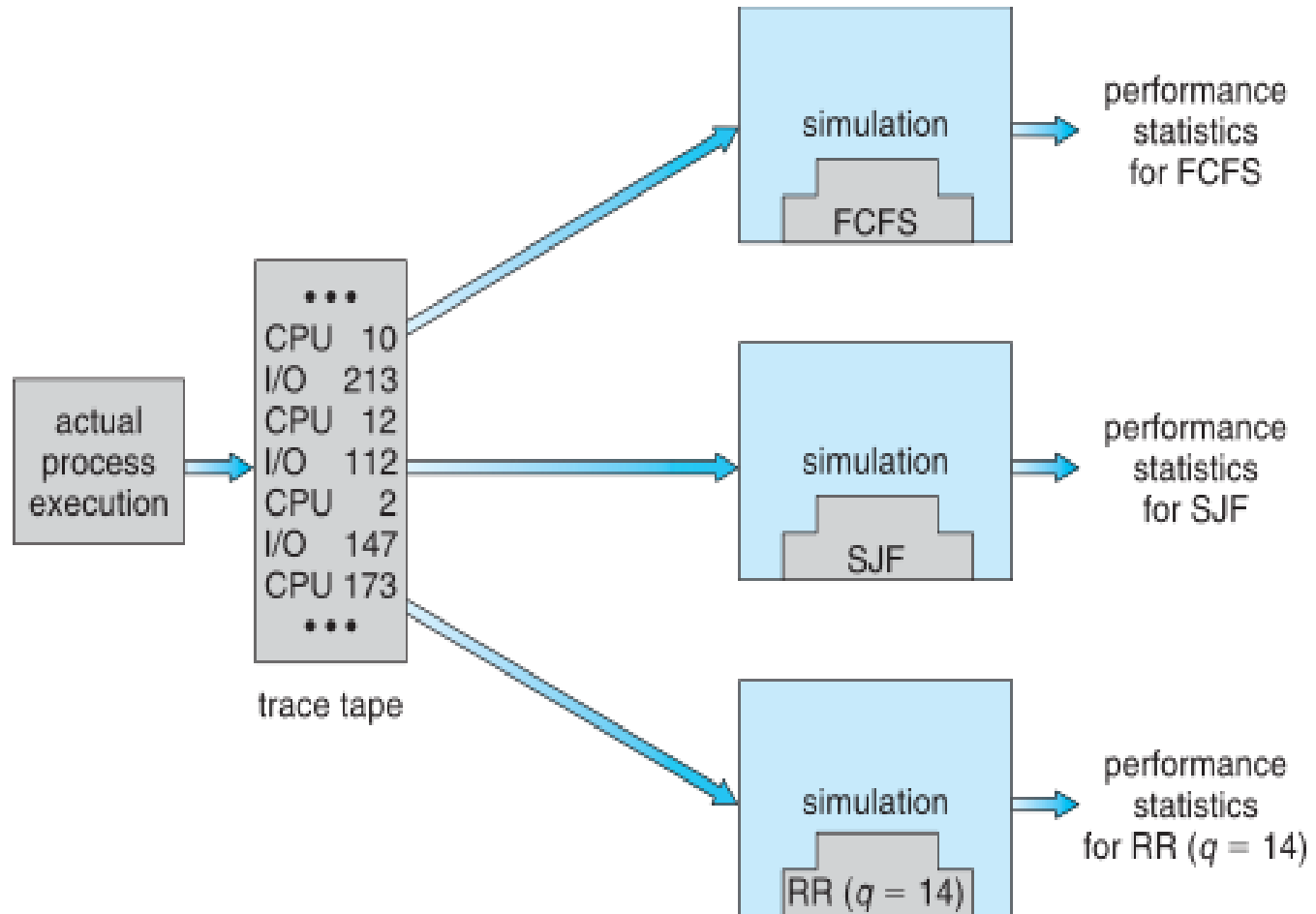
- $n$  = average queue length
- $W$  = average waiting time in queue
- $\lambda$  = average arrival rate into queue
- Little's law – in steady state, processes leaving queue must equal processes arriving, thus:  
$$n = \lambda \times W$$
  - Valid for any scheduling algorithm and arrival distribution
- For example, if on average 7 processes arrive per second, and normally 14 processes in queue, then average wait time per process = 2 seconds

# Simulations

---

- Queueing models limited
- **Simulations** more accurate
  - Programmed model of computer system
  - Clock is a variable
  - Gather statistics indicating algorithm performance
  - Data to drive simulation gathered via
    - Random number generator according to probabilities
    - Distributions defined mathematically or empirically
    - Trace tapes record sequences of real events in real systems

# Evaluation of CPU Schedulers by Simulation



# Implementation

---

- Even simulations have limited accuracy
- Just implement new scheduler and test in real systems
  - High cost, high risk
  - Environments vary
- Most flexible schedulers can be modified per-site or per-system
- Or APIs to modify priorities
- But again environments vary

# End of Chapter 6

