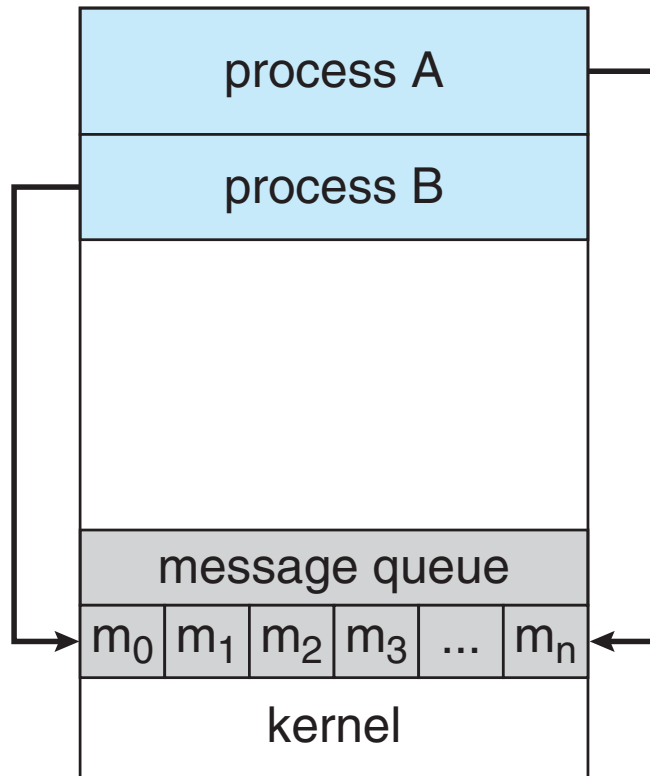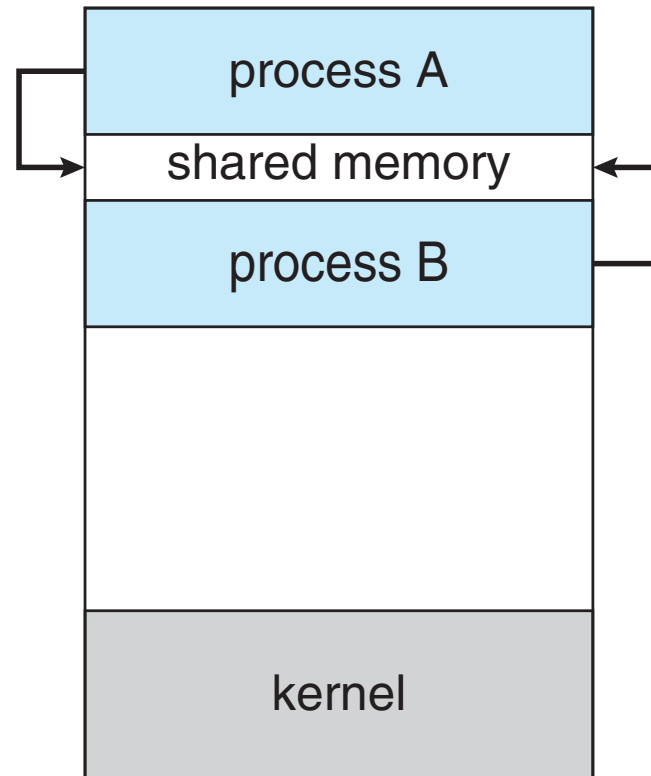# COMP 3511
# Operating Systems

Lab #4

# Interprocess Communication - IPC

- 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)                              (b)

# 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

# Bounded-Buffer – Shared-Memory Solution

- Shared data

```
#define BUFFER_SIZE 10
typedef struct {
  . . .
} item;

item buffer[BUFFER_SIZE];
int in = 0;
int out = 0;
```

- Solution is correct, but can only use BUFFER_SIZE-1 elements

# Bounded-Buffer – Producer

```
item next produced;
while (true) {
  /* produce an item in next produced
   */
  while (((in + 1) % BUFFER SIZE) ==
    out)
      ; /* do nothing */
  buffer[in] = next produced;
  in = (in + 1) % BUFFER SIZE;
}
```

# Bounded Buffer – Consumer
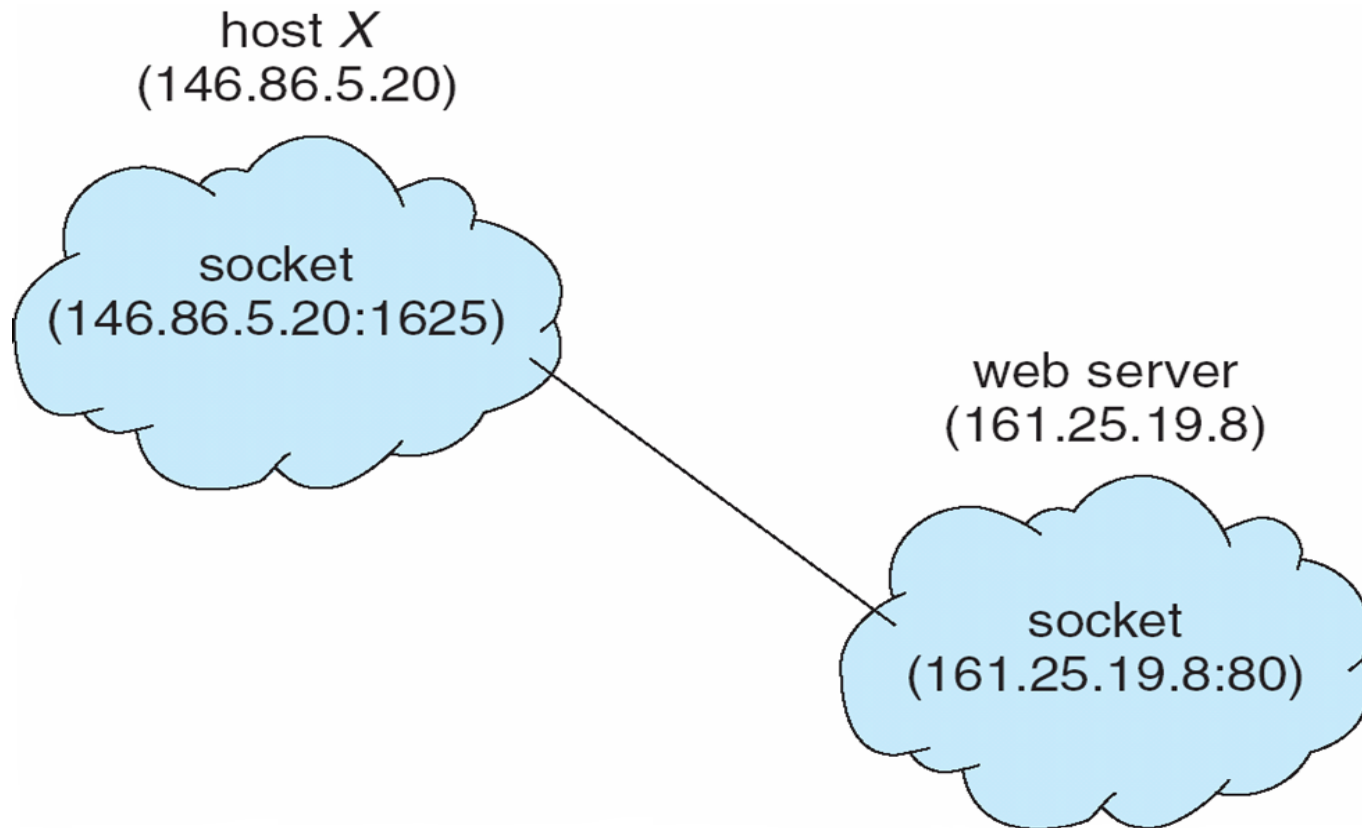
```
item next consumed;
while (true) {
    while (in == out)
        ; /* do nothing */
    next consumed = buffer[out];
    out = (out + 1) % BUFFER SIZE;

    /* consume the item in next
 consumed */
}
```

# Sockets

- A **socket** is defined as an endpoint for communication

- Concatenation of IP address and **port** – a number included at start of message packet to differentiate network services on a host

- The socket **161.25.19.8:1625** refers to port **1625** on host **161.25.19.8**

- Communication consists between a pair of sockets

- All ports below 1024 are ***well known***, used for standard services

- Special IP address 127.0.0.1 (**loopback**) to refer to system on which process is running

# Socket Communication

# Sockets in Java

- Three types of sockets
  - **Connection-oriented** (**TCP**)
  - **Connectionless** (**UDP**)
  - `MulticastSocket` class– data can be sent to multiple recipients

- Consider this "Date" server:

```java
import java.net.*;
import java.io.*;

public class DateServer
{
  public static void main(String[] args) {
    try {
      ServerSocket sock = new ServerSocket(6013);

      /* now listen for connections */
      while (true) {
        Socket client = sock.accept();

        PrintWriter pout = new
          PrintWriter(client.getOutputStream(), true);

        /* write the Date to the socket */
        pout.println(new java.util.Date().toString());

        /* close the socket and resume */
        /* listening for connections */
        client.close();
      }
    }
    catch (IOException ioe) {
      System.err.println(ioe);
    }
  }
}
```
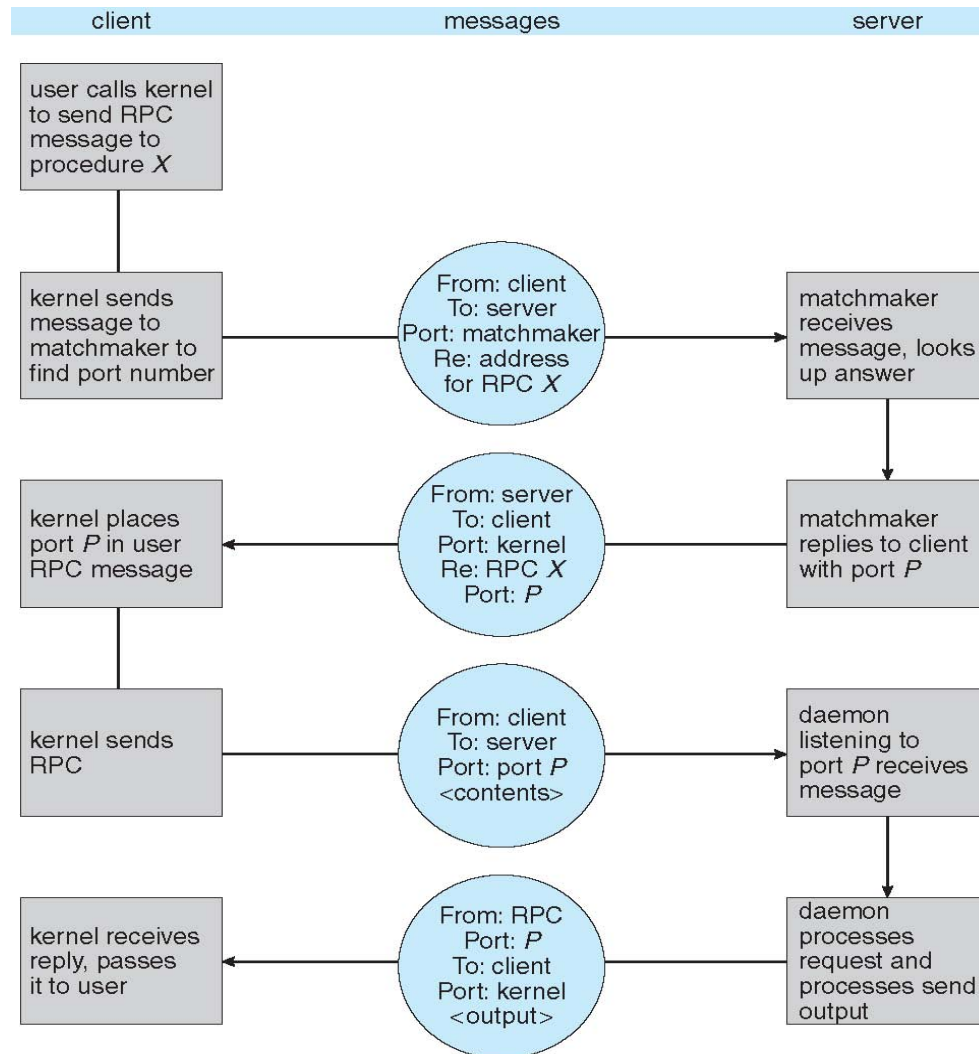
# Remote Procedure Calls

- Remote procedure call (RPC) abstracts procedure calls between processes on networked systems
  - Again uses ports for service differentiation
- **Stubs** – client-side proxy for the actual procedure on the server
- The client-side stub locates the server and **marshalls** the parameters
- The server-side stub receives this message, unpacks the marshalled parameters, and performs the procedure on the server
- On Windows, stub code compile from specification written in **Microsoft Interface Definition Language** (**MIDL**)
- Data representation handled via **External Data Representation** (**XDL**) format to account for different architectures
  - **Big-endian** and **little-endian**
- Remote communication has more failure scenarios than local
  - Messages can be delivered *exactly once* rather than *at most once*
- OS typically provides a rendezvous (or **matchmaker**) service to connect client and server

# Execution of RPC



| client | messages | server |
|---|---|---|
| user calls kernel to send RPC message to procedure X | | |
| kernel sends message to matchmaker to find port number | From: client<br>To: server<br>Port: matchmaker<br>Re: address<br>for RPC X | matchmaker receives message, looks up answer |
| kernel places port P in user RPC message | From: server<br>To: client<br>Port: kernel<br>Re: RPC X<br>Port: P | matchmaker replies to client with port P |
| kernel sends RPC | From: client<br>To: server<br>Port: port P<br><contents> | daemon listening to port P receives message |
| kernel receives reply, passes it to user | From: RPC<br>Port: P<br>To: client<br>Port: kernel<br><output> | daemon processes request and processes send output |

# Thread Libraries

- **Thread library** provides programmer with API for creating and managing threads

- Two primary ways of implementing
  - Library entirely in user space
  - Kernel-level library supported by the OS

# Pthreads

- May be provided either as user-level or kernel-level
- A POSIX standard (IEEE 1003.1c) API for thread creation and synchronization
- ***Specification***, not ***implementation***
- API specifies behavior of the thread library, implementation is up to development of the library
- Common in UNIX operating systems (Solaris, Linux, Mac OS X)

# Pthreads Example

```c
#include <pthread.h>
#include <stdio.h>

int sum; /* this data is shared by the thread(s) */
void *runner(void *param); /* threads call this function */

int main(int argc, char *argv[])
{
    pthread_t tid; /* the thread identifier */
    pthread_attr_t attr; /* set of thread attributes */

    if (argc != 2) {
        fprintf(stderr,"usage: a.out <integer value>\n");
        return -1;
    }
    if (atoi(argv[1]) < 0) {
        fprintf(stderr,"%d must be >= 0\n",atoi(argv[1]));
        return -1;
    }
```

# Pthreads Example (Cont.)

```
        /* get the default attributes */
        pthread_attr_init(&attr);
        /* create the thread */
        pthread_create(&tid,&attr,runner,argv[1]);
        /* wait for the thread to exit */
        pthread_join(tid,NULL);

        printf("sum = %d\n",sum);
}

/* The thread will begin control in this function */
void *runner(void *param)
{
        int i, upper = atoi(param);
        sum = 0;

        for (i = 1; i <= upper; i++)
                sum += i;

        pthread_exit(0);
}
```

**Figure 4.9** Multithreaded C program using the Pthreads API.

# Pthreads Code for Joining 10 Threads

```
#define NUM_THREADS 10

/* an array of threads to be joined upon */
pthread_t workers[NUM_THREADS];

for (int i = 0; i < NUM_THREADS; i++)
    pthread_join(workers[i], NULL);
```

**Figure 4.10**  Pthread code for joining ten threads.

# Win32 API  Multithreaded C Program

```c
#include <windows.h>
#include <stdio.h>
DWORD Sum; /* data is shared by the thread(s) */

/* the thread runs in this separate function */
DWORD WINAPI Summation(LPVOID Param)
{
   DWORD Upper = *(DWORD*)Param;
   for (DWORD i = 0; i <= Upper; i++)
      Sum += i;
   return 0;
}

int main(int argc, char *argv[])
{
   DWORD ThreadId;
   HANDLE ThreadHandle;
   int Param;

   if (argc != 2) {
      fprintf(stderr,"An integer parameter is required\n");
      return -1;
   }
   Param = atoi(argv[1]);
   if (Param < 0) {
      fprintf(stderr,"An integer >= 0 is required\n");
      return -1;
   }
```

# Win32 API  Multithreaded C Program (Cont.)

```c
/* create the thread */
ThreadHandle = CreateThread(
    NULL, /* default security attributes */
    0, /* default stack size */
    Summation, /* thread function */
    &Param, /* parameter to thread function */
    0, /* default creation flags */
    &ThreadId); /* returns the thread identifier */

if (ThreadHandle != NULL) {
    /* now wait for the thread to finish */
    WaitForSingleObject(ThreadHandle,INFINITE);

    /* close the thread handle */
    CloseHandle(ThreadHandle);

    printf("sum = %d\n",Sum);
}
}
```
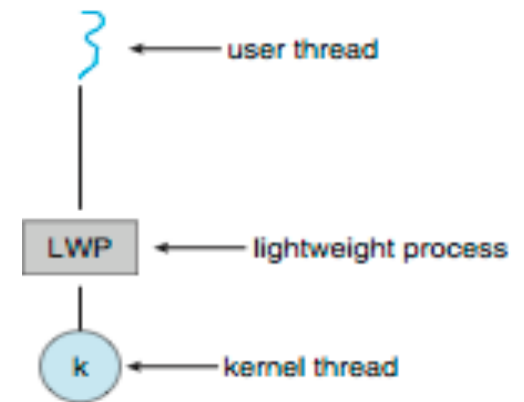
# Scheduler Activations

- Both M:M and Two-level models require communication to maintain the appropriate number of kernel threads allocated to the application
- Typically use an intermediate data structure between user and kernel threads – **lightweight process** (**LWP**)
  - Appears to be a virtual processor on which process can schedule user thread to run
  - Each LWP attached to kernel thread
  - How many LWPs to create?
- Scheduler activations provide **upcalls** - a communication mechanism from the kernel to the **upcall handler** in the thread library
- This communication allows an application to maintain the correct number kernel threads

# Operating System Examples

- Linux scheduling

- Windows scheduling

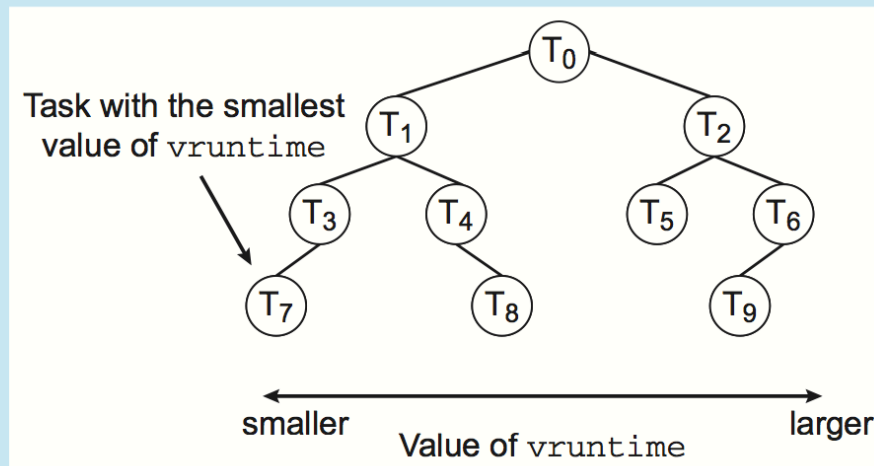- Solaris scheduling

# 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
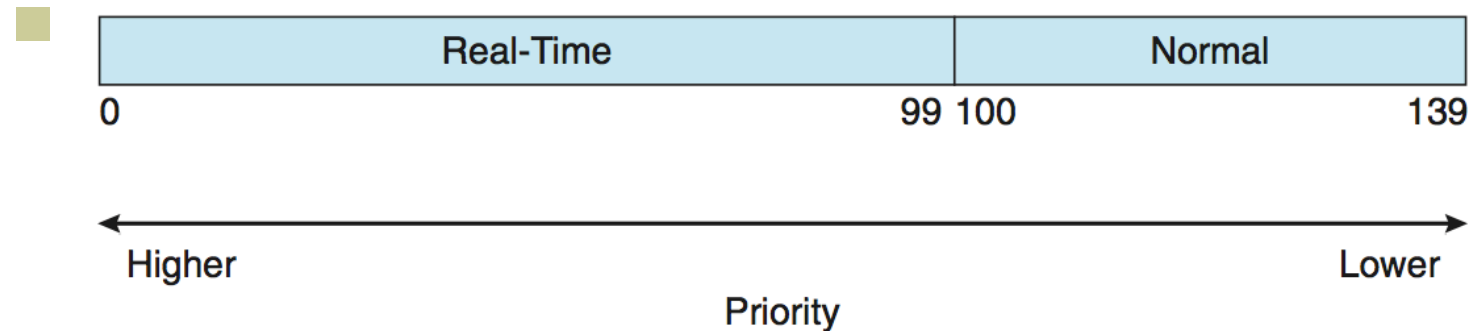
# CFS Performance

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

| Real-Time | Normal |
|---|---|
| 0                          99 | 100                        139 |

Higher ◄─────────────────────────────► Lower

Priority

# 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
- 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

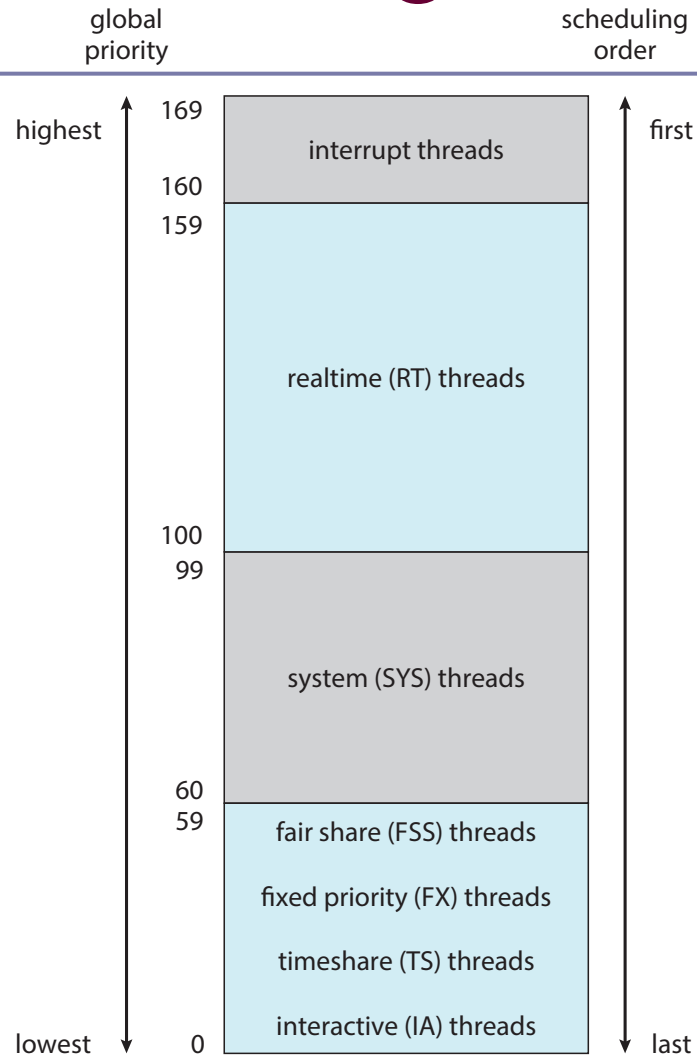|  | real-time | high | above normal | normal | below normal | idle priority |
|---|---|---|---|---|---|---|
| time-critical | 31 | 15 | 15 | 15 | 15 | 15 |
| highest | 26 | 15 | 12 | 10 | 8 | 6 |
| above normal | 25 | 14 | 11 | 9 | 7 | 5 |
| normal | 24 | 13 | 10 | 8 | 6 | 4 |
| below normal | 23 |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

# CPU Scheduler & Dispatcher

- Selects from among the processes in memory that are ready to execute, and allocates the CPU to one of them
- 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

- Dispatcher module gives control of the CPU to the process selected by the CPU scheduler; this involves
    - switching context
    - switching to user mode
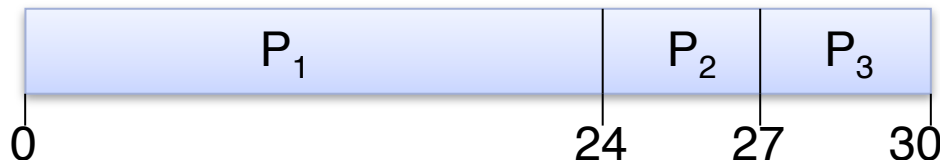    - jumping to the proper location in the user program to restart that program

# Scheduling Criteria

- To Maximize:
  - **CPU utilization** – keep the CPU as busy as possible
  - **Throughput** – # of processes that complete their execution per time unit
- To Minimize:
  - **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)
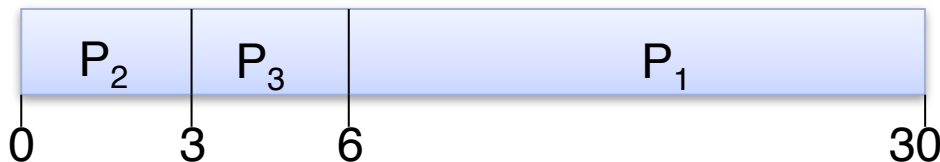
# First-Come, First-Served (FCFS) Scheduling

- First-Come-First-Serve (FCFS)
- Pro: easy to implement
- Con: potentially bad for short jobs
-

| Process | Burst Time |
|---------|------------|
| $P_1$   | 24         |
| $P_2$   | 3          |
| $P_3$   | 3          |

| $P_1$ | $P_2$ | $P_3$ |
|-------|-------|-------|

0　　　　　　　　　　24　　27　　30

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

Average Turnaround time:
(24+27+30)/3 = 27

- If arrive in the order $P_2$ , $P_3$ , $P_1$

| $P_2$ | $P_3$ | $P_1$ |
|-------|-------|-------|

0　　3　　6　　　　　　　　30

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

Average Turnaround time:
(30+3+6)/3 = 13
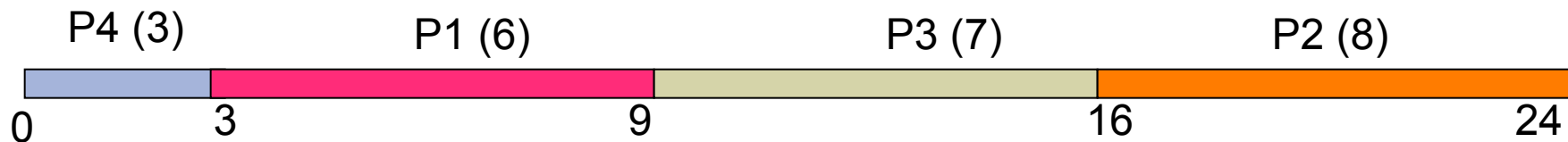
# Shortest-Job-First (SJF) Scheduling

- Associate with each process the length of its next CPU burst. Use these lengths to schedule the process with the shortest time

- Two schemes:
  - **Non-preemptive** – once CPU given to the process it cannot be preempted until completes its current CPU burst
  - **Preemptive** – if a new process arrives with CPU burst length less than remaining time of current executing process, preempt. This scheme is know as the **Shortest-Remaining-Time-First (SRTF)**

| | Long | Short | |
|---|---|---|---|

| | Short | Long | |
|---|---|---|---|

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

# Non-preemptive SJF: Example

| Process | Burst Time | Arrival Time |
|---------|------------|--------------|
| P1 | 6 | 0 |
| P2 | 8 | 0 |
| P3 | 7 | 0 |
| P4 | 3 | 0 |

| P4 (3) | P1 (6) | P3 (7) | P2 (8) |
|--------|--------|--------|--------|

0    3                9                    16              24

P4 waiting time: 0
P1 waiting time: 3
P3 waiting time: 9
P2 waiting time: 16

The average waiting time
$(0+3+9+16)/4 = 7$

# Comparing to FCFS

| Process | Burst Time | Arrival Time |
|---------|------------|--------------|
| P1 | 6 | 0 |
| P2 | 8 | 0 |
| P3 | 7 | 0 |
| P4 | 3 | 0 |

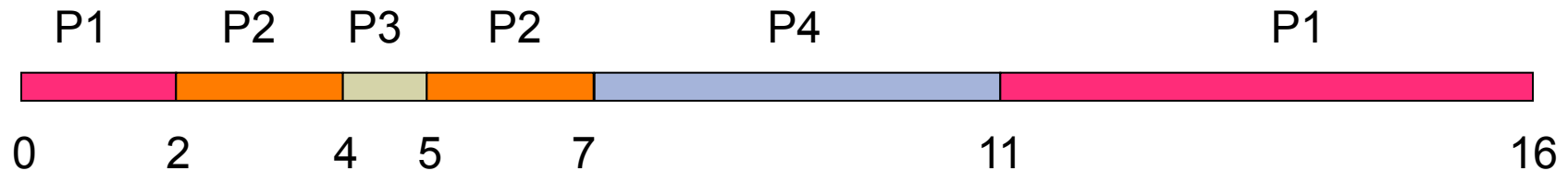P1 (6)    P2 (8)    P3 (7)    P4 (3)

0    6    14    21    24

P1 waiting time: 0
P2 waiting time: 6
P3 waiting time: 14
P4 waiting time: 21

Assume execution order is
P1, P2, P3, P4

The average waiting time
(0+6+14+21)/4 = 10.25 > 7

# Preemptive SJF: Example

| Process | Burst Time | Arrival Time |
|---------|------------|--------------|
| P1 | 7 | 0 |
| P2 | 4 | 2 |
| P3 | 1 | 4 |
| P4 | 4 | 5 |

P1      P2      P3      P2              P4                      P1

```
0       2       4   5       7                   11                      16
```

P1 gets preempted at time 2

P2 gets preempted at time 4

The average waiting time
(9 + 1 + 0 +2)/4 = 3

# Priority Scheduling

- A priority number (integer) is associated with each process
- The CPU is allocated to the process with the highest priority (smallest integer ≡ highest priority)
  - Preemptive
  - Nonpreemptive

- SJF is a priority scheduling where priority is the predicted next CPU burst time
- Problem ≡ **Starvation** – low priority processes may never execute
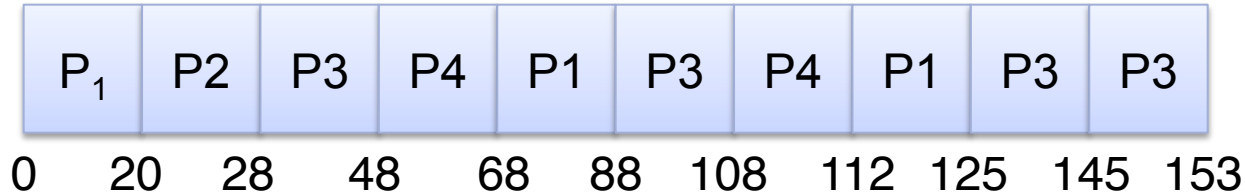- Solution ≡ **Aging** – as time progresses increase the priority of the process

# Round Robin (RR) Scheduling

- Each process gets a small unit of CPU time (*time quantum*), usually 10-100 milliseconds.  After quantum expires, the process is preempted (OS timer interrupt) and added to the end of the ready queue.

- n processes in the ready queue and the time quantum is q
  - 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

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

# Example of RR in Notes

- Example:

| Process | Burst Time |
|---------|------------|
| P1 | 53 |
| P2 | 8 |
| P3 | 68 |
| P4 | 24 |

- The Gantt chart is:    turnaround time = completion time-arrival time

| P$_1$ | P2 | P3 | P4 | P1 | P3 | P4 | P1 | P3 | P3 |
|------|----|----|----|----|----|----|----|----|----|

```
0    20   28    48    68   88   108   112  125   145  153
```

- Waiting time for

  P1=(68-20)+(112-88)=72                          P2=(20-0)=20
  
  P3=(28-0)+(88-48)+(125-108)=85
  
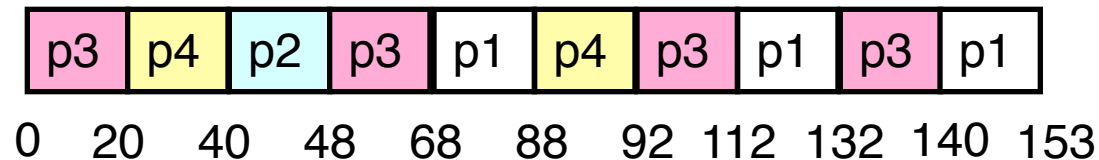  P4=(48-0)+(108-68)=88

- Average waiting time = (72+20+85+88)/4=66¼

- Average completion time = (125+28+153+112)/4 = 104½

# Consider Different Arrival Time

- Example:

| Process | Burst Time | Arrival Time |
|---------|------------|--------------|
| P1 | 53 | 28 |
| P2 | 8 | 16 |
| P3 | 68 | 0 |
| P4 | 24 | 10 |

- The Gantt chart is:

| p3 | p4 | p2 | p3 | p1 | p4 | p3 | p1 | p3 | p1 |
|----|----|----|----|----|----|----|----|----|----|

0   20   40   48   68   88   92   112   132   140   153

- Waiting time for

P1=(68-28)+(112-88)+(140-132)=72
P2=(40-16)=24
P3=(48-20)+(92-68)+(132-112)=72
P4=(20-10)+(88-40)=58

- Average waiting time = (72+24+72+58)/4=56½

- Average completion time = (153+48+140+92)/4 = 108¼

# Turnaround Time Varies With The Time Quantum



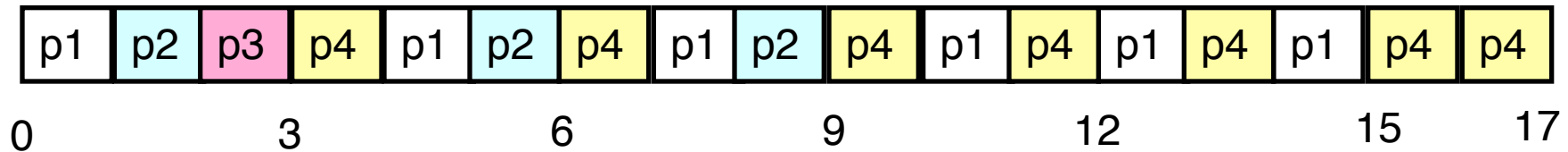| process | time |
|---------|------|
| $P_1$   | 6    |
| $P_2$   | 3    |
| $P_3$   | 1    |
| $P_4$   | 7    |

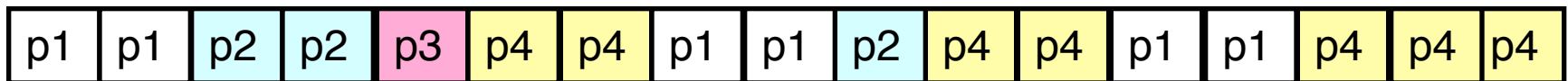**Turnaround time (waiting time + burst)**

**the total amount of time to execute a particular process**

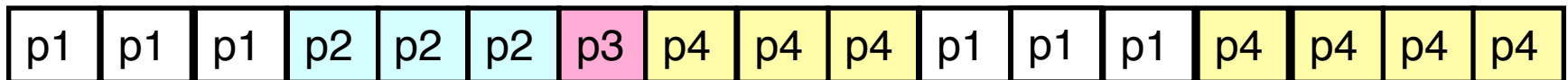# Turnaround Time Varies With The Time Quantum: Example

P1: 6; P2: 3
P3: 1; P4: 7

| p1 | p2 | p3 | p4 | p1 | p2 | p4 | p1 | p2 | p4 | p1 | p4 | p1 | p4 | p1 | p4 | p4 |

0       3       6       9       12       15       17

Q=1      Avg Turnaround Time = (15+9+3+17)/4 = 11

| p1 | p1 | p2 | p2 | p3 | p4 | p4 | p1 | p1 | p2 | p4 | p4 | p1 | p1 | p4 | p4 | p4 |

Q=2      Avg Turnaround Time = (14+10+5+17)/4 = 11.5

| p1 | p1 | p1 | p2 | p2 | p2 | p3 | p4 | p4 | p4 | p1 | p1 | p1 | p4 | p4 | p4 | p4 |

Q=3      Avg Turnaround Time = (13+6+7+17)/4 = 10.75

Q=4   Avg Turnaround Time = (14+7+8+17)/4 = 11.5

| p1 | p1 | p1 | p1 | p2 | p2 | p2 | p3 | p4 | p4 | p4 | p4 | p1 | p1 | p4 | p4 | p4 |

0            3                 6                 9                12              15      17

Q=5    Avg Turnaround Time = (15+8+9+17)/4 = 12.25

| p1 | p1 | p1 | p1 | p1 | p2 | p2 | p2 | p3 | p4 | p4 | p4 | p4 | p4 | p1 | p4 | p4 |

Q=6    Avg Turnaround Time = (6+9+10+17)/4 = 10.5

| p1 | p1 | p1 | p1 | p1 | p1 | p2 | p2 | p2 | p3 | p4 | p4 | p4 | p4 | p4 | p4 | p4 |

Q=7    Avg Turnaround Time = (6+9+10+17)/4 = 10.5

| p1 | p1 | p1 | p1 | p1 | p1 | p2 | p2 | p2 | p3 | p4 | p4 | p4 | p4 | p4 | p4 | p4 |