

ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA Department of computer science and engineering

Real-Time Operating Systems M

3. IPC • Threads • Process Scheduling

Notice

The course material includes slides downloaded from:

http://codex.cs.yale.edu/avi/os-book/

(slides by Silberschatz, Galvin, and Gagne, associated with Operating System Concepts, 9th Edition, Wiley, 2013)

and

http://retis.sssup.it/~giorgio/rts-MECS.html

(slides by Buttazzo, associated with Hard Real-Time Computing Systems, 3rd Edition, Springer, 2011)

which has been edited to suit the needs of this course.

The slides are authorized for personal use only.

Any other use, redistribution, and any for profit sale of the slides (in any form) requires the consent of the copyright owners.





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



Operating System Concepts – 9th Edition



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



Operating System Concepts – 9th Edition



```
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;
```



}



```
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 */





POSIX Shared Memory

POSIX Shared Memory (access by name):

 Create a new shared memory object (or open an existing object and share it) using a name

```
shm_fd = shm_open("object_1", O_CREAT | O_RDRW, 0666);
```

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

```
• Use shared memory
char msg_0="Writing to shared memory";
sprintf(ptr, msg_0);
ptr += strlen(msg_0);
```





POSIX Producer

```
#include <stdio.h>
#include <stlib.h>
#include <string.h>
#include <fcntl.h>
#include <svs/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 obect */
void *ptr;
   /* create the shared memory object */
   shm_fd = shm_open(name, O_CREAT | O_RDRW, 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:
}
```





POSIX Consumer

```
#include <stdio.h>
#include <stlib.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 obect */
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;
```



Programming assignment (2)

Read the online manual (man) of the following system calls:

```
shm_open(), mmap(), shm_unlink()
```

or

```
shmget(), shmat(), shmdt(), shmctl()
```

- Implement a bounded buffer (both producer and consumer) using the POSIX API
 - Extra: without BUFFER_SIZE-1 elements limitation







- 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*) message size fixed or variable
 - receive(*message*)
- If P and Q wish to communicate, they need to:
 - establish a *communication link* between them
 - exchange messages via send/receive
- Implementation of communication link
 - physical (e.g., shared memory, hardware bus)
 - logical (e.g., direct or indirect, synchronous or asynchronous, automatic or explicit buffering)



Implementation Questions

- 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?





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 (must know identity of other)
 - 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
 - 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.

Who owns A?





Synchronization

Message passing may be either blocking or non-blocking

- Blocking is considered synchronous
 - Blocking send has the sender block until the message is received
 - Blocking receive has the receiver block until a message is available
- Non-blocking is considered asynchronous
 - Non-blocking send has the sender send the message and continue
 - Non-blocking receive has the receiver receive a valid message or null





Synchronization

- Different combinations possible
 - If both send and receive are blocking, we have a rendezvous
- Producer-consumer becomes trivial

```
message next_produced;
while (true) {
    /* produce an item in next produced */
    send(next_produced);
}
    message next_consumed;
    while (true) {
        receive(next_consumed);
        /* consume the item in next consumed */
    }
```





Buffering

- Queue of messages attached to the link; implemented in one of three ways
 - Zero capacity 0 messages 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





Sockets

Remote Procedure Calls









- 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





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
- 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





Operating System Concepts – 9th Edition

Silberschatz, Galvin and Gagne ©2013





Acts as a conduit allowing two processes to communicate

Issues

- Is communication unidirectional or bidirectional?
- In the case of two-way communication, is it half or full-duplex?
- Must there exist a relationship (i.e. *parent-child*) between the communicating processes?
- Can the pipes be used over a network?





Ordinary Pipes

- Ordinary Pipes allow communication in standard producer-consumer style
- Producer writes to one end (the write-end of the pipe)
- Consumer reads from the other end (the read-end of the pipe)
- Ordinary pipes are therefore unidirectional
 - See pipes in shell
- Require parent-child relationship between communicating processes



- Remember to close unused end of pipe from start, and other end once done
- See Unix code samples in textbook





Named Pipes

- Named Pipes (FIFO in Unix) are more powerful than ordinary pipes
- Communication is bidirectional
- No parent-child relationship is necessary between the communicating processes (but must be on same machine)
- Several processes can use the named pipe for communication
- Provided on both UNIX and Windows systems



Programming assignment (3)

- Read the online manual (man) of the following system calls:
 - pipe(), read(), write(), close()
- Implement a Producer-Consumer two-process interaction using a pipe
 - Producer: "produces" lower-case characters (taken from user input or generated in any other way)
 - Consumer: converts each character into upper case, and displays it
- **Extra: use** mkfifo() and implement a FIFO
 - Try with two processes that do not belong to the same subtree





Chapter 4: Threads

- 1. Overview
 - Motivation
 - Benefits
 - Multicore Programming
- 2. Multithreading Models
 - Many-to-one, one-to-one, many-to-many, two-level
- 3. Threading Issues
 - Creation, cancellation, signal handling, data, scheduling
- 4. Thread Libraries (POSIX)
- 5. Threads in Linux (NPTL)





Objectives

- To introduce the notion of a thread—a fundamental unit of CPU utilization that forms the basis of multithreaded computer systems
- To examine issues related to multithreaded programming
- To discuss the APIs for the Pthreads thread library
- To cover operating system support for threads in Linux



Single and Multithreaded Processes







Motivation

- Thread: basic unit of CPU utilization
 - Threads run within application
 - Most modern applications are multithreaded
- Multiple tasks with the application can be implemented by separate threads
 - Update display
 - Fetch data
 - Spell checking
 - Answer a network request
- Process creation is heavy-weight while thread creation is light-weight
- Can simplify code, increase efficiency
 - Web servers, RPC servers
 - Kernels are generally multithreaded



Multithreaded Server Architecture





Operating System Concepts – 9th Edition



- Responsiveness may allow continued execution if part of process is blocked, especially important for user interfaces
- Resource Sharing threads share resources of process, easier than shared memory or message passing
- Economy cheaper than process creation, thread switching lower overhead than context switching
- **Scalability** process can take advantage of multiprocessor architectures


User Threads and Kernel Threads

- Support for threads provided at user level or at kernel level
- **User threads** management done by user-level threads library
 - Three primary thread libraries:
 - POSIX Pthreads
 - Win32 threads
 - Java threads
- Kernel threads Supported by the Kernel
 - Examples virtually all general purpose operating systems, including:
 - Windows
 - Solaris
 - ▶ Linux
 - Tru64 UNIX
 - Mac OS X





Multithreading Models

Relationship between user-level and kernel-level threads?

- Many-to-One
- One-to-One
- Many-to-Many
- Two-level





Many-to-One

- Many user-level threads mapped to single kernel thread
- © Efficient thread management in user space
- ^(S) One thread blocking causes all to block
- Multiple threads may not run in parallel on multicore system because only one may be in kernel at a time
- Few systems currently use this model
- Examples:
 - Solaris Green Threads
 - GNU Portable Threads







One-to-One

- Each user-level thread maps to kernel thread
- © More concurrency than many-to-one
- ^(S) Creating a user-level thread creates a kernel thread (overhead)
 - Number of threads per process sometimes restricted due to overhead
- Examples
 Windows NT/XP/2000
 Linux
 Solaris 9 and later
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k
 k





Many-to-Many Model

- Multiplexes many user-level threads to a smaller or equal number of kernel threads
- Allows developer to create as many threads as she wishes, and threads can run concurrently
- Allows the operating system to create a sufficient number of kernel threads
 - A thread blocking does not block other threads



- Examples:
 - Solaris prior to version 9
 - Windows NT/2000 with the ThreadFiber package





Two-level Model

Similar to M:M, except that it allows a user thread to be **bound** to kernel thread







Threading Issues

- Semantics of **fork()** and **exec()** system calls
- Signal handling
 - Synchronous and asynchronous
- Thread cancellation of target thread
 - Asynchronous or deferred
- Thread-local storage





Semantics of fork() and exec()

- Does fork() duplicate only the calling thread or all threads?
 - Some UNIXes have two versions of fork
- Exec() usually works as normal replace the running process including all threads





Signal Handling

- Signals are used in UNIX systems to notify a process that a particular event has occurred.
- A signal handler is used to process signals
 - 1. Signal is generated by particular event
 - 2. Signal is delivered to a process
 - 3. Signal is handled by one of two signal handlers:
 - 1. default
 - 2. user-defined
- Every signal has default handler that kernel runs when handling signal
 - User-defined signal handler can override default
 - For single-threaded, signal delivered to process
- Where should a signal be delivered for multi-threaded?
 - Deliver the signal to the thread to which the signal applies
 - Deliver the signal to every thread in the process
 - Deliver the signal to certain threads in the process
 - Assign a specific thread to receive all signals for the process





Thread Cancellation

- Terminating a thread before it has finished
- Thread to be canceled is target thread
- Two general approaches:
 - Asynchronous cancellation terminates the target thread immediately
 - Deferred cancellation allows the target thread to periodically check if it should be cancelled
- Pthread code to create and cancel a thread:

```
pthread_t tid;
/* create the thread */
pthread_create(&tid, 0, worker, NULL);
. . .
/* cancel the thread */
pthread_cancel(tid);
```





Thread Cancellation (Cont.)

Invoking thread cancellation requests cancellation, but actual cancellation depends on thread state

Mode	State	Туре
Off	Disabled	-
Deferred	Enabled	Deferred
Asynchronous	Enabled	Asynchronous

- If thread has cancellation disabled, cancellation remains pending until thread enables it
- Default type is deferred
 - Cancellation only occurs when thread reaches cancellation point
 - > l.e. pthread_testcancel()
 - Then cleanup handler is invoked
- On Linux systems, thread cancellation is handled through signals





Thread-Local Storage

- Thread-local storage (TLS) allows each thread to have its own copy of data
- Useful when you do not have control over the thread creation process (i.e., when using a thread pool)
- Different from local variables
 - Local variables visible only during single function invocation
 - TLS visible across function invocations
- Similar to static data
 - TLS is unique to each thread





Thread Libraries

- Thread library provides programmer with API for creating and managing threads
- Two primary ways of implementing
 - Library entirely in user space
 - No kernel support
 - No system calls
 - Kernel-level library supported by the OS
 - Code and data structures exist in kernel space
 - API function \rightarrow system call





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

```
#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[])
ſ
  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++)</pre>
     sum += i;
  pthread_exit(0);
```

Figure 4.9 Multithreaded C program using the Pthreads API.





Linux Threads

- NPTL (Native POSIX Thread Library; formerly: LinuxThreads)
- Linux refers to them as tasks rather than threads
- Thread creation is done through clone () system call
- clone() allows a child task to share the address space of the parent task (process)
 - Flags control behavior

flag	meaning
CLONE_FS	File-system information is shared.
CLONE_VM	The same memory space is shared.
CLONE_SIGHAND	Signal handlers are shared.
CLONE_FILES	The set of open files is shared.

struct task_struct points to process data structures (shared or unique)



2.53

Programming assignment (4)

- Read the online manual (man) of:
 - pthreads

and of the following system calls:

- pthread_create(), pthread_join(), pthread_exit
- olone(), waitpid()
- Implement a bounded buffer using shared memory (same as in assignment
 2) and the pthreads library
 - Extra: implement the same, without pthreads (use clone())





Chapter 5: CPU Scheduling

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- Thread Scheduling
- Multiple-Processor Scheduling
- Operating Systems Examples
- Algorithm Evaluation







- 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





Histogram of CPU-burst Times





Operating System Concepts – 9th Edition



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
- If scheduling only under 1 and 4: nonpreemptive (cooperative)
- Otherwise: preemptive scheduling
 - 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 timesharing environment)



Scheduling Algorithm Optimization Criteria

- Max CPU utilization
- Max throughput
- Min turnaround time
- Min waiting time
- Min response time (variance)



First-Come, First-Served (FCFS) Scheduling

(Note: for simplicity, illustrations show only one CPU-burst per process)

Nonpreemptive

Process	Burst Time
P_1	24
P_2	3
P_3	3

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

Waiting time?

Average waiting time?



First-Come, First-Served (FCFS) Scheduling

(Note: for simplicity, illustrations show only one CPU-burst per process)

Nonpreemptive

Process	Burst Time
P_1	24
P_2	3
P_3	3

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



2.64

- 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

What is the Gantt chart for the schedule?

Waiting time?

Average waiting time?





FCFS Scheduling (Cont.)

Suppose that the processes arrive in the order:

$$P_2, P_3, P_1$$

Gantt chart for the schedule:



- 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
 - In batch systems (long-term scheduling)
 - Good estimate: fast response (incentive)
 - Time-limit-exceeded error + resubmission (penalty)







<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?



Operating System Concepts – 9th Edition

Silberschatz, Galvin and Gagne ©2013





<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



Determining Length of Next CPU Burst

CPU scheduling: can only estimate – similar to the previous ones?

- 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 $\frac{1}{2}$
- Preemptive version called shortest-remaining-time-first





Prediction of the Length of the Next CPU Burst



Examples of Exponential Averaging

- α =0
 - $\tau_{n+1} = \tau_n$
 - Recent history does not count
- α =1
 - $\tau_{n+1} = \alpha t_n$
 - Only the actual last CPU burst counts
- Since both α and (1 α) 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

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

Preemptive SJF Gantt Chart

Average waiting time?



Example of Shortest-remaining-time-first

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

Process	<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 = highest priority)
 - Preemptive
 - Nonpreemptive
- SJF is priority scheduling where priority is the inverse of predicted next CPU burst time
- Priorities can be defined internally or externally (wrt the OS)
 - Time limits, memory requirements, number of open files, politics
- Problem = Starvation low priority processes may never execute
- Solution = Aging as time progresses increase the priority of the process





Example of Priority Scheduling

Process	Burst Time	Priority
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?



Operating System Concepts – 9th Edition

Silberschatz, Galvin and Gagne ©2013



Example of Priority Scheduling

Process	Burst Time	<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
 - $q \text{ large} \Rightarrow \text{FCFS}$
 - *q* small ⇒ processor sharing. *q* must be large wrt context switch, otherwise overhead is too high



Example of RR with Time Quantum = 4

Process	Burst Time
P_1	24
P_2	3
P_3	3

Gantt chart? (q=4)



Example of RR with Time Quantum = 4

Process	<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











Turnaround Time Varies With The Time Quantum



- Increasing q does not necessarily improve average turnaround
- In general: improvement whenever most processes finish their next CPU burst in 1q



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
 - ▶ E.g., 80% to foreground in RR, 20% to background in FCFS.



Multilevel Queue Scheduling

highest 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



- $Q_0 RR$ with time quantum 8 msec
- $Q_1 RR$ time quantum 16 msec
- $Q_2 FCFS$



Scheduling

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



Exercise

Process	Burst Time	Priority
P_1	10	3
P_2	1	1
P_3	2	3
P_4	1	4
P_5	5	2

Arrived in the top-to-bottom order.

• Illustrate FCFS, SJF, nonpreemptive priority, RR (q=1) using Gantt

- Turnaround times?
- Waiting times?
- Which algorithm results in the minimum average waiting time?

