# Processes and Threads

Alessio Vecchio
alessio.vecchio@unipi.it
Dip.di Ingegneria dell'Informazione
Università di Pisa

# Outline

- Processes

- Threads

- Scheduling algorithms

# Process Concept

- Program is *passive* entity stored on disk (**executable file**), process is *active*

  - Program becomes originates when executable file loaded into memory and run

- 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** – a program in execution; process execution must progress in sequential fashion
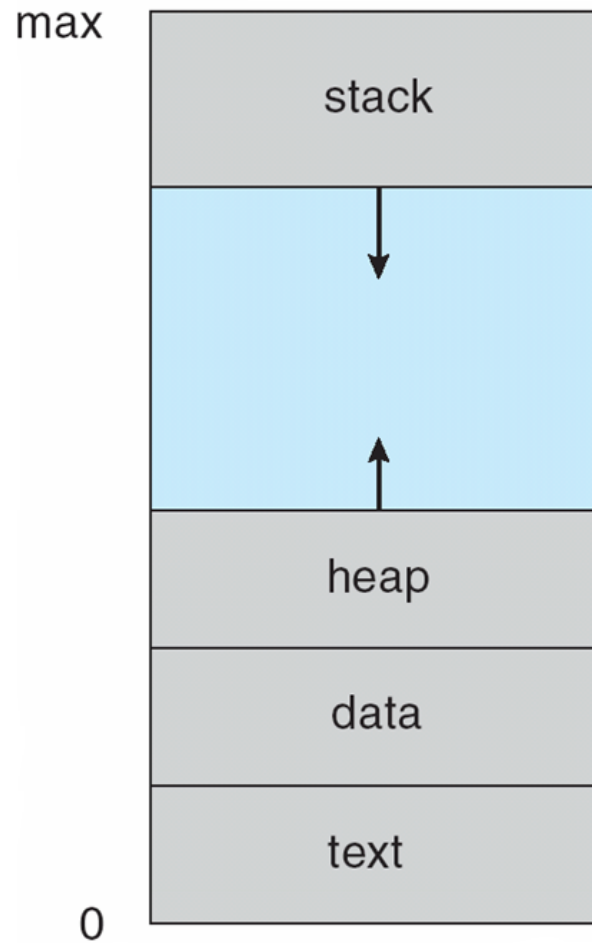
# Process Concept

- 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

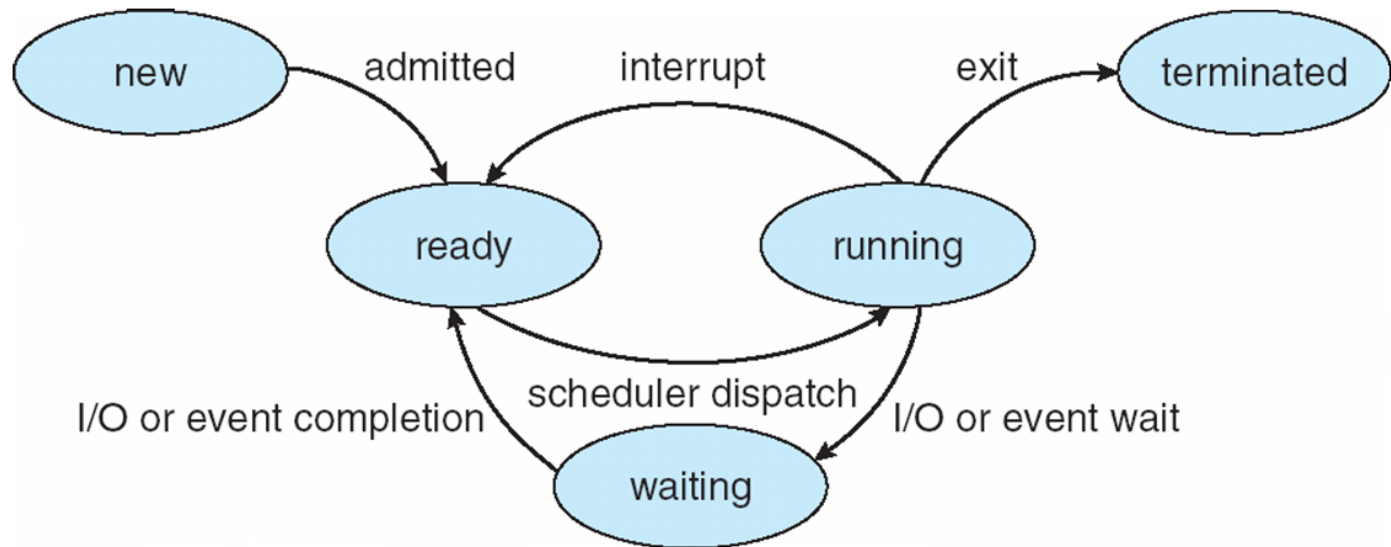  - **Heap** containing memory dynamically allocated during run time

# Process 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
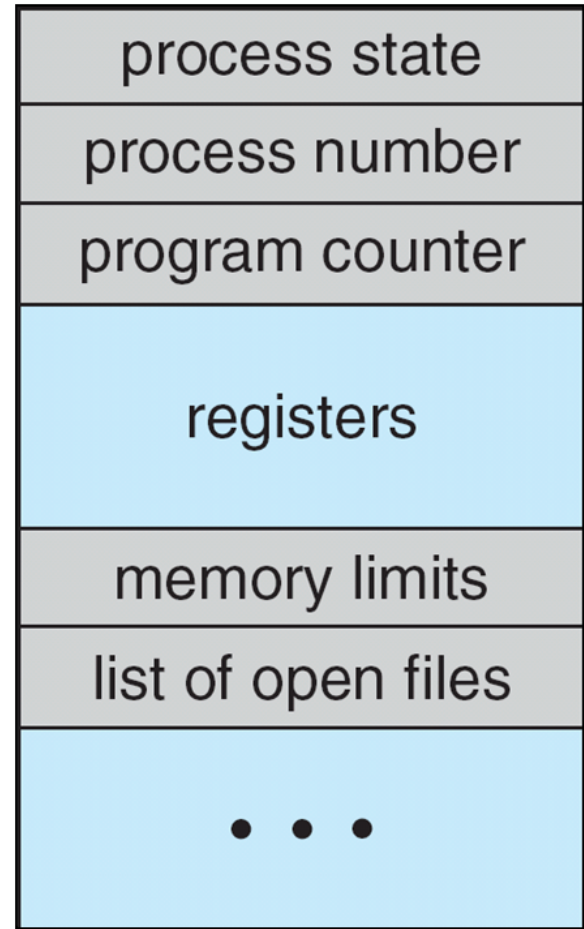
# Diagram of Process State

# 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

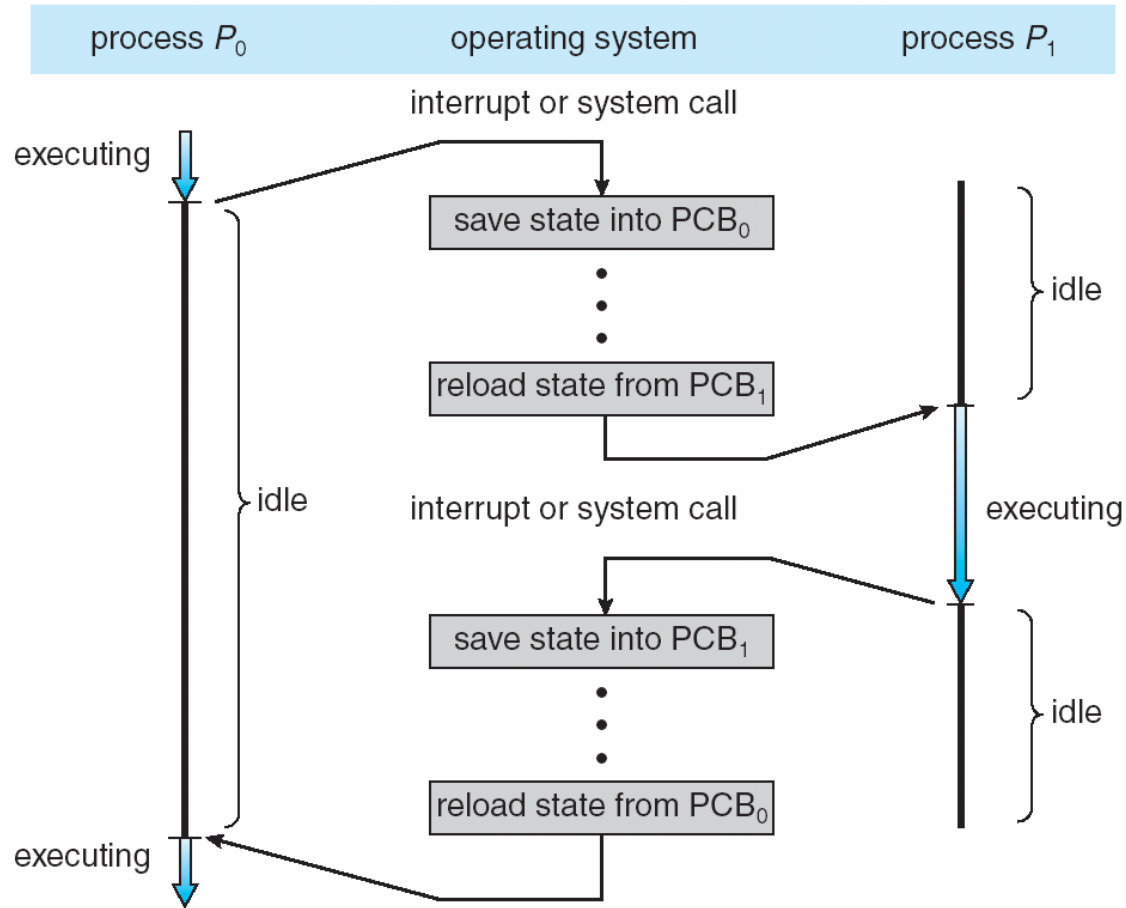- I/O status information – I/O devices allocated to process, list of open files

| |
|---|
| process state |
| process number |
| program counter |
| registers |
| memory limits |
| list of open files |
| • • • |

# 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

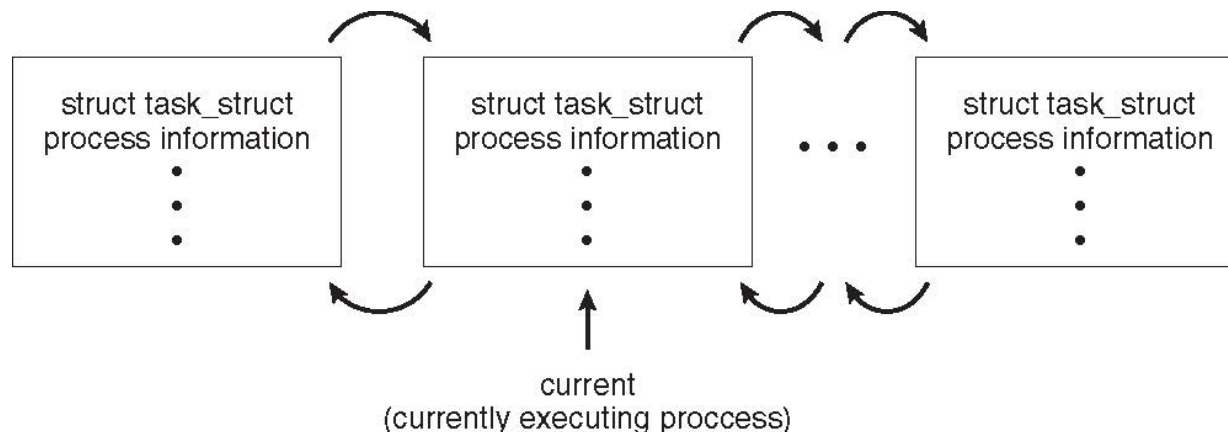  - Some hardware provides multiple sets of registers per CPU ➔ multiple contexts loaded at once

# CPU Switch From Process to Process

# Process Representation in Linux

Represented by the C structure `task_struct`

```
pid t_pid; /* process identifier */
long state; /* state of the process */
unsigned int time_slice /* scheduling information */
struct task_struct *parent; /* this process's parent */
struct list_head children; /* this process's children */
struct files_struct *files; /* list of open files */
struct mm_struct *mm; /* address space of this process */
```
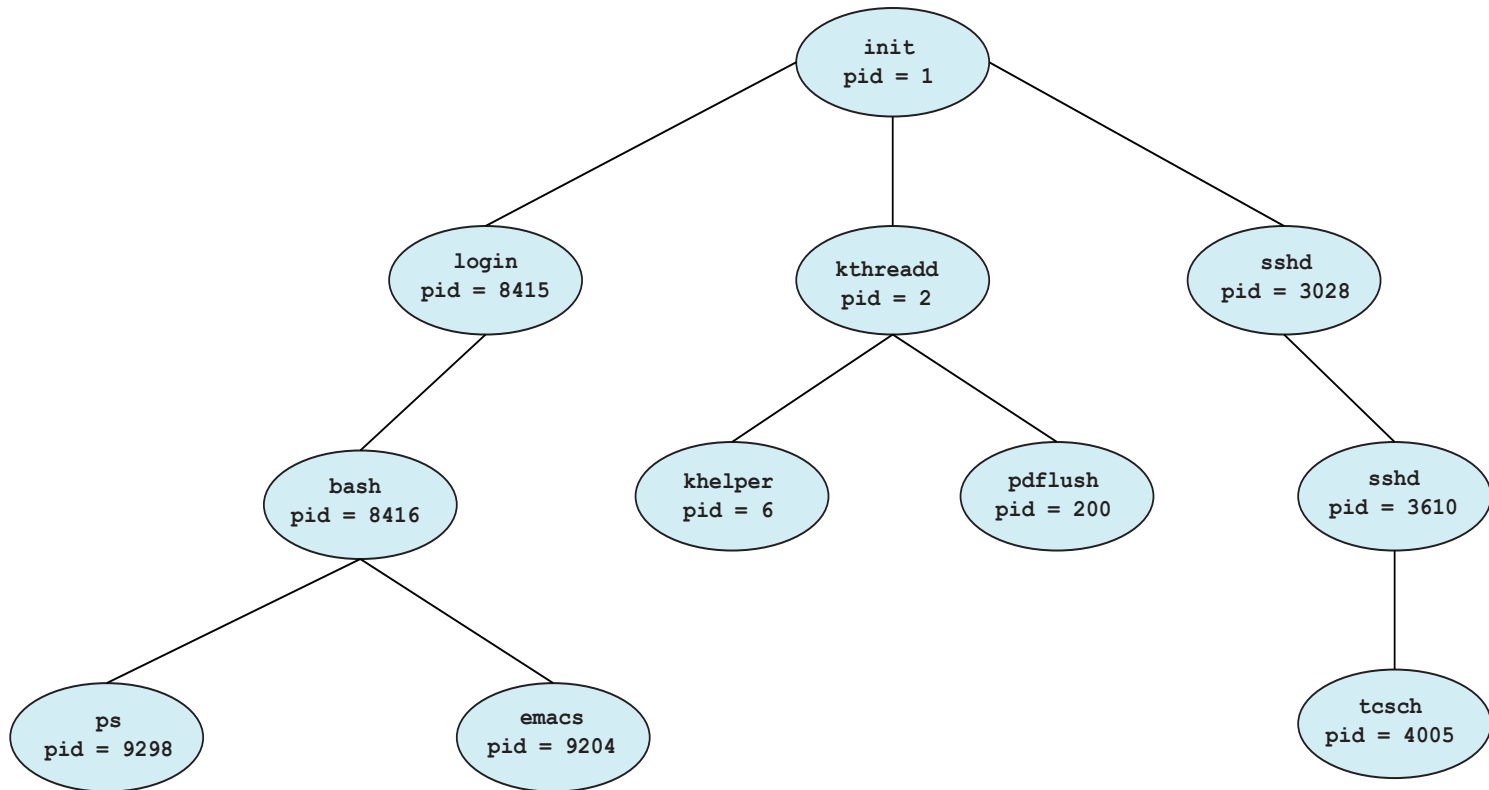
# Process Creation

- **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
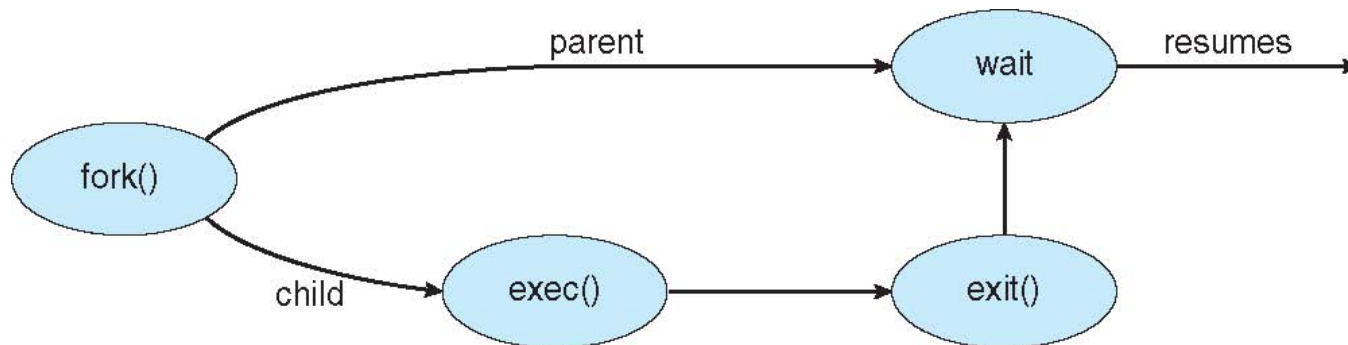  - Parent waits until children terminate

# A Tree of Processes in Linux

# Process Creation (Cont.)

- Address space
  - Child duplicate of parent
  - Child has a program loaded into it
- UNIX examples
  - `fork()` system call creates new process
  - `exec()` system call used after a `fork()` to replace the process' memory space with a new program

# 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 exists 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**

# Example in UNIX

```cpp
#include <iostream>
#include <unistd.h>
#include <stdlib.h>
#include <sys/types.h>
#include <sys/wait.h>
using namespace std;

int main(int argc, char* argv[]) {
  pid_t pid;
  pid=fork(); /* genera un nuovo processo */
  if(pid<0) { /* errore */
      cout << "Errore nella creazione del processo\n";
      exit(-1);
  } else if(pid==0) { /* processo figlio */
      execlp("/usr/bin/touch", "touch", "my_new_file", NULL);
  } else { /* processo genitore */
      int status;
      pid = wait(&status);
      cout << "Il processo figlio " << pid << " ha terminato\n";
      exit(0);
  }
}
```

# Multiprocess Architecture – Chrome Browser

- Many web browsers ran as single process (some still do)

  - If one web site causes trouble, entire browser can hang or crash

- Google Chrome Browser is multiprocess with 3 different types of processes:

  - **Browser** process manages user interface, disk and network I/O

  - **Renderer** process renders web pages, deals with HTML, Javascript. A new renderer created for each website opened

    ‣ Runs in **sandbox** restricting disk and network I/O, minimizing effect of security exploits

  - **Plug-in** process for each type of plug-in



*Each tab represents a separate process*
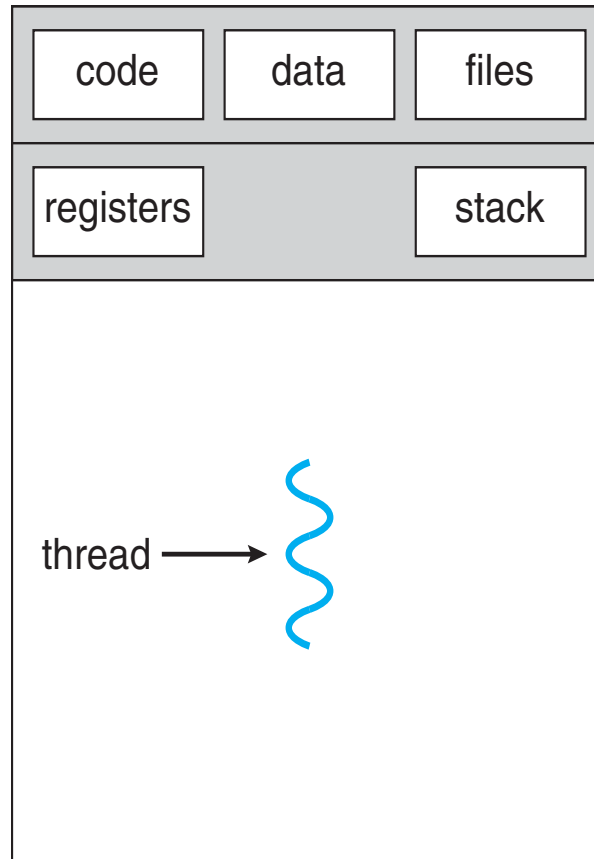
# Multitasking in Mobile Systems

- Some mobile systems (e.g., early version of iOS) allow only one process to run, others suspended

- Due to screen real estate, user interface limits iOS provides for a
  - Single **foreground** process- controlled via user interface
  - Multiple **background** processes– in memory, running, but not on the display, and with limits
  - Limits include single, short task, receiving notification of events, specific long-running tasks like audio playback

- Android runs foreground and background, with fewer limits
  - Background process uses a **service** to perform tasks
  - Service can keep running even if background process is suspended
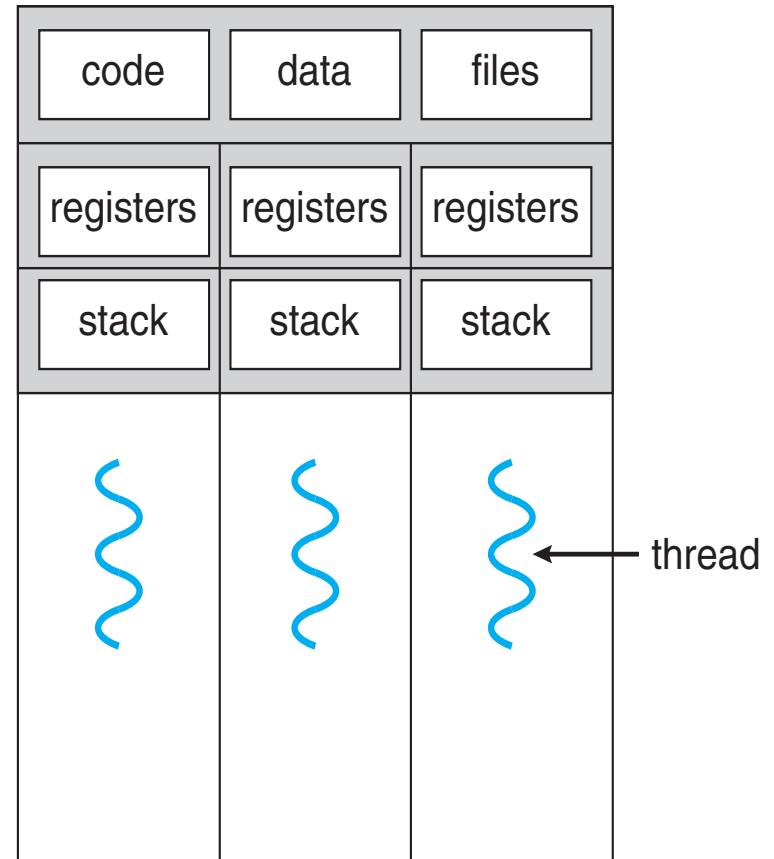  - Service has no user interface, small memory use

# Threads

- Most modern applications are multithreaded

- Threads run within application

- 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

- Kernels are generally multithreaded
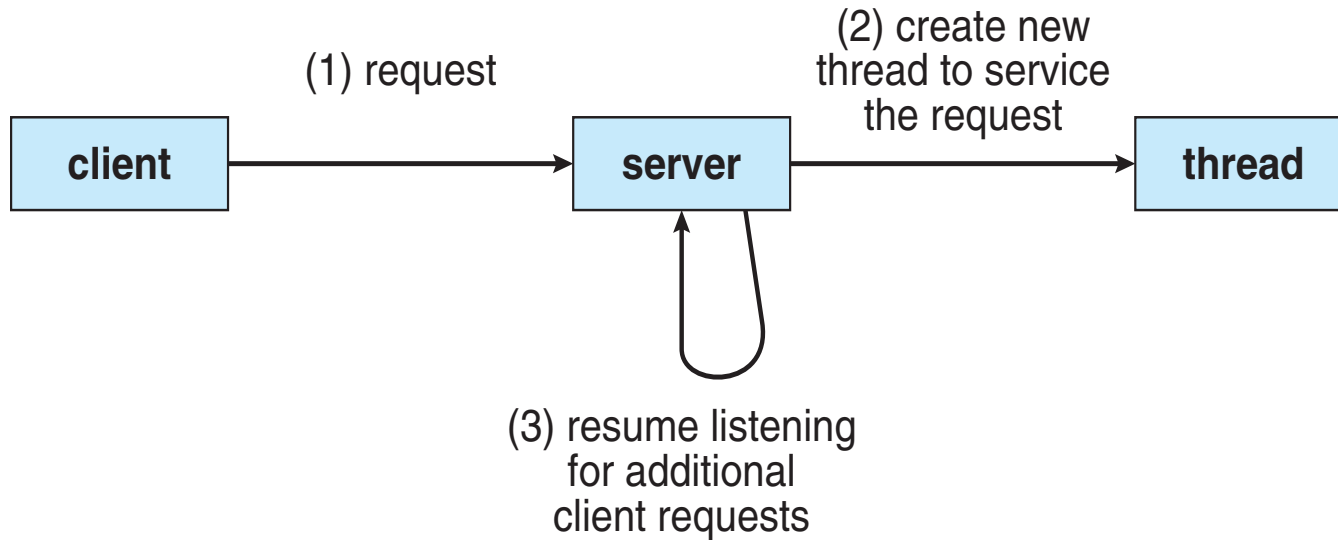
# Single and Multithreaded Processes



single-threaded process　　　　multithreaded process

# Multithreaded Server Architecture

# Benefits

- **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

# Multicore Programming

- **Multicore** or **multiprocessor** systems putting pressure on programmers, challenges include:

  - **Dividing activities**

  - **Balance**

  - **Data splitting**

  - **Data dependency**

  - **Testing and debugging**

- *Parallelism* implies a system can perform more than one task simultaneously

- *Concurrency* supports more than one task making progress

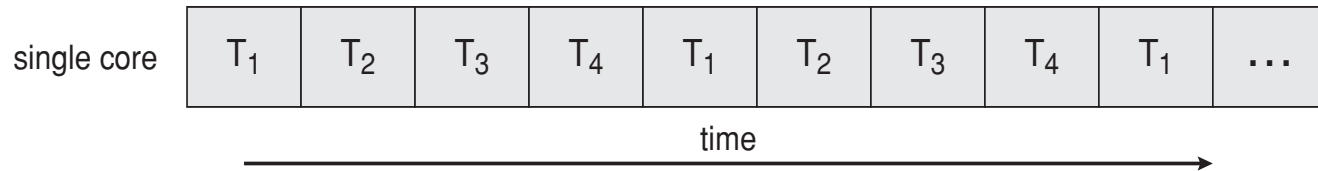  - Single processor / core, scheduler providing concurrency
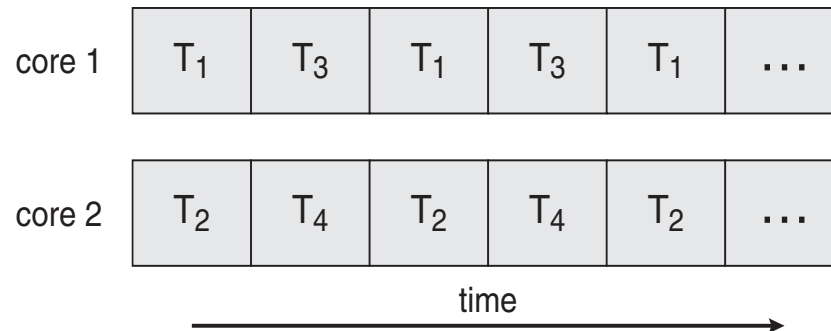
# Multicore Programming (Cont.)

- Types of parallelism

  - **Data parallelism** – distributes subsets of the same data across multiple cores, same operation on each

  - **Task parallelism** – distributing threads across cores, each thread performing unique operation

- As # of threads grows, so does architectural support for threading

  - CPUs have cores as well as ***hardware threads***

  - Consider Oracle SPARC T4 with 8 cores, and 8 hardware threads per core

# Concurrency vs. Parallelism

- **Concurrent execution on single-core system:**

| single core | $T_1$ | $T_2$ | $T_3$ | $T_4$ | $T_1$ | $T_2$ | $T_3$ | $T_4$ | $T_1$ | ... |

time →

- **Parallelism on a multi-core system:**

| core 1 | $T_1$ | $T_3$ | $T_1$ | $T_3$ | $T_1$ | ... |

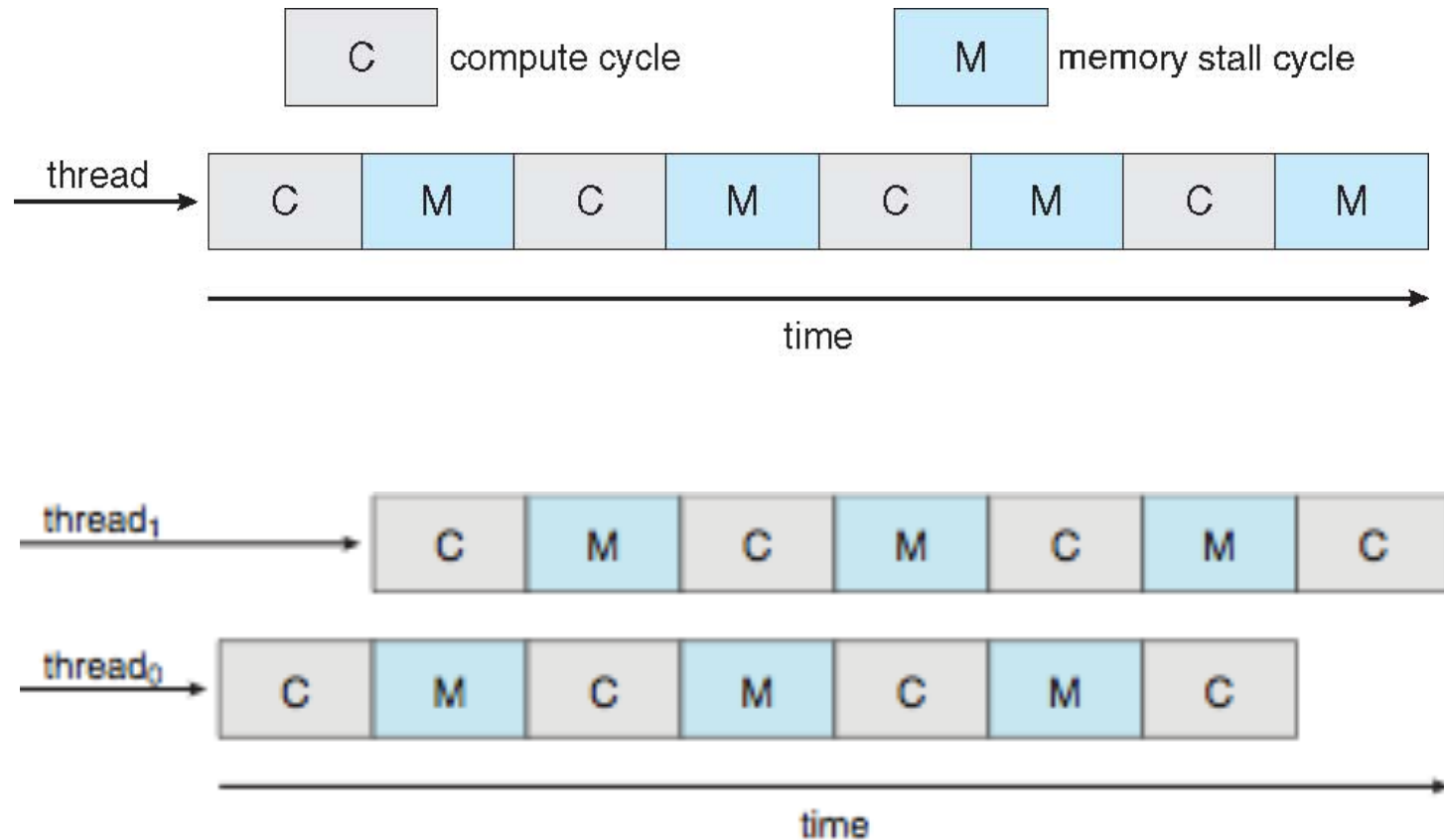| core 2 | $T_2$ | $T_4$ | $T_2$ | $T_4$ | $T_2$ | ... |

time →

# 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

# Amdahl's Law

- Identifies performance gains from adding additional cores to an application that has both serial and parallel components

- $S$ is serial portion

- $N$ processing cores

$$speedup \leq \frac{1}{S + \frac{(1-S)}{N}}$$

- That is, if application is 75% parallel / 25% serial, moving from 1 to 2 cores results in speedup of 1.6 times

- As $N$ approaches infinity, speedup approaches $1 / S$

- Serial portion of an application has important effect on performance gained by adding additional cores

# User Threads and Kernel Threads

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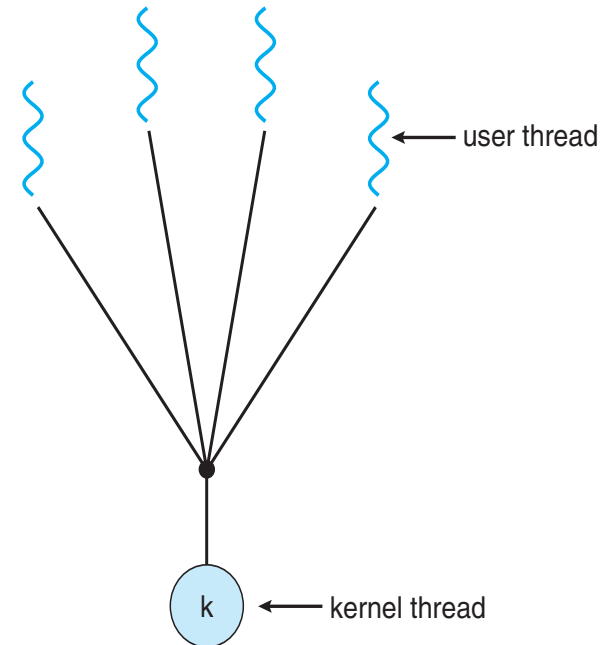  - Windows
  - Solaris
  - Linux
  - Tru64 UNIX
  - Mac OS X

# Multithreading Models

- Many-to-One

- One-to-One

- Many-to-Many

# Many-to-One

■ Many user-level threads mapped to single kernel thread

■ One thread blocking causes all to block

■ Multiple threads may not run in parallel on muticore system because only one may be in kernel at a time

■ Few systems currently use this model

■ Examples:

   ● **Solaris Green Threads**

   ● **GNU Portable Threads**
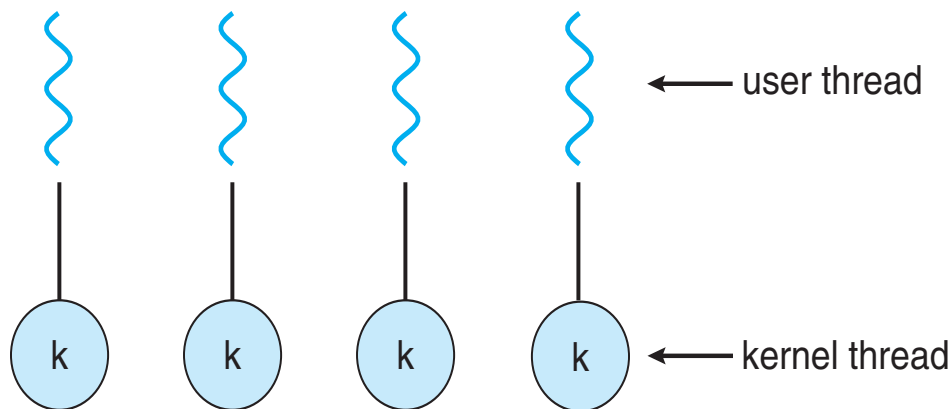
user thread

k ← kernel thread

# One-to-One

- Each user-level thread maps to kernel thread

- Creating a user-level thread creates a kernel thread

- More concurrency than many-to-one

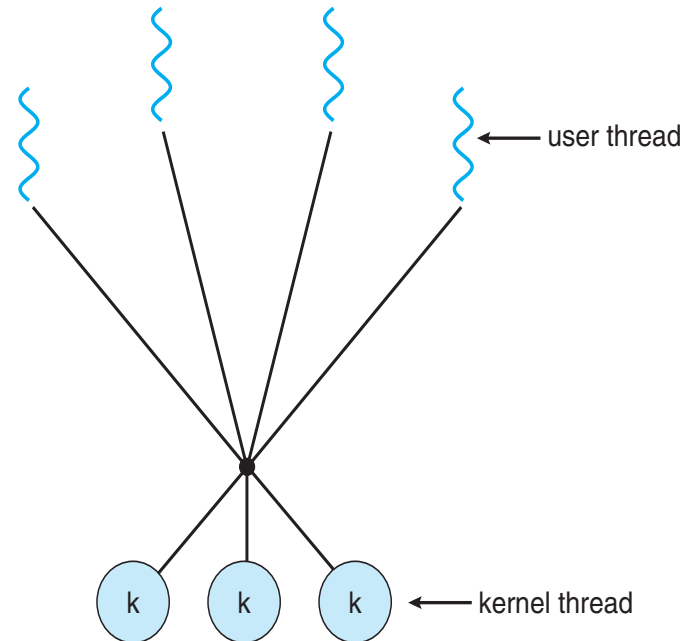- Number of threads per process sometimes restricted due to overhead

- Examples
  - Windows
  - Linux
  - Solaris 9 and later
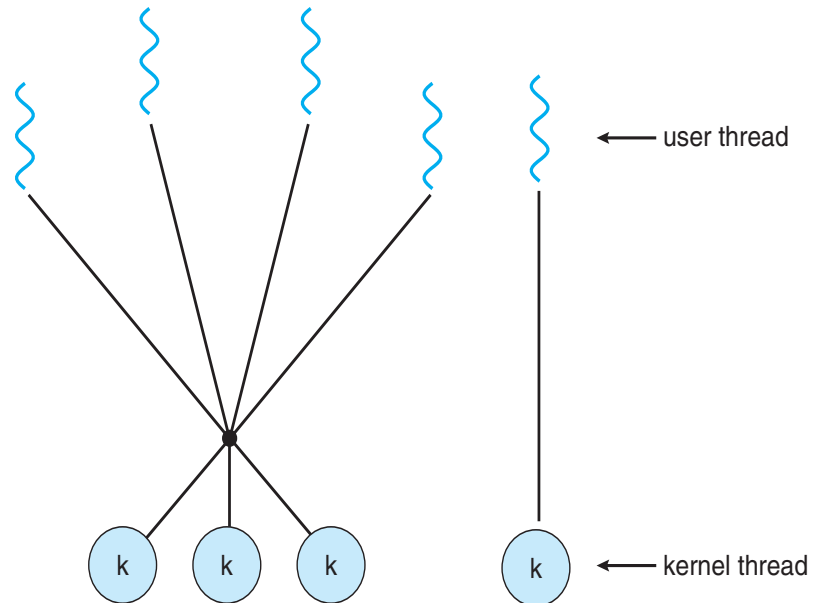
← user thread

← kernel thread

# Many-to-Many Model

- Allows many user level threads to be mapped to many kernel threads

- Allows the operating system to create a sufficient number of kernel threads

- Solaris prior to version 9

- Windows with the *ThreadFiber* package



user thread

kernel thread

# Two-level Model

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

■ Examples

- IRIX

- HP-UX

- Tru64 UNIX

- Solaris 8 and earlier

user thread

kernel thread

k   k   k        k

# 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)

# Threading Issues

- Semantics of **fork()** and **exec()** system calls

- Signal handling
  - Synchronous and asynchronous

- Thread cancellation of target thread
  - Asynchronous or deferred

# 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

# Signal Handling (Cont.)

■ 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 | Type |
|------|-------|------|
| 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**

▸ I.e. `pthread_testcancel()`

▸ Then **cleanup handler** is invoked

■ On Linux systems, thread cancellation is handled through signals

# Operating System Examples
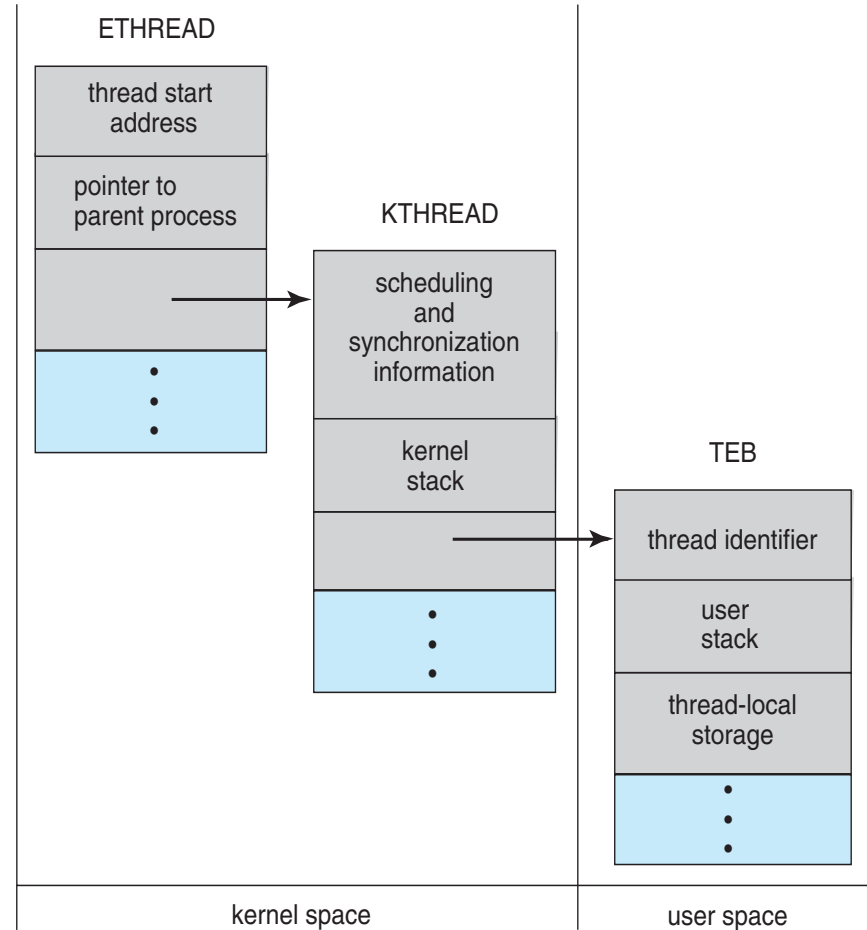
- Windows Threads
- Linux Threads

# Windows Threads

- Windows implements the Windows API – primary API for Win 98, Win NT, Win 2000, Win XP, and Win 7

- Implements the one-to-one mapping, kernel-level

- Each thread contains

  - A thread id

  - Register set representing state of processor

  - Separate user and kernel stacks for when thread runs in user mode or kernel mode

  - Private data storage area used by run-time libraries and dynamic link libraries (DLLs)

- The register set, stacks, and private storage area are known as the **context** of the thread

# Windows Threads Data Structures

The primary data structures of a thread include:

- ETHREAD (executive thread block) – includes pointer to process to which thread belongs and to KTHREAD, in kernel space
- KTHREAD (kernel thread block) – scheduling and synchronization info, kernel-mode stack, pointer to TEB, in kernel space
- TEB (thread environment block) – thread id, user-mode stack, thread-local storage, in user space

ETHREAD

| thread start address |
| pointer to parent process |
| |
| •••• |

KTHREAD

| scheduling and synchronization information |
| kernel stack |
| |
| •••• |

TEB

| thread identifier |
| user stack |
| thread-local storage |
| •••• |

kernel space | user space

# Linux Threads

■ 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)
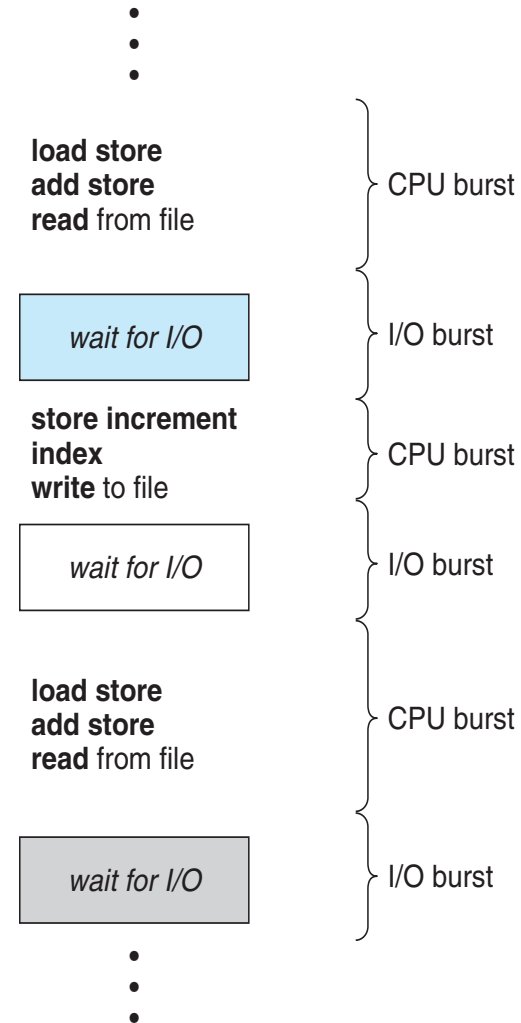
# Process Scheduling

- 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
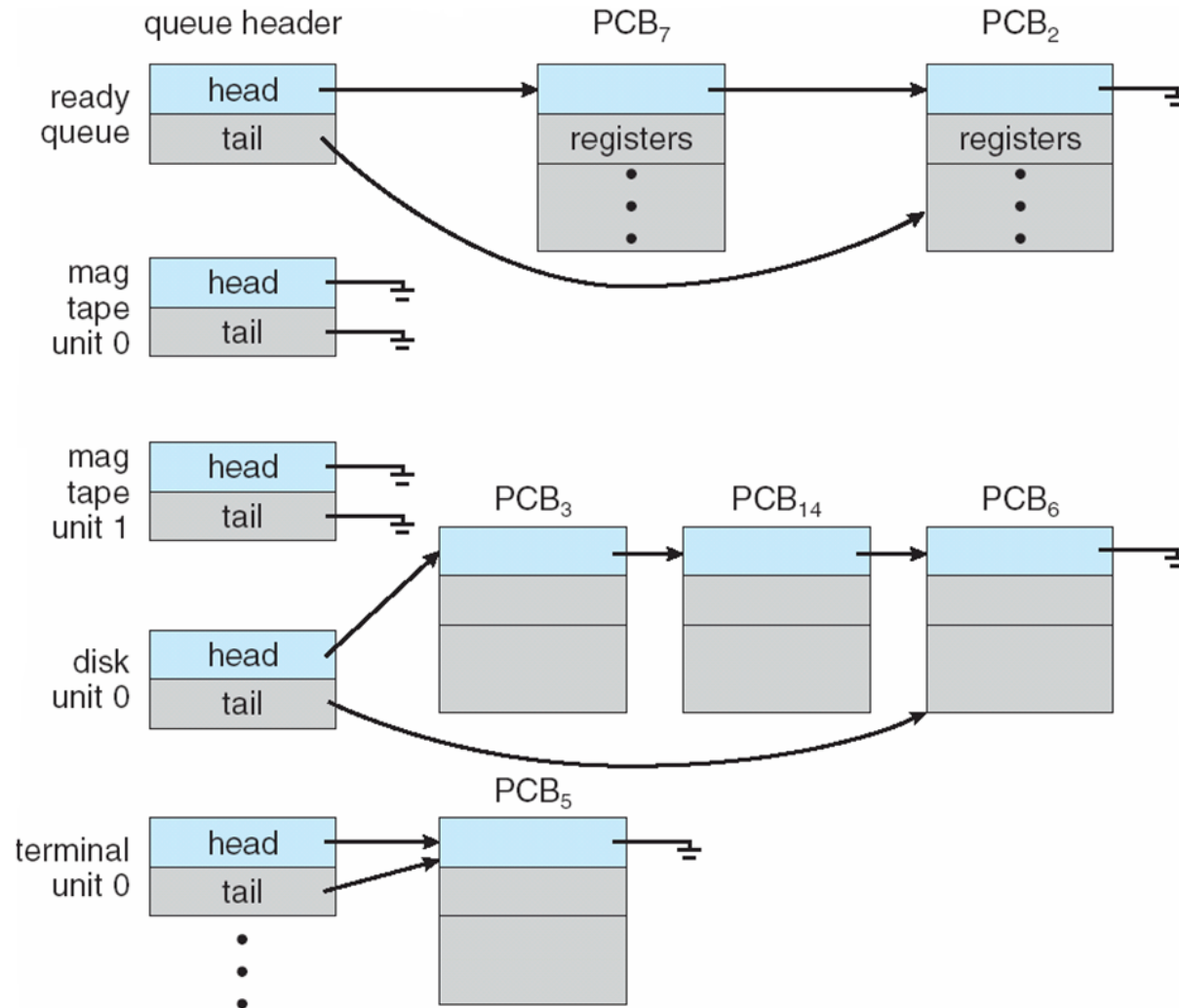
  - Processes migrate among the various queues

# 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**
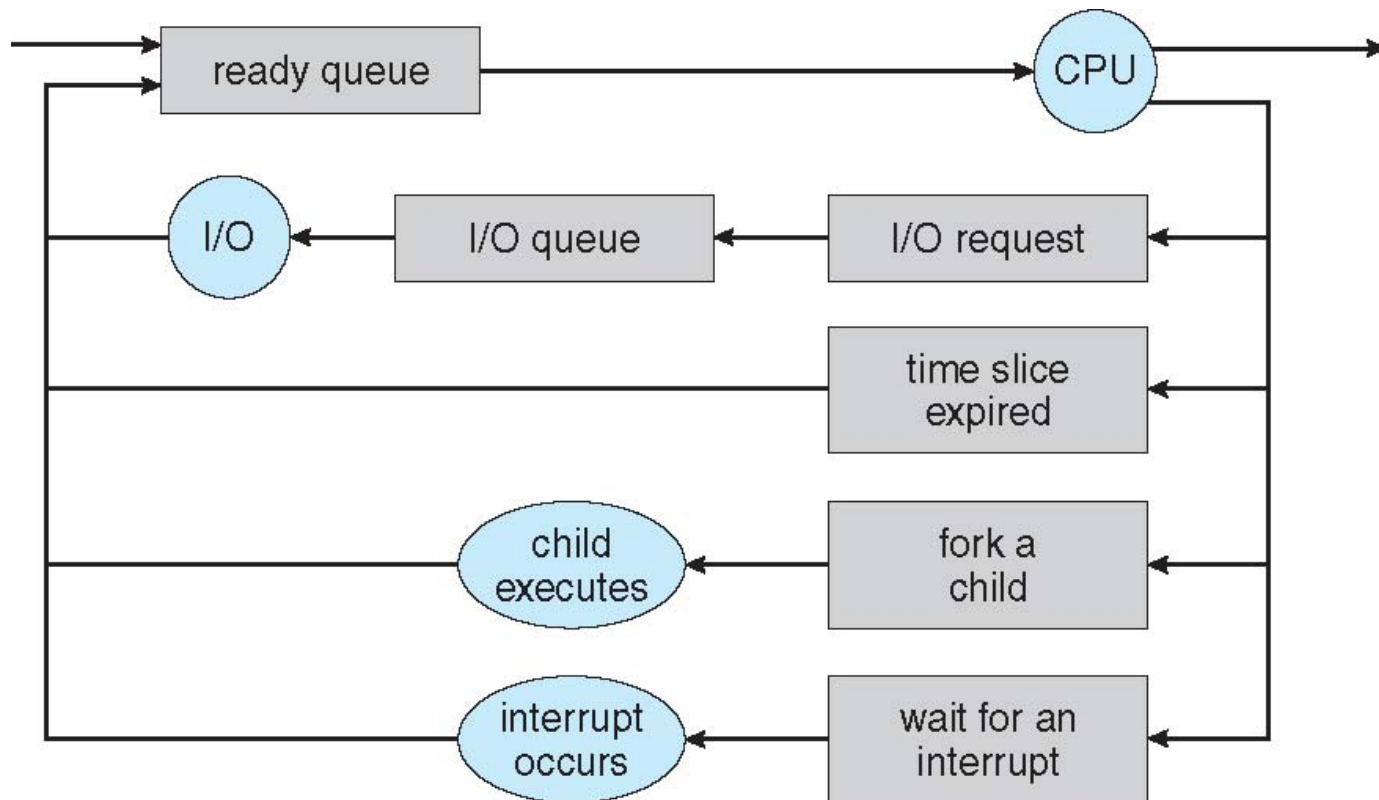
- CPU burst distribution is of main concern

```
•
•
•
```

**load store**
**add store**
**read** from file                    CPU burst

| *wait for I/O* |                      I/O burst

**store increment**
**index**
**write** to file                      CPU burst

| *wait for I/O* |                      I/O burst

**load store**
**add store**
**read** from file                     CPU burst

| *wait for I/O* |                      I/O burst

```
•
•
•
```

# Representation of Process Scheduling

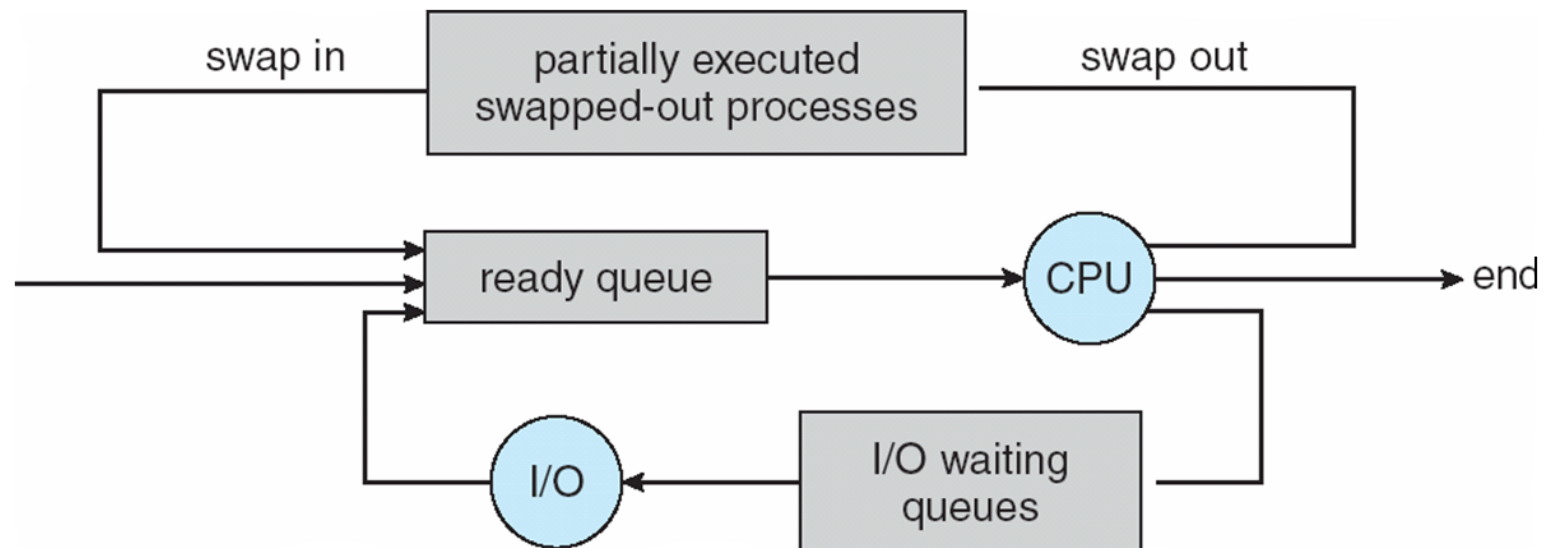■ **Queueing diagram** represents queues, resources, flows

# 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**

# 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

# 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)

# First- Come, First-Served (FCFS) Scheduling

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

- Suppose that the processes arrive in the order: $P_1$ , $P_2$ , $P_3$
The schedule is:

| $P_1$ | | $P_2$ | $P_3$ |
|---|---|---|---|
| 0 | 24 | 27 | 30 |

- 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 schedule is:

| P$_2$ | P$_3$ | P$_1$ |
|---|---|---|

0         3         6         30

■ 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

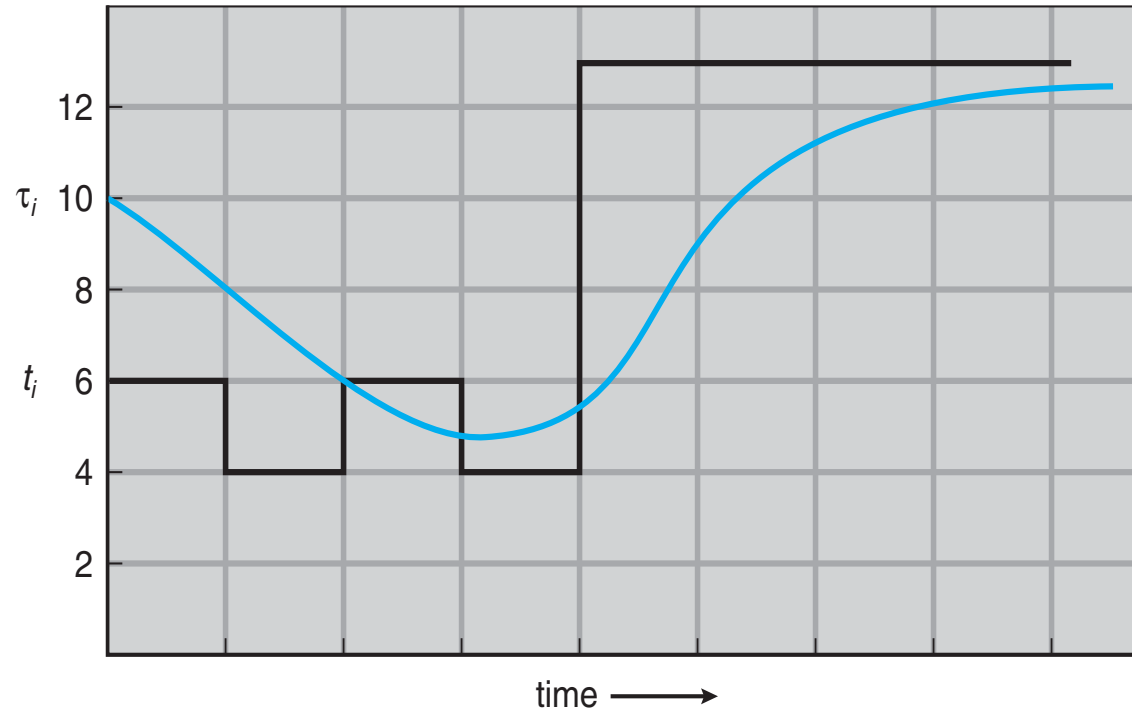| Process | Burst Time |
|---------|------------|
| $P_1$ | 6 |
| $P_2$ | 8 |
| $P_3$ | 7 |
| $P_4$ | 3 |

■ SJF scheduling chart

| $P_4$ | $P_1$ | $P_3$ | $P_2$ |
|-------|-------|-------|-------|

0       3               9                       16                      24

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

# Determining Length of Next CPU Burst

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

  - Then pick process with shortest predicted next CPU burst

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

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

- Commonly, α set to ½

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

# Prediction of the Length of the Next CPU Burst



|  | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| CPU burst ($t_i$) | | 6 | 4 | 6 | 4 | 13 | 13 | 13 | ... |
| "guess" ($\tau_i$) | 10 | 8 | 6 | 6 | 5 | 9 | 11 | 12 | ... |

61

# Examples of Exponential Averaging

- $\alpha = 0$
  - $\tau_{n+1} = \tau_n$
  - Recent history does not count
- $\alpha = 1$
  - $\tau_{n+1} = \alpha\, t_n$
  - Only the actual last CPU burst counts
- If we expand the formula, we get:

$$\tau_{n+1} = \alpha\, t_n + (1 - \alpha)\alpha\, t_{n-1} + \ldots$$
$$+ (1 - \alpha)^j \alpha\, t_{n-j} + \ldots$$

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

# Example of Shortest-remaining-time-first

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

| Process | *Arrival* Time | Burst Time |
|---------|----------------|------------|
| $P_1$ | 0 | 8 |
| $P_2$ | 1 | 4 |
| $P_3$ | 2 | 9 |
| $P_4$ | 3 | 5 |

■ *Preemptive* SJF

| $P_1$ | $P_2$ | $P_4$ | $P_1$ | $P_3$ |
|---|---|---|---|---|

0    1         5              10              17              26

■ 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

| 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:

| P$_1$ | P$_2$ | P$_1$ | P$_3$ | P$_4$ |
|---|---|---|---|---|

0   1            6                              16      18  19

■ 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

# Example of RR with Time Quantum = 4

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

- The execution is:

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

0     4     7     10     14     18     22     26     30

- 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

# 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 $Q_2$

# Operating System Examples

- Windows XP scheduling
- Linux scheduling

# Windows XP Scheduling

- Thread scheduling based on
  - Priority
  - Preemption
  - Time slice

- A thread is executed until one of the following event occurs
  - The thread has terminated its execution
  - The thread has exhausted its assigned time slice
  - The has executed a blocking system call
  - A higher-priority thread has entered the ready queue

# Kernel Priorities

- Kernel priority scheme: 32 priority levels

  - Real-time class (16-31)

  - Variable class (1-15)

  - Memory management thread (0)

- A different queue for each priority level

  - Queues are scanned from higher levels to lower levels

  - When no thread is found a special thread (idle thread) is executed

# Win32 API priorities

- **API Priority classes**
  - REALTIME_PRIORITY_CLASS       -> Real-time Class
  - HIGH_PRIORITY_CLASS       -> Variable Class
  - ABOVE_NORMAL_PRIORITY_CLASS       -> Variable Class
  - NORMAL_PRIORITY_CLASS       -> Variable Class
  - BELOW_NORMAL_PRIORITY_CLASS       -> Variable Class
  - IDLE_PRIORITY_CLASS       -> Variable Class

- **Relative Priority**
  - TIME_CRITICAL
  - HIGHEST
  - ABOVE_NORMAL
  - NORMAL
  - BELOW_NORMAL
  - LOWEST
  - IDLE

# Windows XP 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 |

**Default Base Priority**

# Class Priority Management

- A thread is stopped as soon as its time slice is exhausted
- Variable Class
  - If a thread stops because time slice is exhausted, its priority level is decreased
  - If a thread exits a waiting operation, its priority level is increased
    - waiting for data from keyboard, mouse -> significant increase
    - waiting for disk operations -> moderate increase
- Background/Foreground processes
  - The time slice of the foreground window is increased (typically by a factor 3)

# 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

# Priorities and Time-slice length



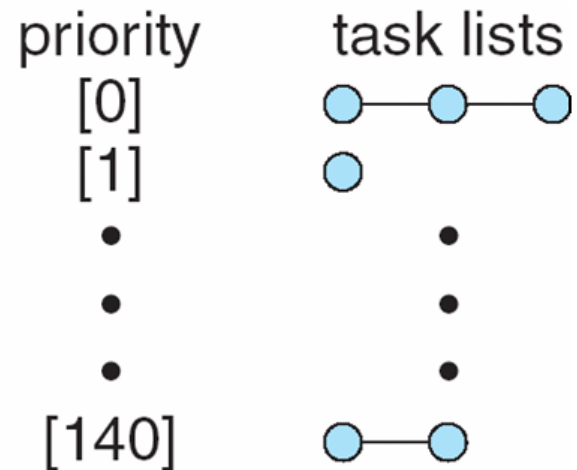| numeric priority | relative priority | | time quantum |
|---|---|---|---|
| 0 | highest | real-time tasks | 200 ms |
| • | | | |
| • | | | |
| • | | | |
| 99 | | | |
| 100 | | other tasks | |
| • | | | |
| • | | | |
| • | | | |
| 140 | lowest | | 10 ms |

# RunQueue

- The runqueue consists of two different arrays
    - Active array
    - Expired array

# Priority Calculation

- Real time tasks have static priority

- Time-sharing tasks have dynamic priority

  - Based on nice value +/- 5

  - +/- 5 depends on how much the task is interactive

    ▸ Tasks with low waiting times are assumed to be scarcely interactive

    ▸ Tasks with large waiting times are assumed to be highly interactive

- Priority re-computation is carried out every time a task has exhausted its time slice
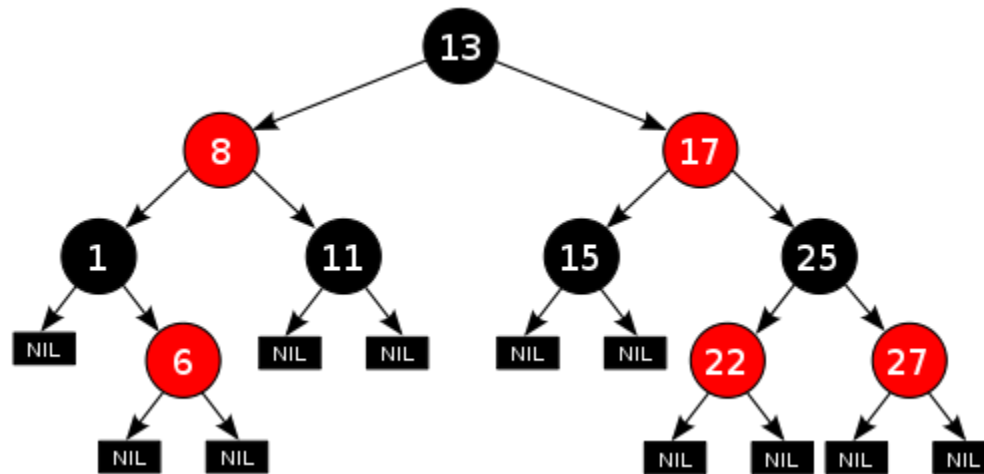
# Linux 2.6+ Scheduling

■ Recent versions of Linux include a new scheduler: Completely Fair Scheduler (CFS)

● Idea: when the time for tasks is not balanced (one or more tasks are not given a fair amount of time relative to others), then these tasks should be given time to execute.

■ CFS registers the amount of time provided to a given task (the virtual runtime)

■ The smaller a task's virtual runtime—meaning the smaller amount of time a task has been granted the CPU—the higher its need for the processor.

# Linux 2.6+ Scheduling

■ Tasks are stored in a red-black tree (not a queue) ordered in terms of virtual time

- A red-black tree is roughly balanced: any path in the tree will never be more than twice as long as any other path.

- Insert and deletion are O(log n)

# Linux 2.6+ Scheduling

- The scheduler picks the left-most node of the red-black tree. The task accounts for its time with the CPU by adding its execution time to the virtual runtime and is then inserted back into the tree if runnable.

- CFS doesn't use priorities directly but instead uses them as a decay factor for the time a task is permitted to execute.

  - Lower-priority tasks have higher factors of decay, where higher-priority tasks have lower factors of delay.

  - This means that the time a task is permitted to execute dissipates more quickly for a lower-priority task than for a higher-priority task.

  - This avoids maintaining run queues per priority.