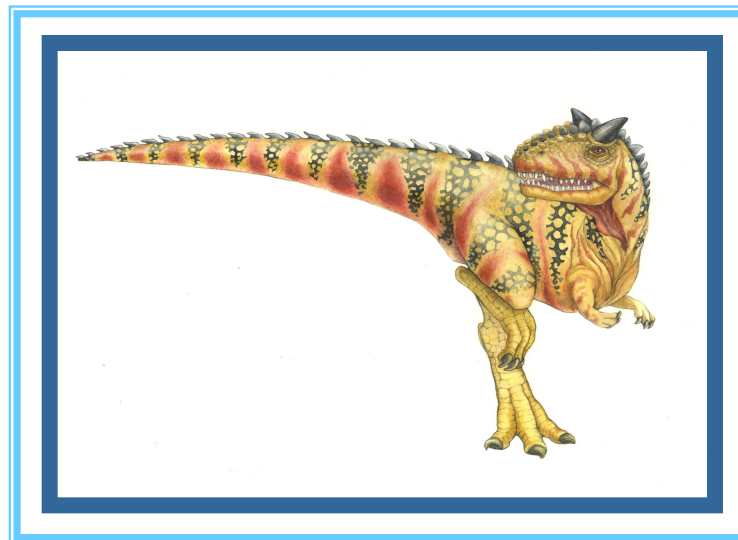
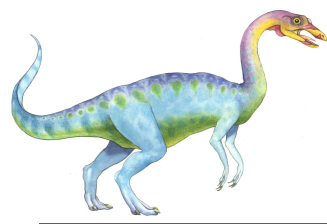


Chapter 4: Threads

(9th Edition)

NARZU TARANNUM(NAT)
LECTURER
DEPT. OF CSE, BRAC UNIVERSITY





Chapter 4: Threads

- Overview
- Multicore Programming
- Multithreading Models
- Thread Libraries
- Threading Issues
- Operating System Examples





Motivation

- 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





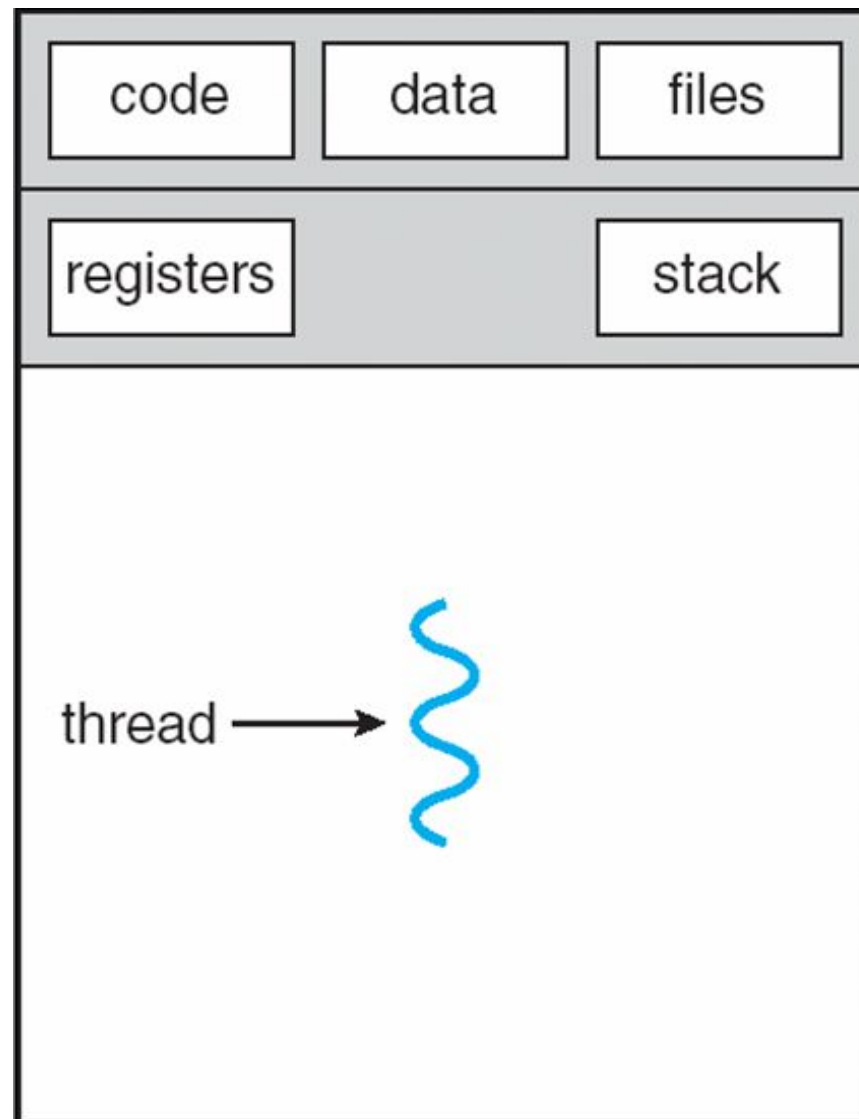
Threads

- A thread is a basic unit of CPU utilization
- It comprise
 - Thread ID
 - Program counter
 - Register set
 - Stack
- It shares with other threads belonging to the same process its code section, data section and other OS resources such as open files and signals.
- A traditional / heavyweight has a single thread of control.
- **If a process has multiple threads of control, it can perform more than one task at a time.**

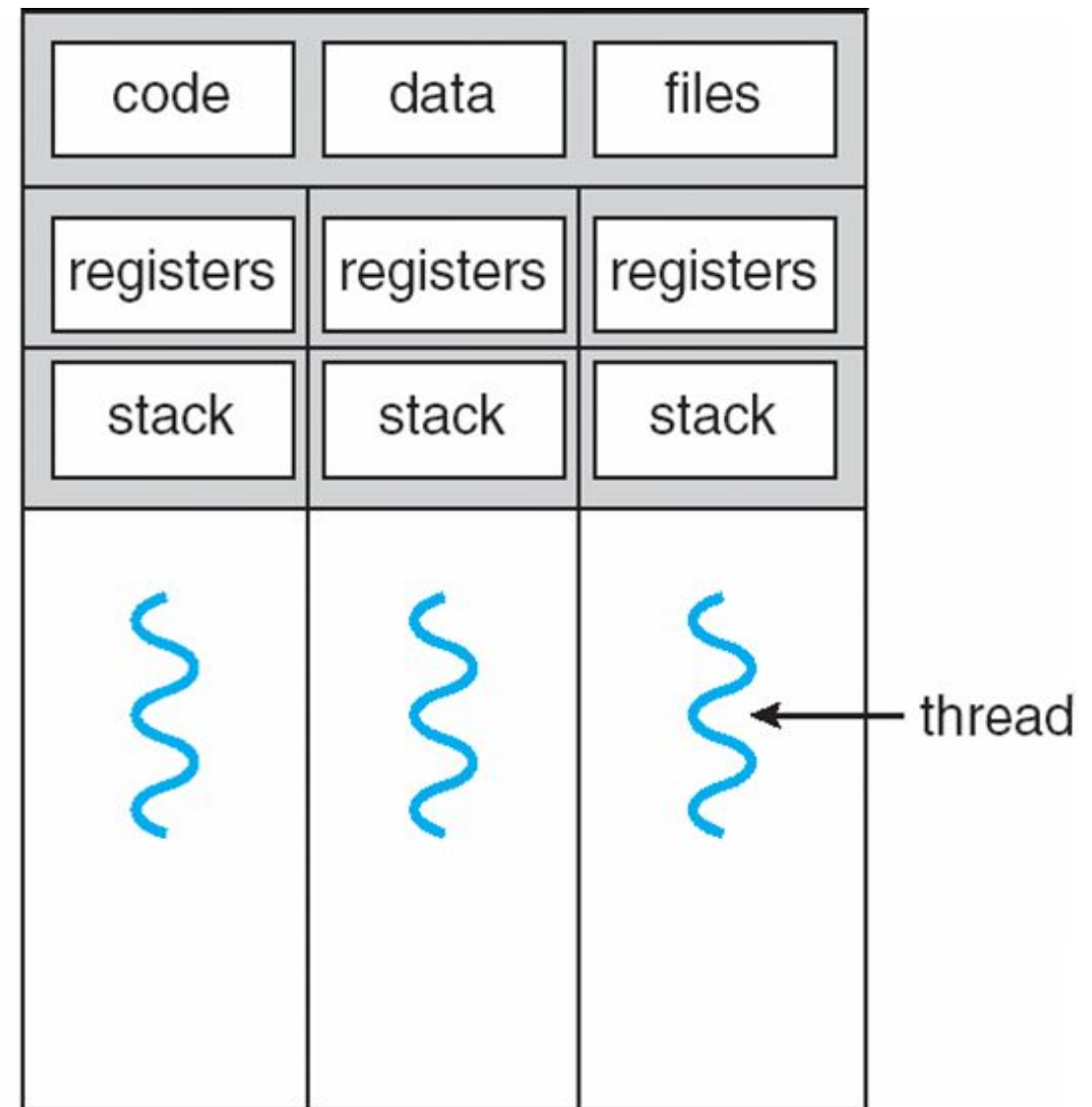




Single and Multithreaded Processes



single-threaded process

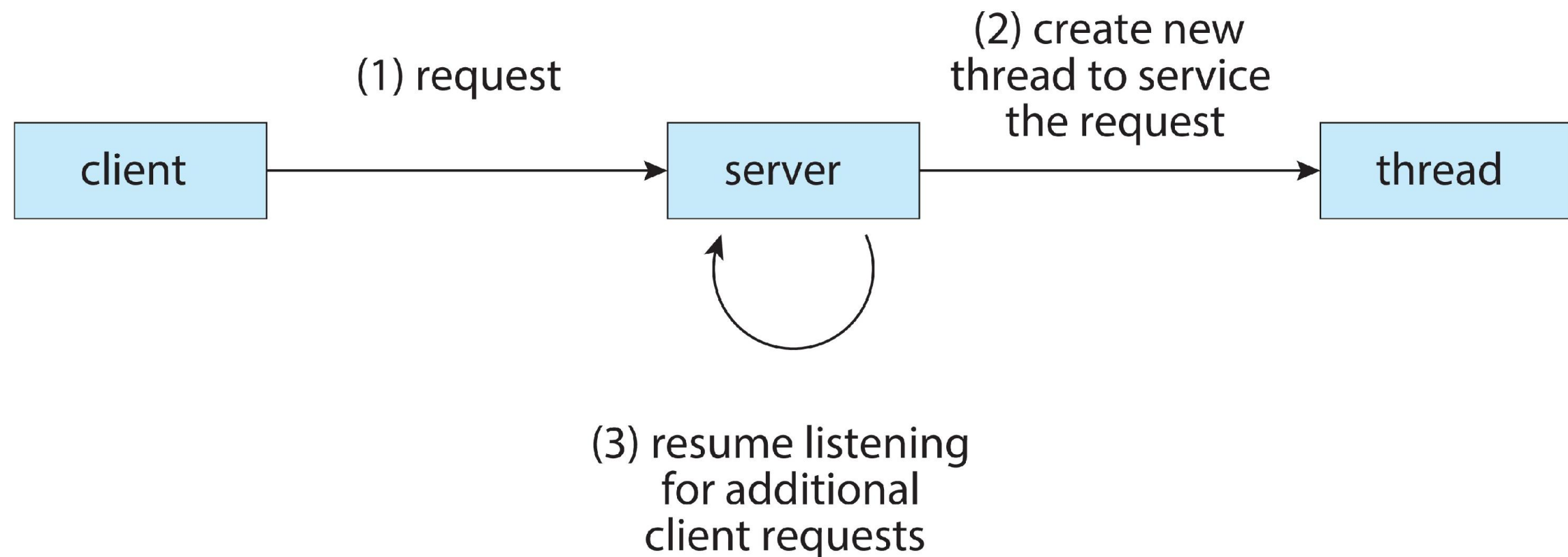


multithreaded process





Multithreaded Server Architecture



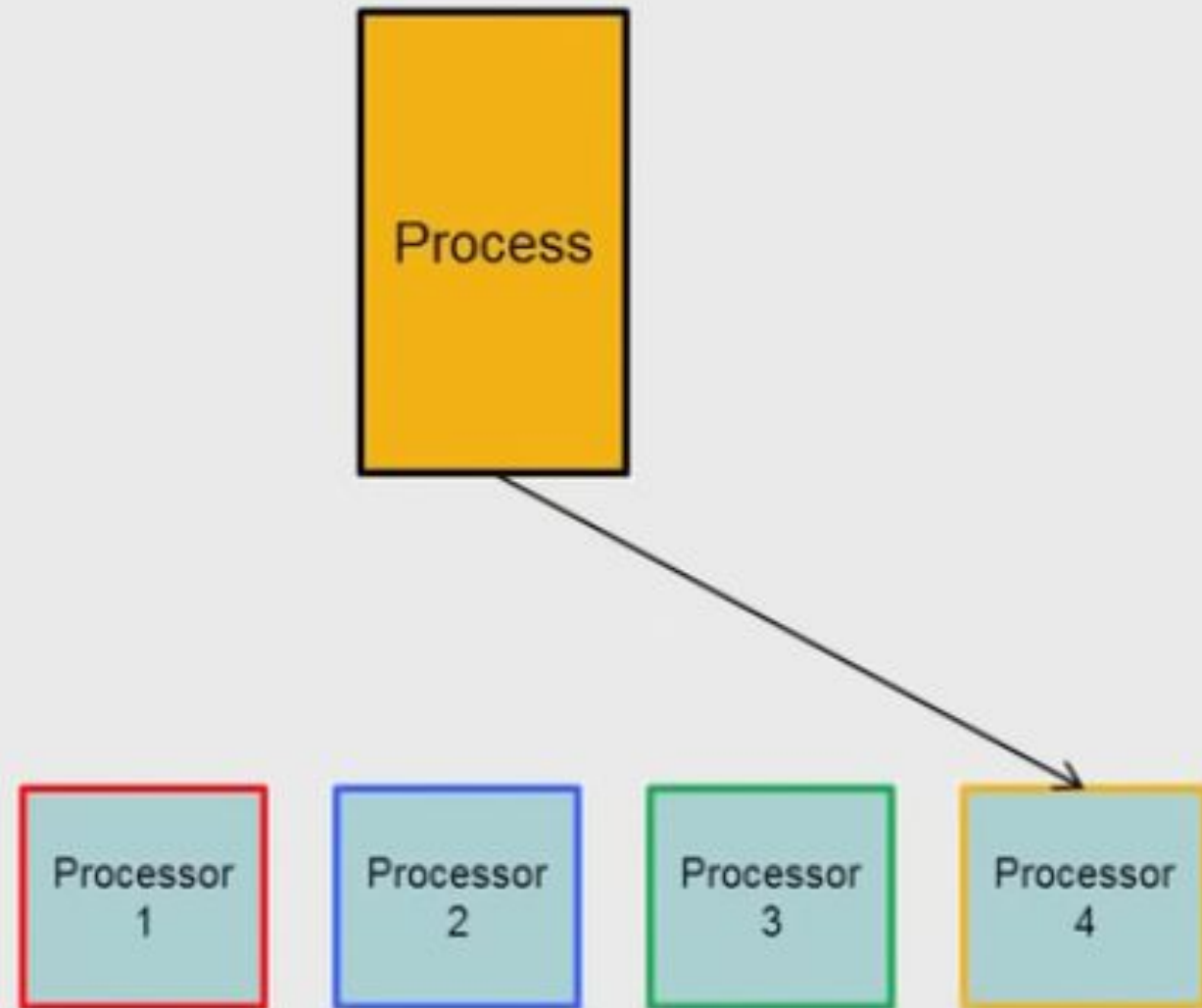
Consider this Example

```
#include <stdio.h>

unsigned long addall(){
    int i=0;
    unsigned long sum=0;

    while (i< 100000000){
        sum += i;
        i++;
    }
    return sum;
}

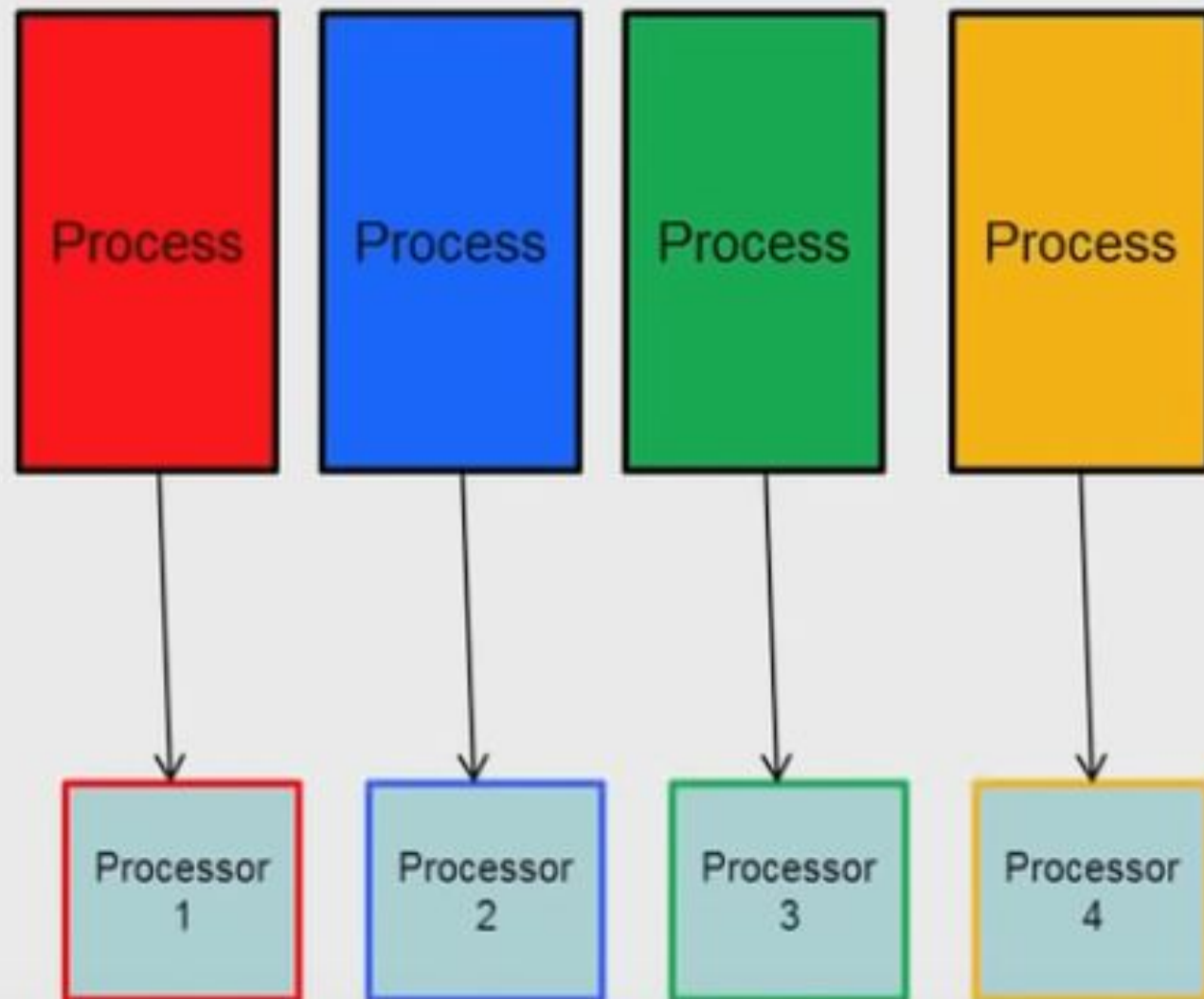
int main()
{
    unsigned long sum;
    srand(time(NULL));
    sum = addall();
    printf("%lu\n", sum);
}
```



Speeding up with multiple processes

$$10000000 / 4 = 2500000$$

Create 4 processes, each loop does 1/4th of the work



Properties:

4 fork system calls needed; one for creating each process

Each process is isolated from each other

IPC mechanisms to communicate – more system calls

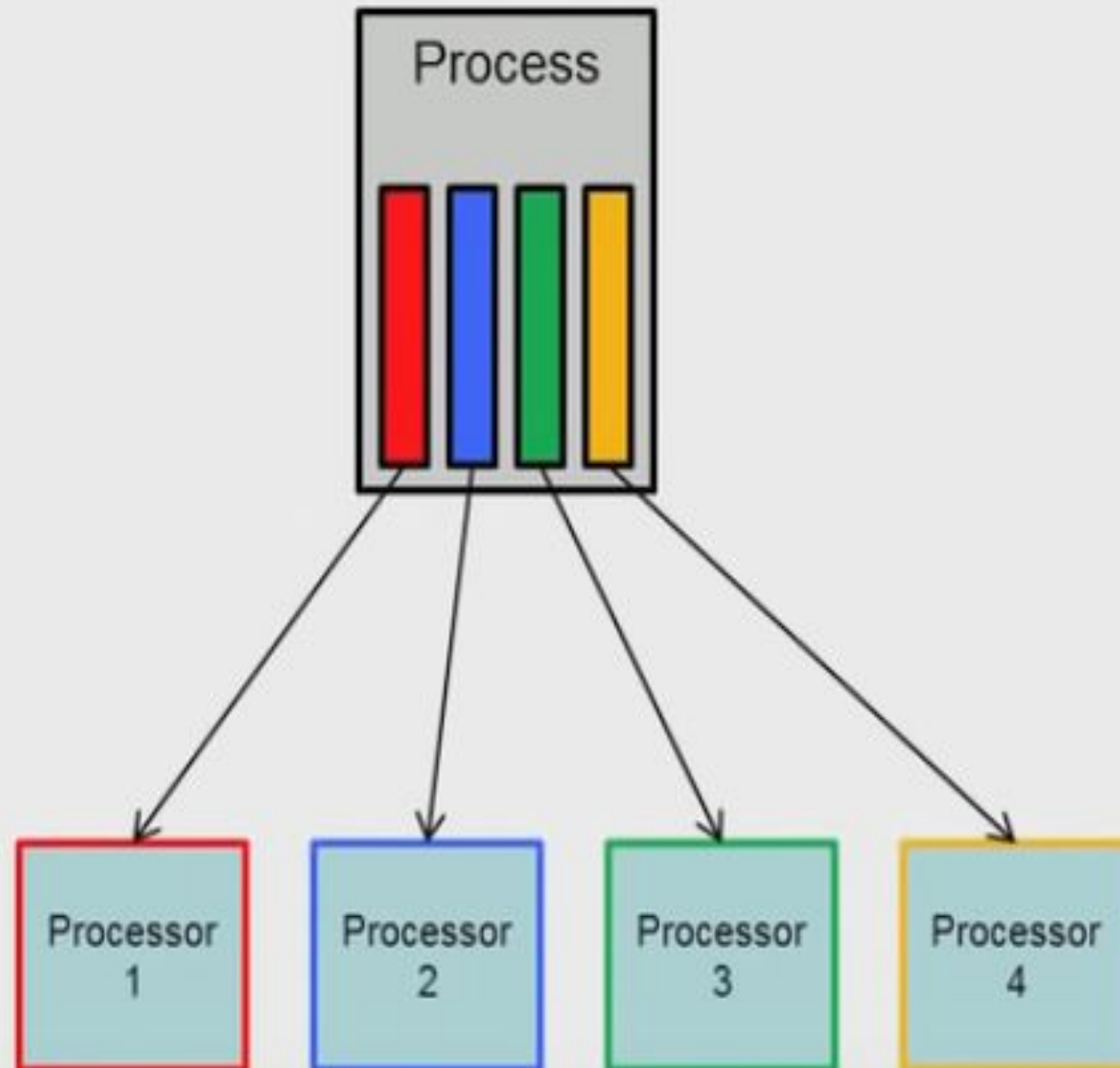
Process management with system calls

Each process has its own memory map – its own

Thread model

$$10000000 / 4 = 2500000$$

Create 1 process with 4 threads; each loop does $1/4^{\text{th}}$ of the work



Properties:

1 fork system call needed;
4 threads need to be
created --- much more
lighter.

Each thread is not isolated
from others

Management of threads
with fewer or no system
calls.

Threads vs Processes

- A thread has no data segment or heap
- A thread cannot live on its own. It needs to be attached to a process
- There can be more than one thread in a process. Each thread has its own stack
- If a thread dies, its stack is reclaimed
- A process has code, heap, stack, other segments
- A process has at-least one thread.
- Threads within a process share the same code, files.
- If a process dies, all threads die.



Benefits

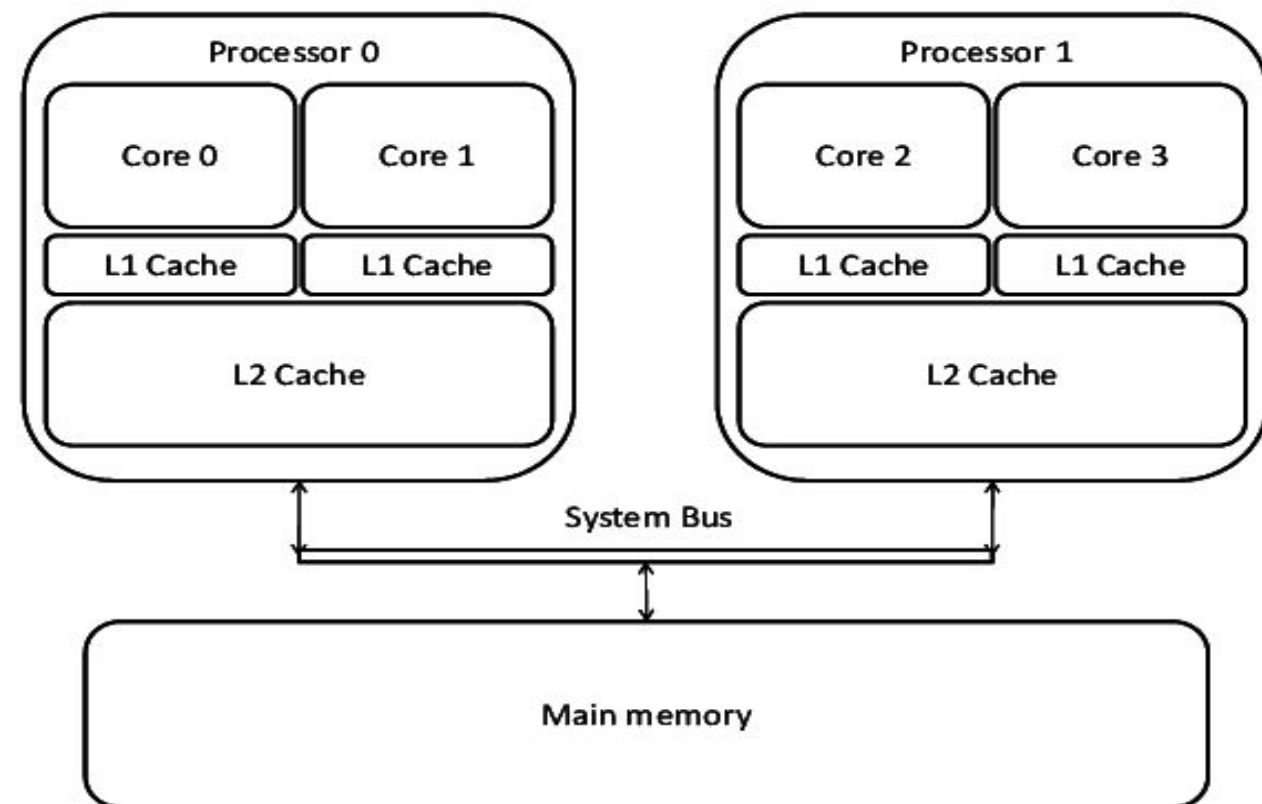
- **Responsiveness** – may allow continued execution if part of process is blocked, especially important for user interfaces
- **Resource Sharing** – threads share resources of process, easier than shared memory or message passing
- **Economy** – cheaper than process creation, thread switching lower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

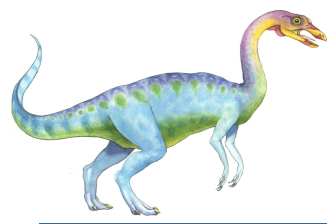




Multicore Programming

- In most recent design whether the cores appear across CPU chips or within CPU chips, we call these systems multicore or multiprocessor systems.
- Multithreaded programming provides a mechanism for more efficient use of these multiple computing cores and improve concurrency.

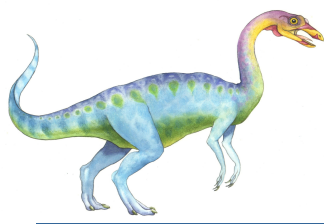




Multicore Programming

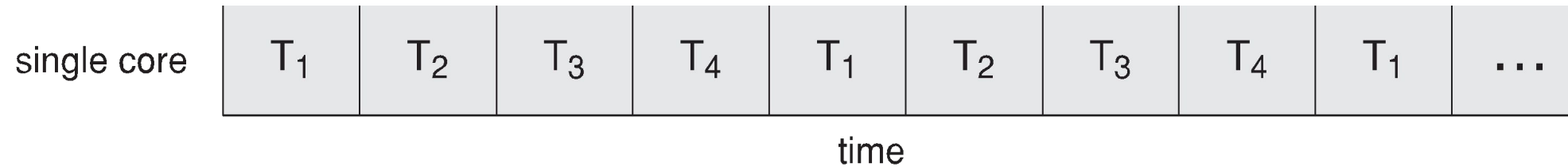
- **Multicore** or **multiprocessor** systems putting pressure on programmers, challenges include:
 - **Dividing activities**
 - **Balance**
 - **Data splitting**
 - **Data dependency**
 - **Testing and debugging**



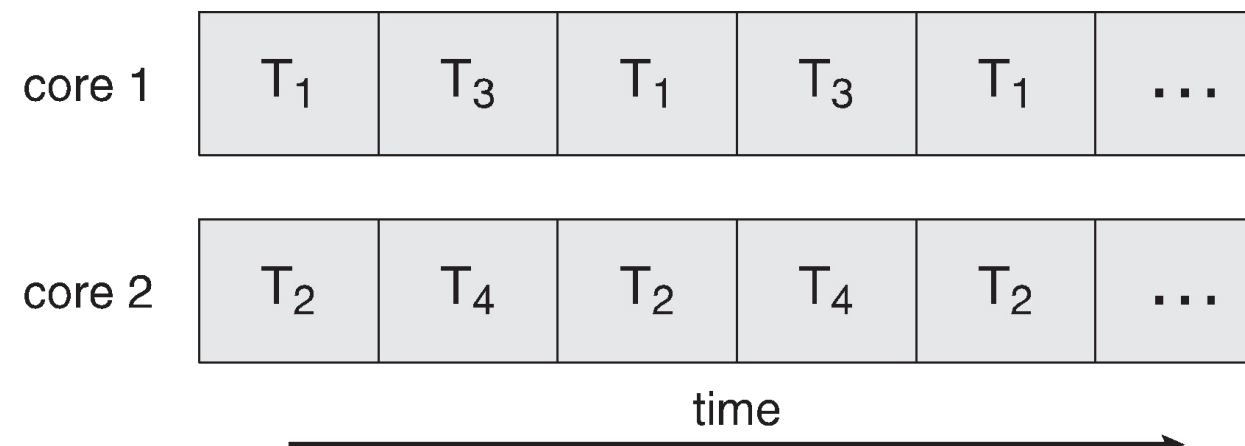


Concurrency vs. Parallelism

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

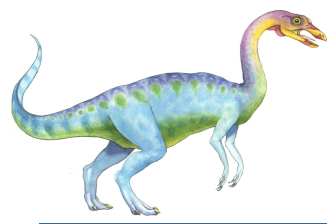


- **Parallelism on a multi-core system:**



- **Parallelism** implies a system can perform more than one task simultaneously
 - Running multiple threads/processes in parallel over different CPU cores
- **Concurrency** supports more than one task making progress
 - Single processor / core, scheduler providing concurrency
 - Running multiple threads/processes at the same time, even on single CPU core, by interleaving their execution.





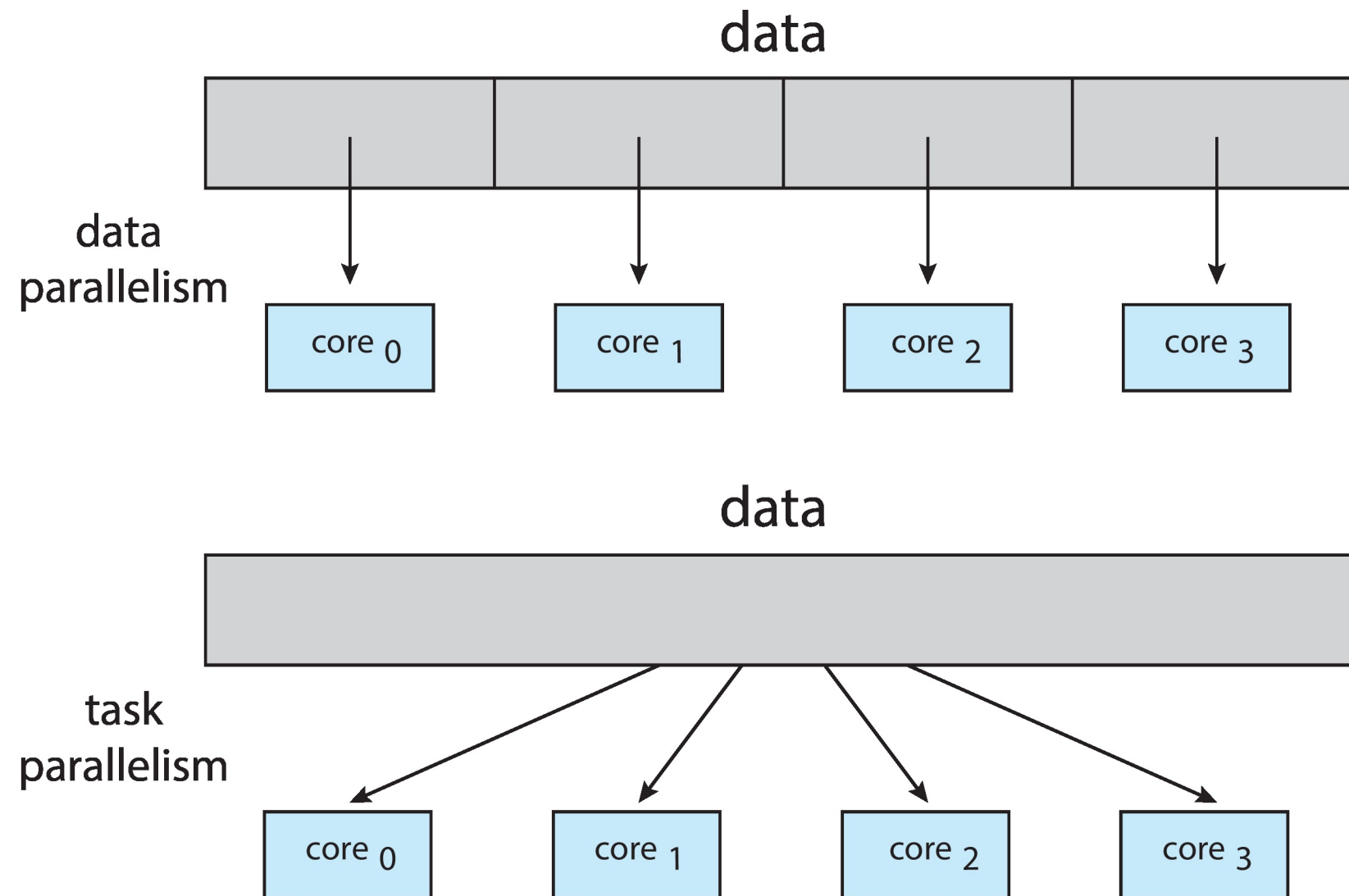
Multicore Programming (Cont.)

- Two Types of parallelism
 - **Data parallelism** – data parallelism is achieved when each processor performs the same task on different pieces of distributed data.
 - **Data parallelism** is a form of **parallel computing** for multiple **processors** using a **technique** for distributing the data across different parallel processor nodes.
 - Synchronous computation is performed.
 - **Task parallelism** – distributing threads across cores, each thread performing unique operation
 - Task Parallelism means concurrent execution of the different task on multiple computing cores on the same or different data.
 - Asynchronous computation is performed.





Data and Task Parallelism





Data vs Task Parallelism

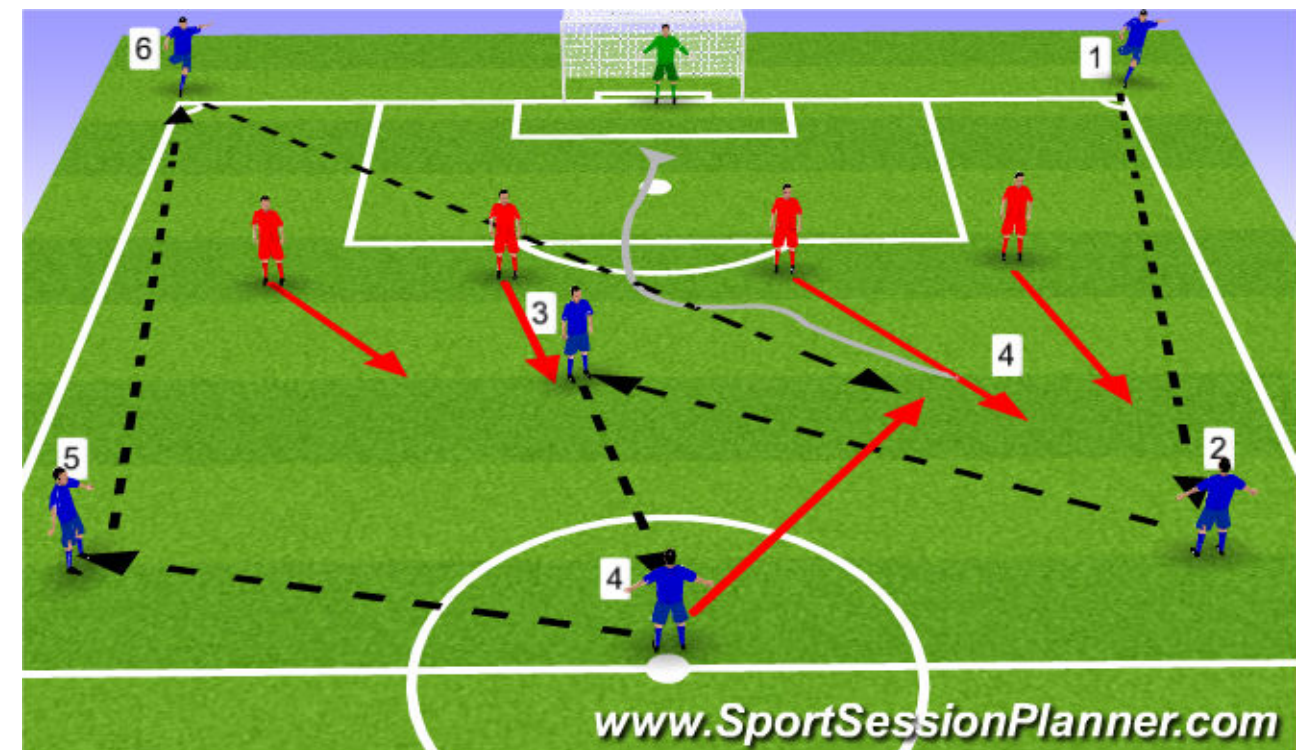
Data Parallelism

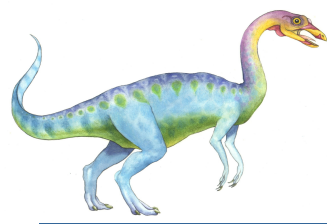
- Single Instruction Multiple data
- Multiple workers all doing same things



Task Parallelism

- Multiple instruction same data or Multiple Instruction Multiple Data
- Works with same objective





Amdahl's Law

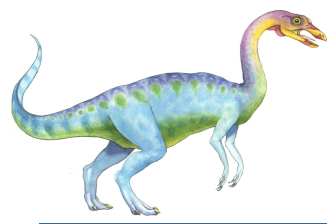
- Identifies performance gains from adding additional cores to an application that has both serial and parallel components

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

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

- That is, if an application is 75% parallel and 25% serial with 2 cores, result speedup by 1.6 times





Amdahl's Law

- Amdahl's Law is mainly used to predict the theoretical maximum speed-up for program processing using multiple-processors
- It is a formula used to find the maximum improvement possible by just improving a particular part of a system
- It is used in parallel computing to predict the theoretical speed-up when using multiple-processor





Amdahl's Law

- Speedup = the time taken for one processor to complete the task / the time taken for 'n' parallel processors to complete the same task

$$\text{Speed up} = \frac{T_1}{T_n} = \frac{1}{S + \frac{P}{n}}$$

- So, in Amdahl's law factors affecting the speedup:
 1. Time taken for serial execution (S)
 2. Time taken for parallel execution (P)
 3. Number of processors (n)





Amdahl's Law

■ If we have one processor to execute a process than time taken for execution is

□ $T_1 = S + P$ [if $n = 1$]

□ Let, execution time $T_1=1$

□ So, $1 = S + P$

□ $P = 1 - S$

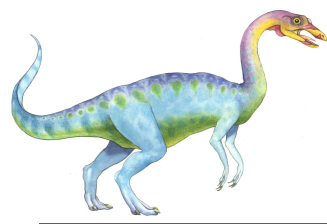
■ If we have n processors to execute a process than time taken for execution is

□ $T_n = S + \frac{P}{n}$

□ $T_n = S + \frac{1-S}{n}$

$$Speed\ up = \frac{T_1}{T_n} = \frac{1}{S + \frac{P}{n}}$$



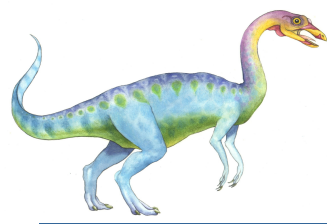




Who Manages Threads?

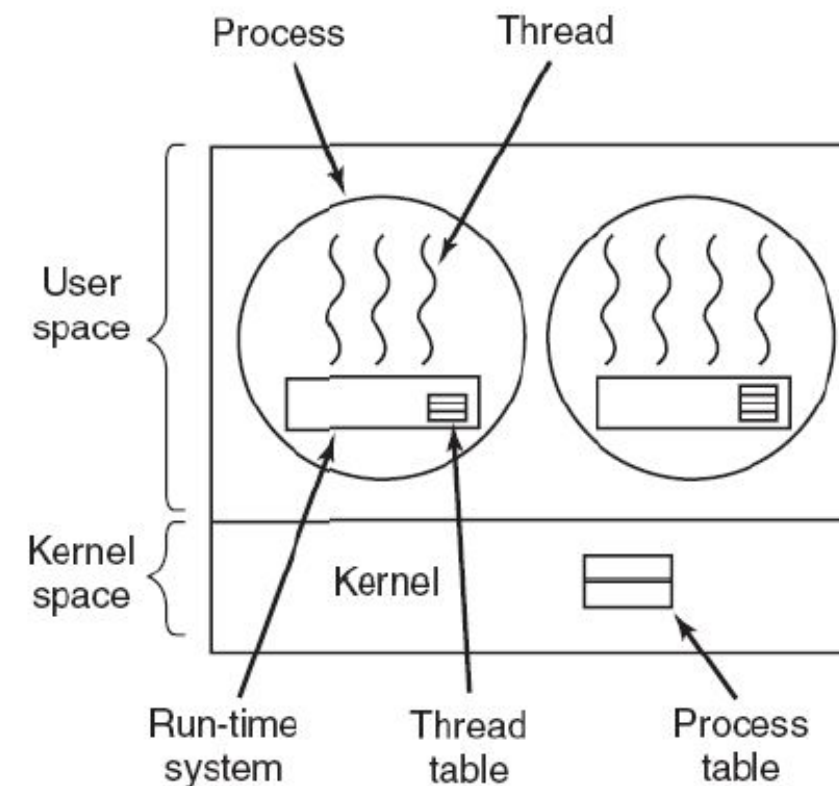
- User Threads
 - Thread management done by user-level threads library
 - User Threads are supported above the Kernel and are managed without kernel support.
 - Kernel knows nothing about User thread
- Kernel Threads
 - Thread directed supported and managed by the OS
 - Many OS kernels are now multi-threaded; several threads operate in the kernel and each thread performs a specific task

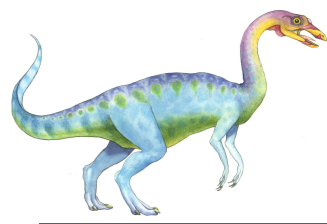




User Level Threads

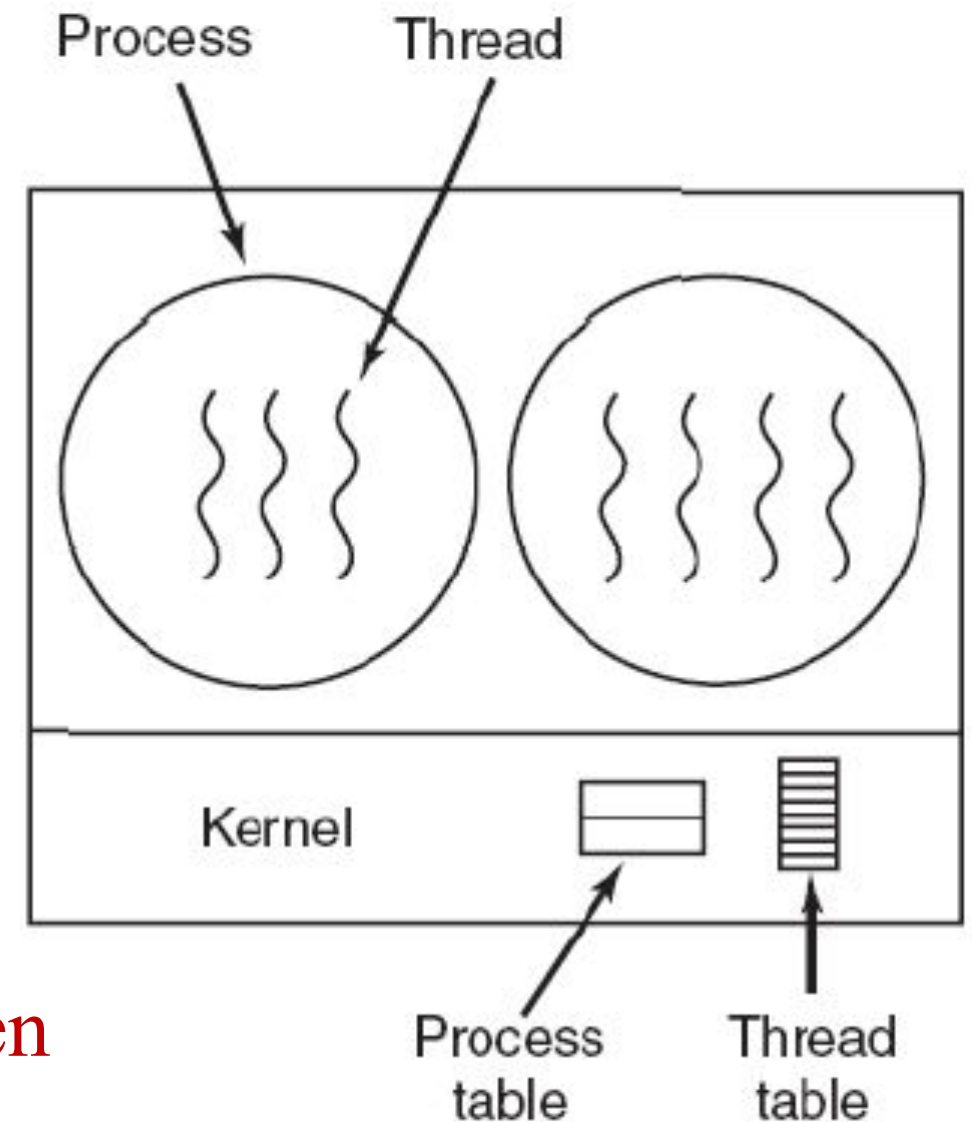
- No system call to manage threads. Thread library does everything.
- Switching is fast, No switching from user to protected mode.
- Can be implemented in an OS that does not support threading.
- **Scheduling can be issue**
- **If one user level thread perform blocking operation then entire process will be blocked.**
- **Lack of coordination between kernel and user level thread.**





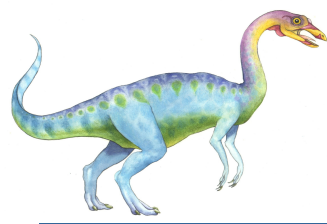
Kernel level Threads

- Scheduler can decide to give more time to a process have large number of threads than process having less number of threads.
- **If one kernel thread perform blocking operation then another thread can continue execution.**



There must exist a relationship between user thread and kernel thread.





Multithreading Models

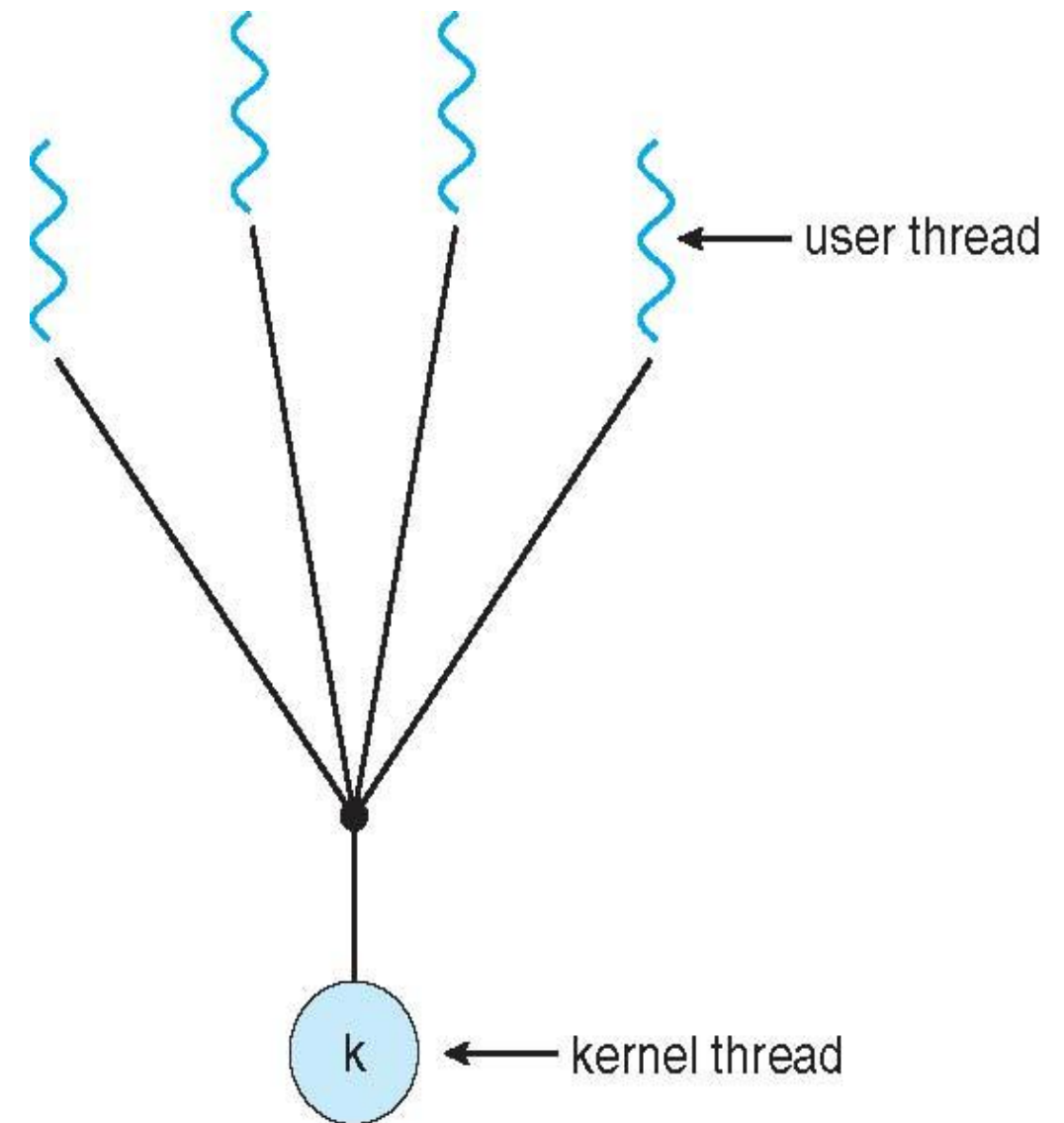
- ❑ **Many-to-One**
- ❑ **One-to-One**
- ❑ **Many-to-Many**





Many-to-One Model

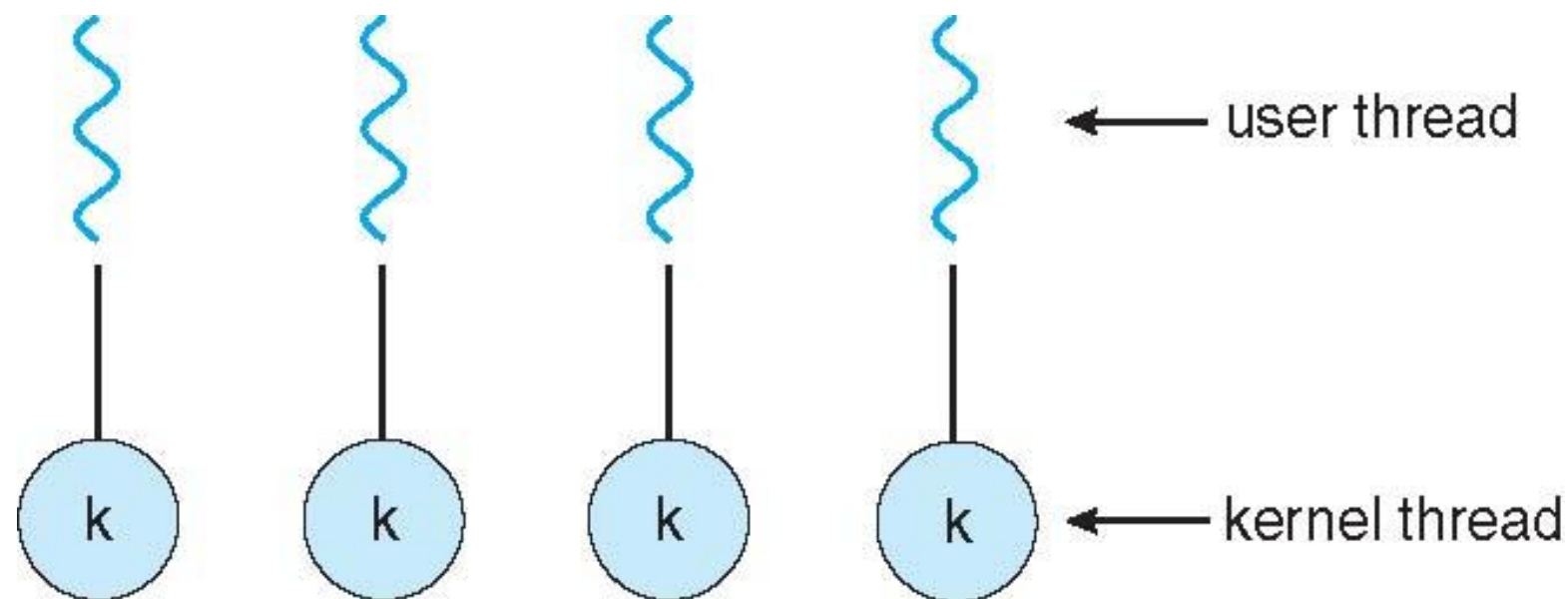
- Many user-level threads mapped to single kernel thread.
- Thread management is done by the thread library in user space, so it is efficient.
- If one thread makes a blocking system call causes the entire process block.
- Because only one thread can access the kernel at a time, multiple threads are unable to run in parallel on multiprocessors.
- Few systems currently use this model
- Examples:
 - **Solaris Green Threads**
 - **GNU Portable Threads**





One-to-one Model

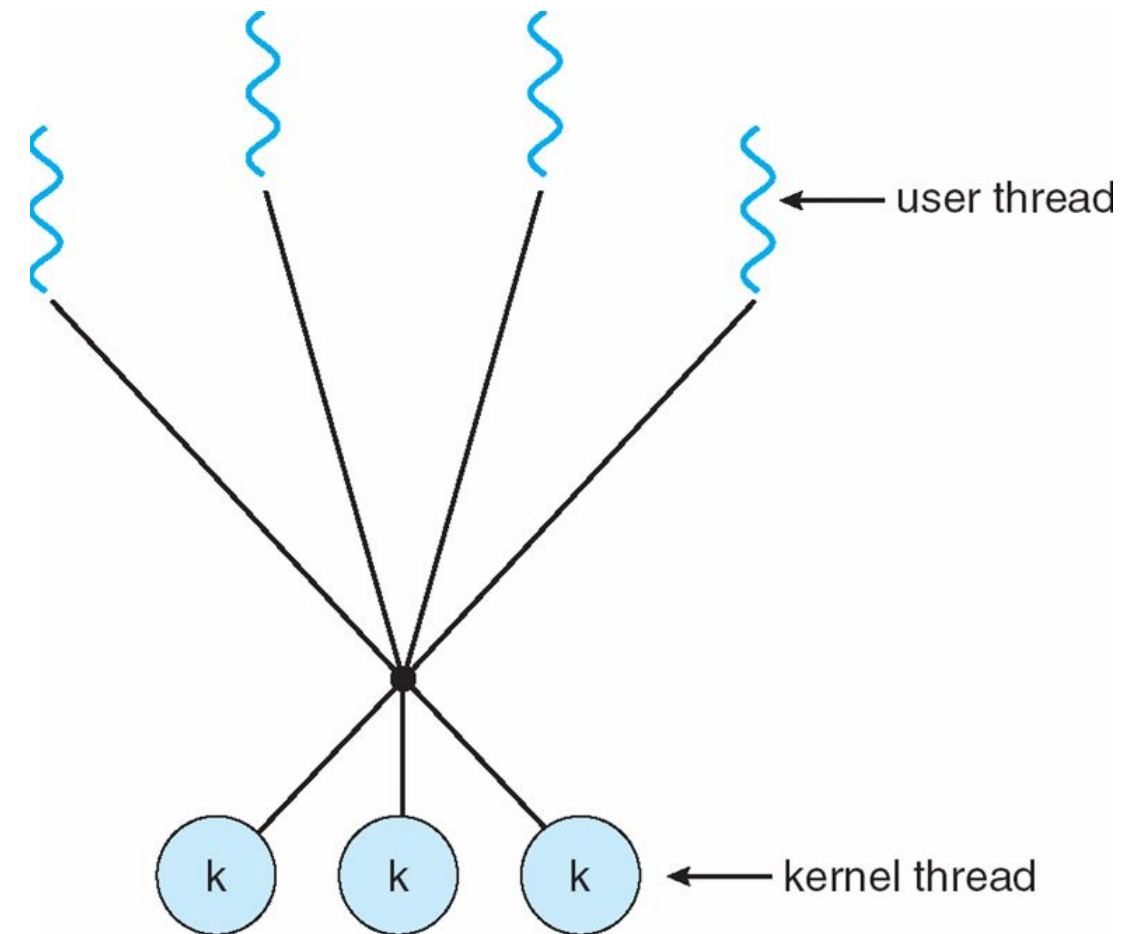
- Each user-level thread maps to kernel thread
- Creating a user-level thread creates a kernel thread
- Provides more concurrency than the many to one model by allowing another thread to run when a thread makes a blocking system call
- Also allows multiple-threads to run in parallel on multi-processors
- Creating a user thread requires creating the corresponding kernel thread
- Because the overhead of creating Kernel threads can Barden the performance of an application, most implementations of this model restrict the number of threads supported by the system.
- Examples: Windows, Linux, Solaris 9 and later

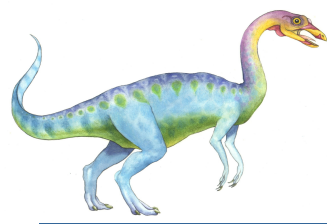




Many-to-Many Model

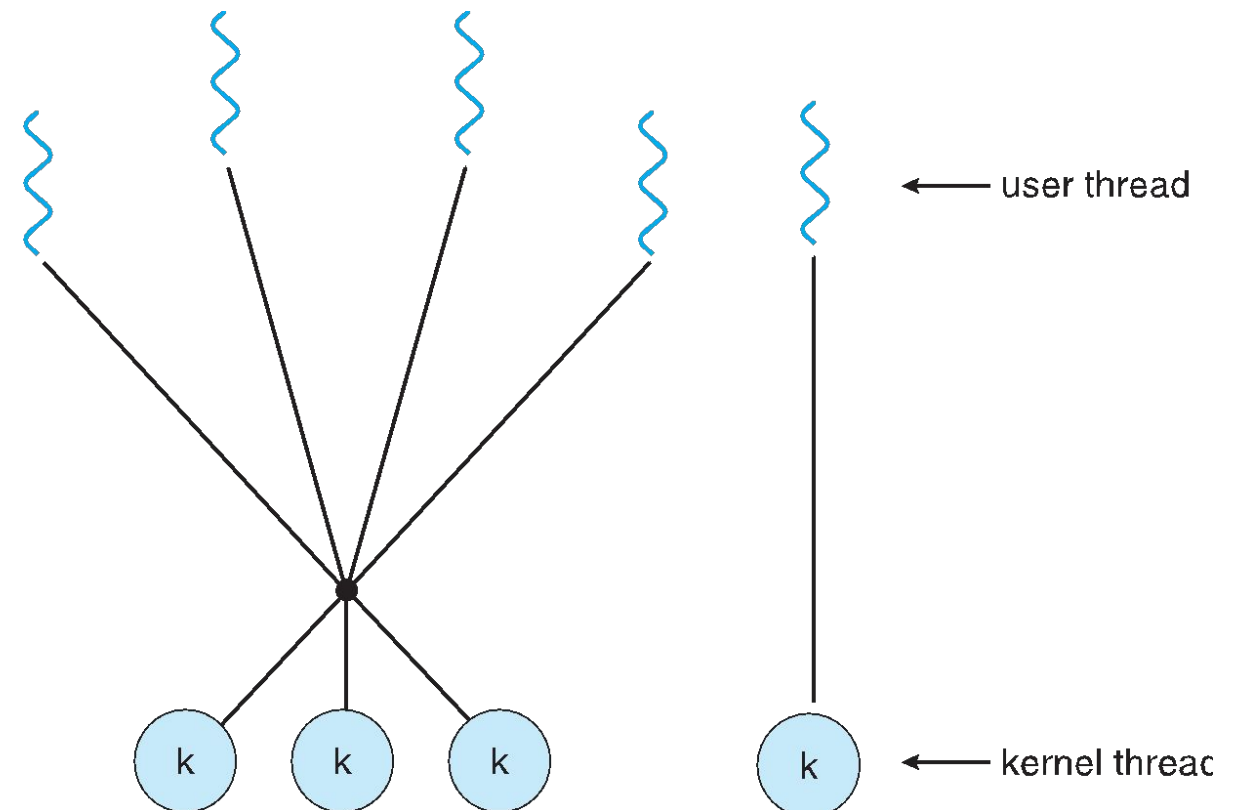
- Multiplexes many user level threads to a smaller or equal number of kernel threads
- Developers can create as many user threads as necessary, and the corresponding kernel threads can run in parallel on a multi-processor
- Also when a thread performs a blocking system call, the kernel can schedule another thread for execution
- The number of kernel threads may be specific to either a particular application or a particular machine

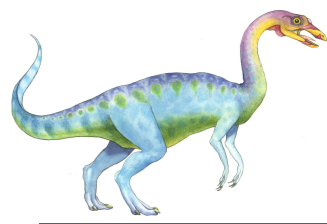




Two-level Model

- 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

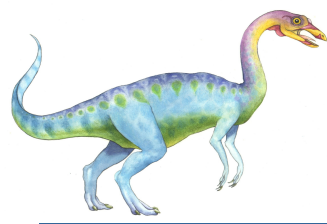




Thread Libraries

- **Thread library** provides programmer with API for creating and managing threads
- Two primary ways of implementing a thread library
 - **First approach** is to implement a user level library with no kernel support. All code and data structures for the library (threads) exist in user space.
 - 4 This means that invoking a function in a library results in a local function calling user space that is not a system call (library entirely in user space) .
 - **Second approach** is to implement a kernel level library supported directly by the OS. In this case code and data structures for the library (threads) exist in kernel space (library supported by OS kernel level) .
 - 4 This means that invoking a function in the library results in a system call to the kernel.
- Three main thread libraries are in use today
 - POSIX Pthreads
 - Win32
 - JAVA

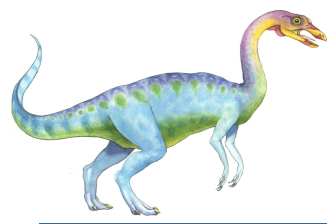




Thread Libraries

- Pthreads
 - Pthreads, is the extension of the POSIX standard may be provided either as user-level or kernel-level library, which defining an API for thread creation and synchronization.
 - This is a specification for thread behavior not an implementation. It is up to developer or OS designer of the library.
 - Common in UNIX operating systems (Solaris, Linux, Mac OS X)
- Win32
 - Win32 is a kernel level library available in windows system
 - These threads are created in the Win32 API using the ***Create Thread()***
 - Just as Pthreads, a set of attributes or parameters is passed to this function





Thread Libraries

- Java Threads

- Java threads API (library) allow JAVA programs directly to create and manage the threads
- Typically implemented using the threads model provided by underlying OS.
- Java threads are managed by the JVM
- Even in single JAVA program which contains only a *main()* method runs as a single thread as the JVM
- There are two techniques for creating threads in JAVA program
 - 4 One approach is to create a new class that is derived from the thread class and to override its *run()* method
 - 4 An alternative and more commonly used technique is to define a class that implements the *Runnable interface*. The runnable interface is defined as follows:

```
public interface Runnable
{
    public abstract void run();
}
```



pthread library

- Create a thread in a process

```
int pthread_create(pthread_t *thread,  
                  const pthread_attr_t *attr,  
                  void *(*start_routine) (void *),  
                  void *arg);
```

Thread identifier (TID) much like

Pointer to a function,
which starts execution in a
different thread

Arguments to the function

- Destroying a thread

```
void pthread_exit(void *retval);
```

Exit value of the thread

- Join : Wait for a specific thread to complete

```
int pthread_join(pthread_t thread, void **retval);
```

TID of the thread to wait for

Exit status of the thread

Example

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

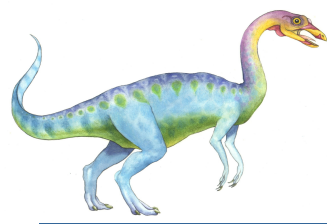
unsigned long sum[4];

void *thread_fn(void *arg){
    long id = (long) arg;
    int start = id * 2500000;
    int i=0;

    while(i < 2500000){
        sum[id] += (i + start);
        i++;
    }
    return NULL;
}

int main(){
    pthread_t t1, t2, t3, t4;

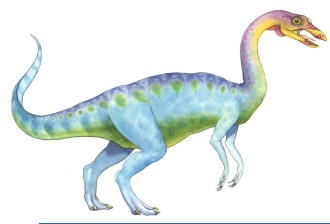
    pthread_create(&t1, NULL, thread_fn, (void *)0);
    pthread_create(&t2, NULL, thread_fn, (void *)1);
    pthread_create(&t3, NULL, thread_fn, (void *)2);
    pthread_create(&t4, NULL, thread_fn, (void *)3);
    pthread_join(t1, NULL);
    pthread_join(t2, NULL);
    pthread_join(t3, NULL);
    pthread_join(t4, NULL);
    printf("%lu\n", sum[0] + sum[1] + sum[2] + sum[3]);
    return 0;
}
```



Implicit Threading

- Implicit multi-threading is concurrent execution of multiple threads extracted from single sequential program.
- Growing in popularity as numbers of threads increase, program correctness more difficult with explicit threads
- Creation and management of threads done by compilers and run-time libraries rather than programmers

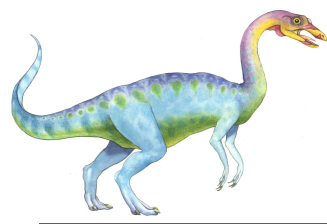




Implicit Threading

- Three methods explored:
 - Thread Pools
 - OpenMP
 - Grand Central Dispatch

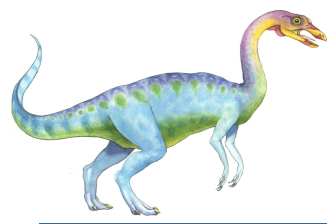




Thread Pools

- If a web server receives a request, it creates a separate thread to serve the request. The issues arise:
 - Thread creation time
 - No bound on concurrent execution of threads
- **Thread pools:**
- Create a number of threads in a pool where they await for work
- Advantages:
 - Usually slightly faster to service a request with an existing thread than create a new thread
 - Allows the number of threads in the application(s) to be bound to the size of the pool





Threading Issues

- ☐ **The fork() and exec() system calls**
- ☐ **Thread Cancellation**
- ☐ **Signal handling**
- ☐ **Thread –Local Storage**





The fork() and exec() system calls

- fork()
 - The fork() system call is used to create a separate, duplicate process
- exec()
 - When a exec() system call is invoked, the program specified in the parameter to exec() will replace the entire process –including all threads
- **Issue:** If one thread in a program calls fork(), does the new process duplicate all threads or the new process single-threaded?
- **Solution:** Some UNIX systems have chosen to have two versions of fork(), one that duplicates all threads and another that duplicates only the thread that invoked the fork() system call

But which version of fork() to use and when?

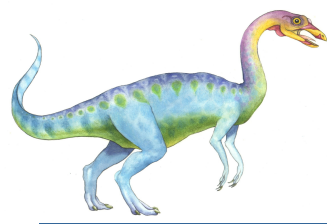




Which version of fork() will to be used depends on the application.

- If exec() is called immediately after forking then duplicating all threads is unnecessary.
- If the separate process does not call exec() after forking, the separate process should duplicate all threads





Thread Cancellation

- Thread cancellation is the task of Terminating a thread before it has completed.
- Thread to be canceled is **target thread**
- Two general approaches:
 - **Asynchronous cancellation:** one thread terminates the target thread immediately.
 - 4 There is an issue if a thread is cancelled while in the midst of updating data it is sharing with other threads.
 - **Deferred cancellation** allows the target thread to periodically check if it should be cancelled
 - 4 Cancellation occurs only after the target thread has checked a flag to determine if it should be cancelled or not. So it should be canceled at a point when it can be cancelled safely.

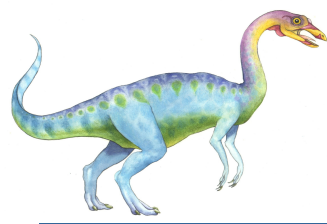




Signal handling

- **Signals** are used in UNIX systems to notify a process that a particular event has occurred.
- Synchronous and asynchronous received signal.
- A **signal handler** is used to process signals
 1. A signal is generated by the occurrence of a particular event
 2. A generated signal is delivered to a process
 3. Once delivered, the signal must be handled
- Every Signal may be handled by one of two possible handlers:
 1. A default signal handler
 2. A user-defined signal handler





Signal handling

- Signals are always delivered to a process but delivering signals is more complicated in multi-threaded program. In that case following options exist
 1. Deliver the signal to the thread to which the signal applies (Synchronous cancellation)
 2. Deliver the signal to every thread in the process (Asynchronous cancellation)
 3. Deliver the signal to certain threads in the process
 4. Assign a specific thread to receive all signals for the process





Thread Local Storage

- ❑ Threads belonging to a process share the data of the process. In some circumstances, each thread might need its own copy of certain data.
- Create Facility needed for data private to thread is called **Thread-local storage (TLS)** or thread-specific data
- Useful when you do not have control over the thread creation process (i.e., when using a thread pool)



End of Chapter 4

NAT

