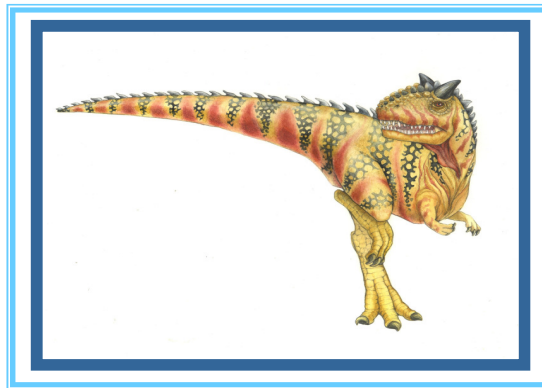


Lecture 3: Thread

From Processes to Threads

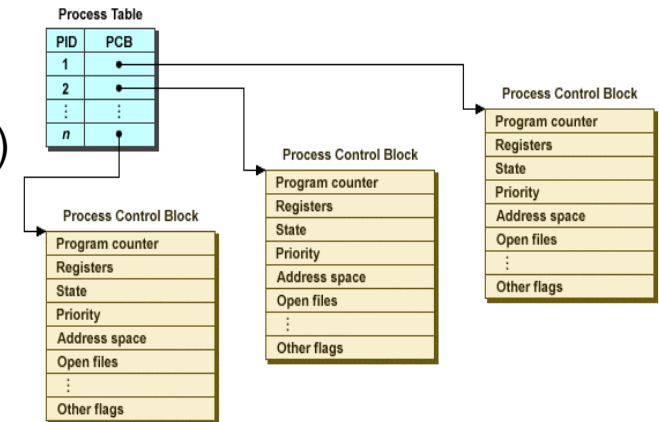




The Soul of a Process

■ What is similar in **cooperating processes**?

- They all share the same code and data (address space)
- They all share the same privileges
- They all share the same resources (files, sockets, etc.)



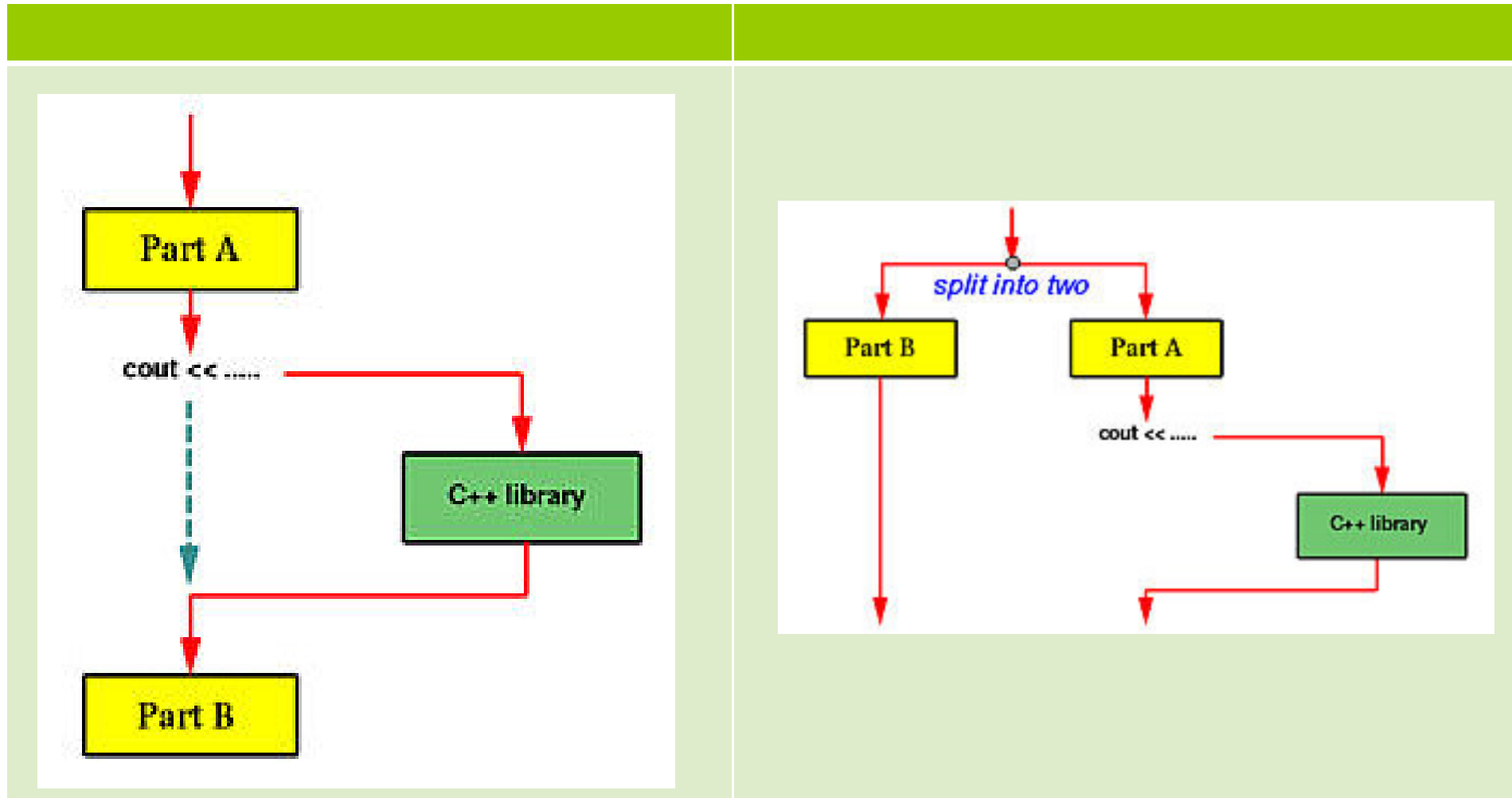
■ What don't they share?

- Each has its own execution state: PC, SP, and registers
- Key idea: **Why don't we separate the concept of a process from its execution state?**
 - ▶ Process: address space, privileges, resources, etc.
 - ▶ Execution state: PC, SP, registers
- Exec state also called **thread of control**, or **thread**





Why Threads?

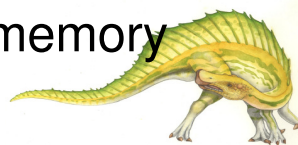




Processes and Threads

- Process abstraction combines two concepts
 - Concurrency
 - ▶ Each process is a - sequential execution stream of instructions
 - Protection
 - ▶ Each process defines an address space
 - ▶ Address space identifies all addresses that can be touched by the program

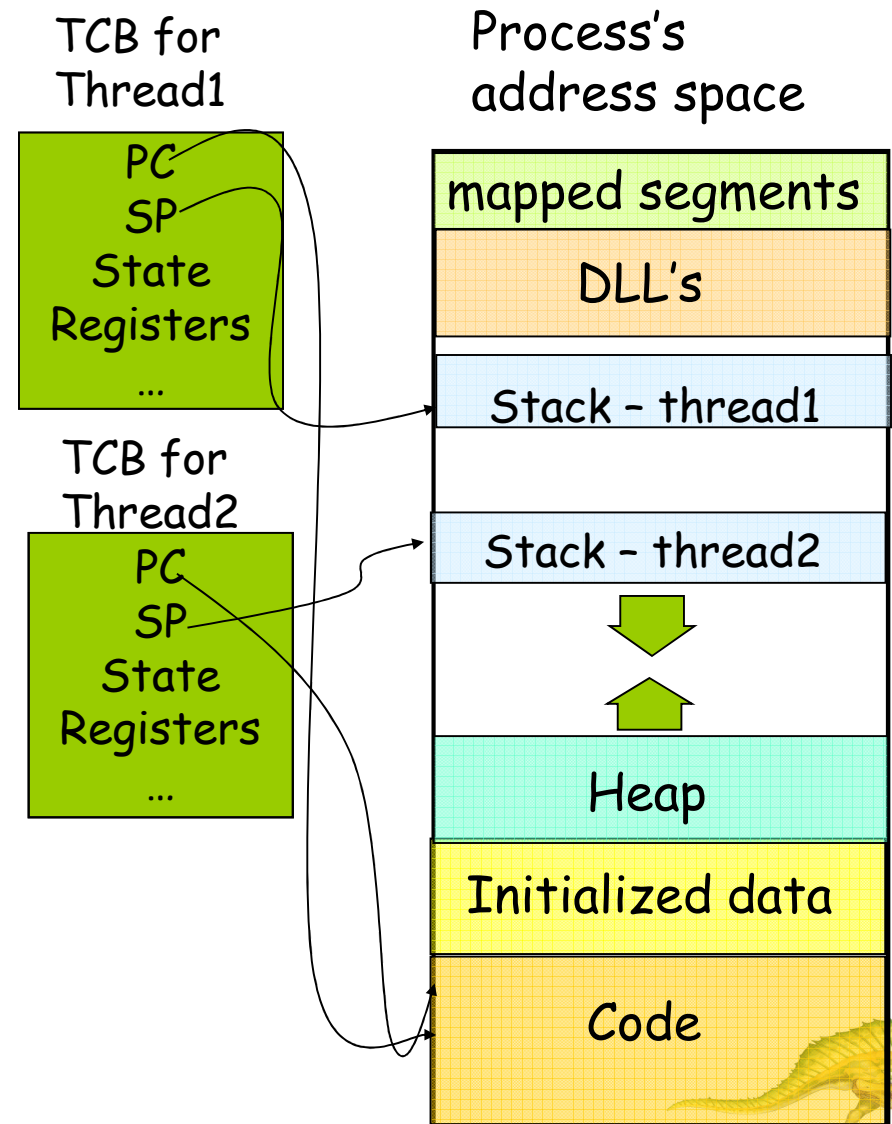
- Threads
 - Key idea: separate the concepts of concurrency from protection
 - A thread is a sequential execution stream of instructions
 - A process defines the address space that may be shared by multiple threads
 - Threads can execute on different cores on a multicore CPU (parallelism for performance) and can communicate with other threads by updating memory





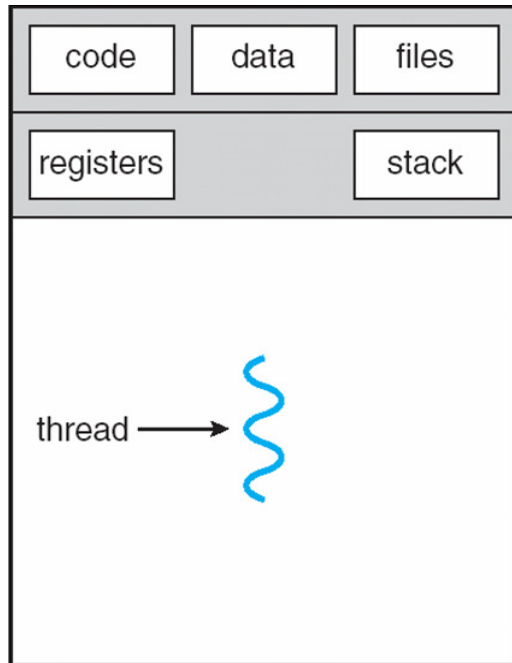
Processes and Threads

- Processes define an address space; threads share the address space
- Process Control Block (PCB) contains process-specific information
 - Owner, PID, heap pointer, priority, active thread, and pointers to thread information
- Thread Control Block (TCB) contains thread-specific information
 - Stack pointer, PC, thread state (running, ...), register values, a pointer to PCB, ...

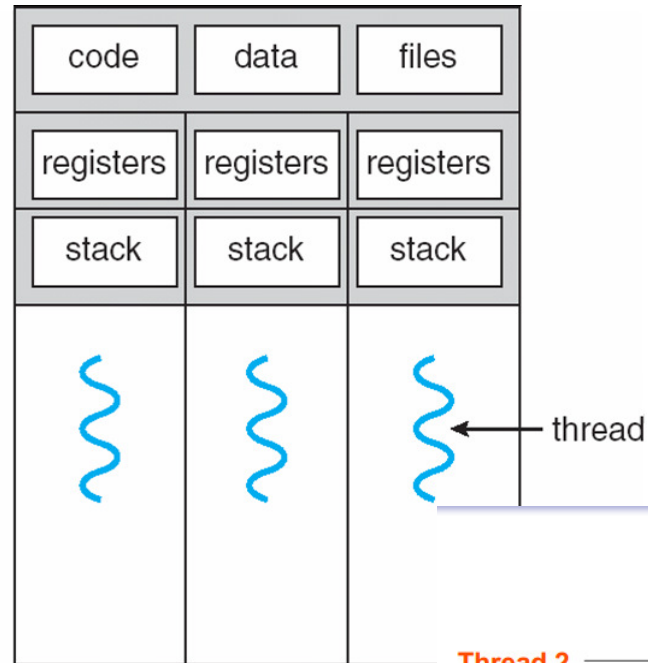




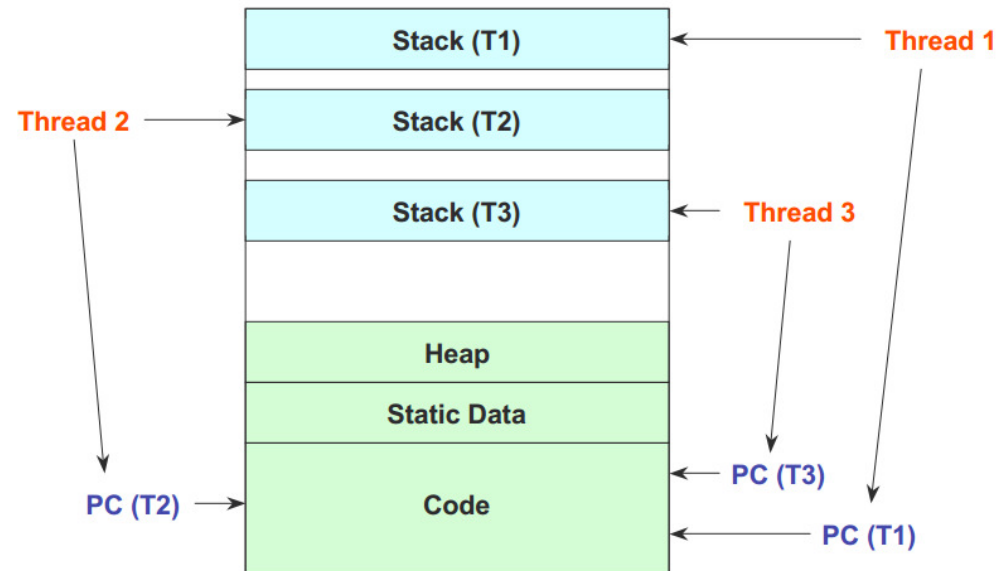
Single and Multithreaded Processes



single-threaded process



multithreaded process





The Case for Threads

Consider the following code fragment

```
for(k = 0; k < n; k++)  
    a[k] = b[k] * c[k] + d[k] * e[k];
```

Is there a missed opportunity here? On a Uni-processor?
On a Multi-processor?





The Case for Threads

Consider a Web server

- get network message (URL) from client

- get URL data from disk

- compose response

- send response

How well does this web server perform?





Introducing Threads

- A thread represents an abstract entity that executes a sequence of instructions
 - It has its own set of CPU registers
 - It has its own stack
 - There is no thread-specific heap or data segment (unlike process)
- Threads are lightweight
 - Creating a thread more efficient than creating a process.
 - Communication between threads easier than processes.
 - Context switching between threads requires fewer CPU cycles and memory references than switching processes.
 - Threads only track a subset of process state (share list of open files, pid, ...)
- Examples:
 - OS-supported: Windows' threads, Sun's LWP, POSIX threads
 - Language-supported: Modula-3, Java





Context switch time for which entity is greater?

1. Process
2. Thread





How Can it Help?

- How can this code take advantage of 2 threads?

```
for(k = 0; k < n; k++)  
    a[k] = b[k] * c[k] + d[k] * e[k];
```

- Rewrite this code fragment as:

```
do_mult(l, m) {  
    for(k = l; k < m; k++)  
        a[k] = b[k] * c[k] + d[k] * e[k];  
}  
  
main() {  
    CreateThread(do_mult, 0, n/2);  
    CreateThread(do_mult, n/2, n);  
}
```

- What did we gain?





How Can it Help?

- Consider a Web server

Create a number of threads, and for each thread do

- ❖ get network message from client
- ❖ get URL data from disk
- ❖ send data over network

- What did we gain?



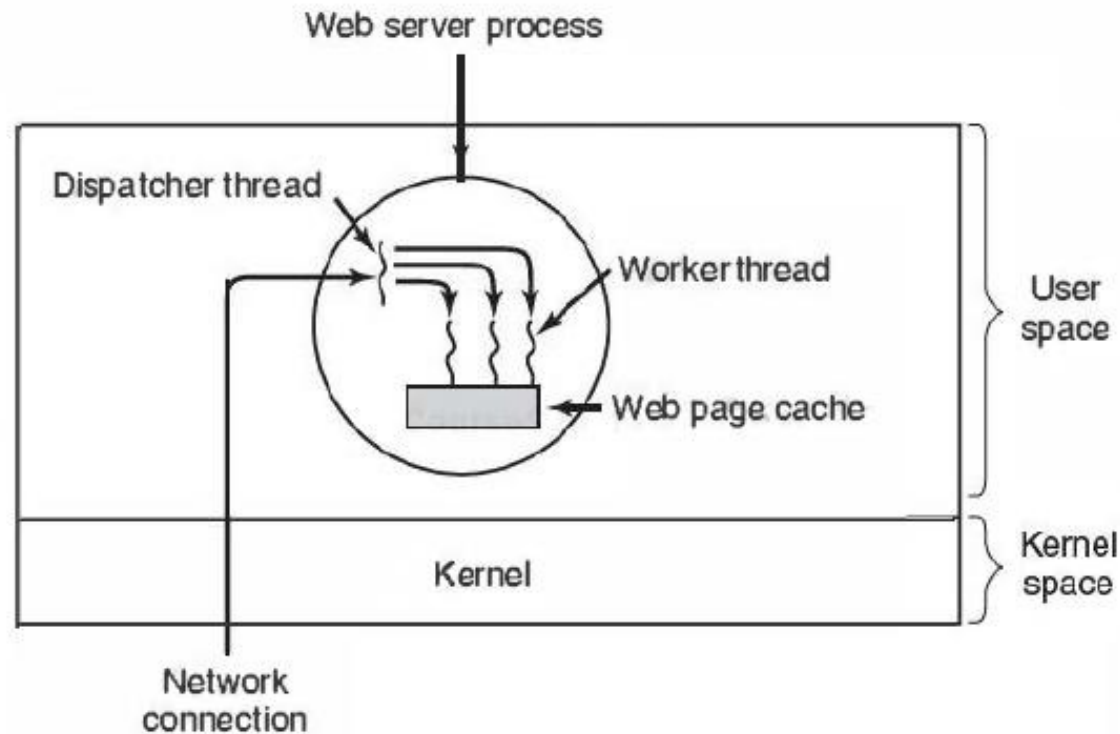


Figure 2. A multithreaded Web server

```
while (TRUE) {  
    get_next_request(&buf);  
    handoff_work(&buf);  
}
```

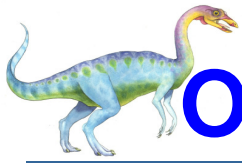
(a)

```
while (TRUE) {  
    wait_for_work(&buf)  
    look_for_page_in_cache(&buf, &page);  
    if (page_not_in_cache(&page))  
        read_page_from_disk(&buf, &page);  
    return_page(&page);  
}
```

(b)

Figure 3. A rough outline of the code for Fig. 2. (a) Dispatcher thread. (b) Worker thread





Overlapping Requests (Concurrency)

Request 1
Thread 1

Request 2
Thread 2

- ❖ get network message (URL) from client
- ❖ get URL data from disk
(disk access latency)

- ❖ get network message (URL) from client
- ❖ get URL data from disk
(disk access latency)

- ❖ send data over network

- ❖ send data over network

- ◆ Total time is less than request 1 + request 2

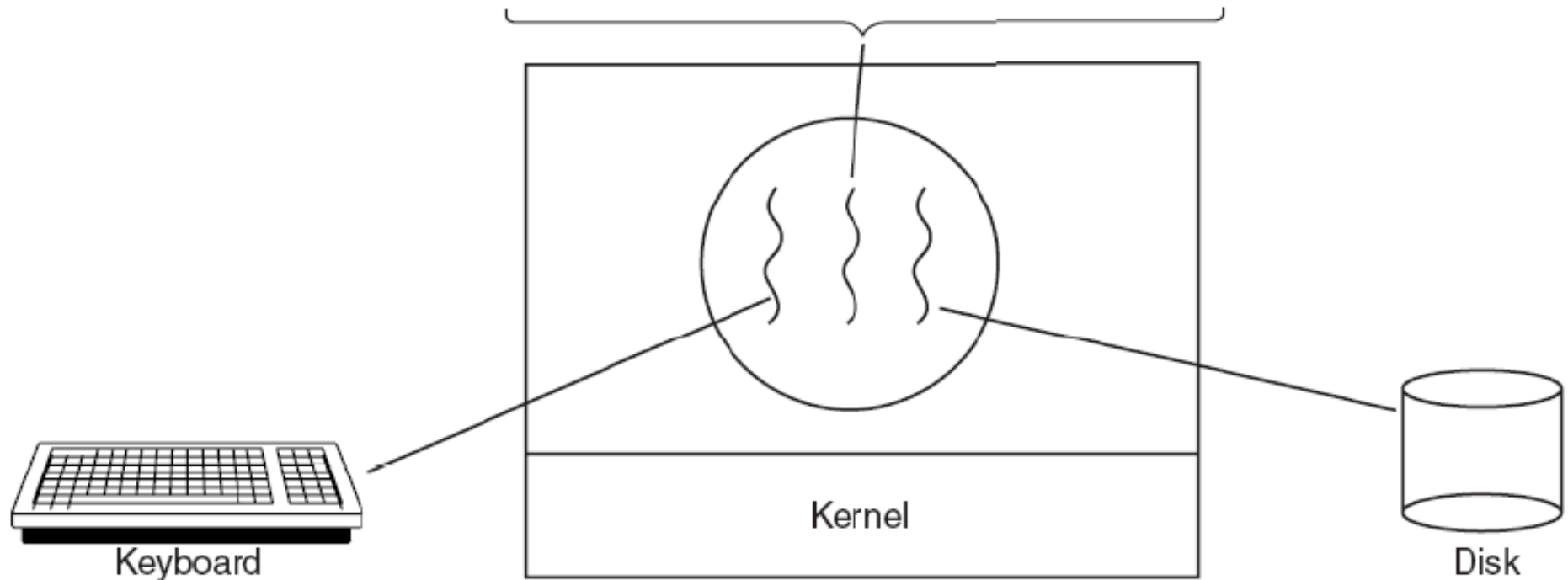
Time





A Word Processor with Three Threads

Four score and seven years ago, our fathers brought forth upon this continent a new nation, conceived in liberty, and dedicated to the proposition that all men are created equal.	ration, or any action so conceived and so dedicated, can long endure. We are met on a great battlefield of this war.	lies that this nation might live. It is altogether fitting and proper that we should do this.	who struggled here have consecrated it, far above our poor power to add or detract. The world will little note, nor long remember, what we say here, but it can never forget what they did here.	here to the unfinished work which they who fought here have thus far so nobly advanced. It is rather for us to be here dedicated to the great task remaining before us, that from these honored dead we take increased devotion to that cause for which	they gave the last full measure of devotion, that we here highly resolve that these dead shall not have died in vain that this nation, under God, shall have a new birth of freedom and that government of the people, by the people
--	--	---	--	---	--





Threads vs. Processes

Threads

- A thread has no data segment or heap
- A thread cannot live on its own, it must live within a process
- There can be more than one thread in a process, the first thread calls main & has the process's stack
- If a thread dies, its stack is reclaimed
- Inter-thread communication via memory.
- Each thread can run on a different physical processor
- Inexpensive creation and context switch

Processes

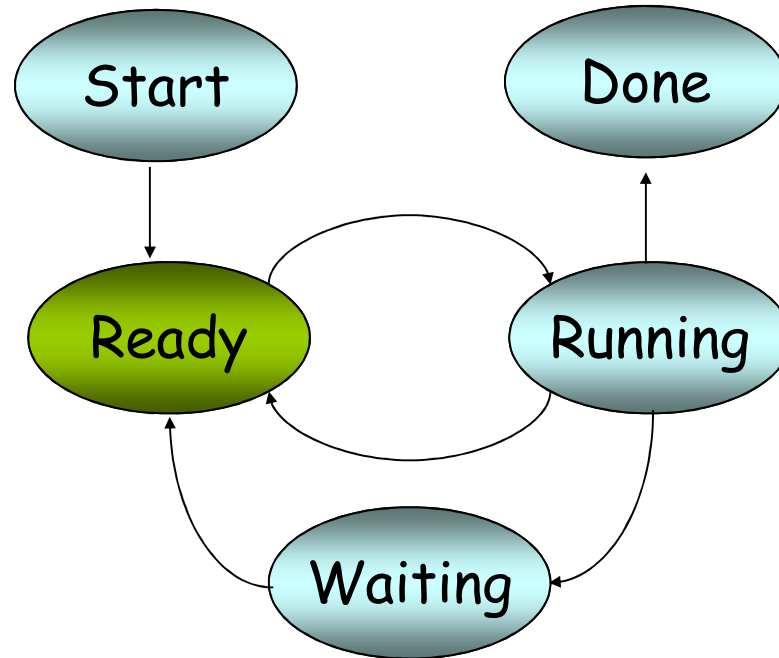
- ◆ A process has code/data/heap & other segments
- ◆ There must be at least one thread in a process
- ◆ Threads within a process share code/data/heap, share I/O, but each has its own stack & registers
- ◆ If a process dies, its resources are reclaimed & all threads die
- ◆ Inter-process communication via OS and data copying.
- ◆ Each process can run on a different physical processor
- ◆ Expensive creation and context switch

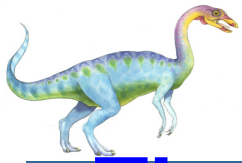




Threads' Life Cycle

- Threads (just like processes) go through a sequence of *start*, *ready*, *running*, *waiting*, and *done* states





Threads have the same scheduling states as processes

1. True
2. False



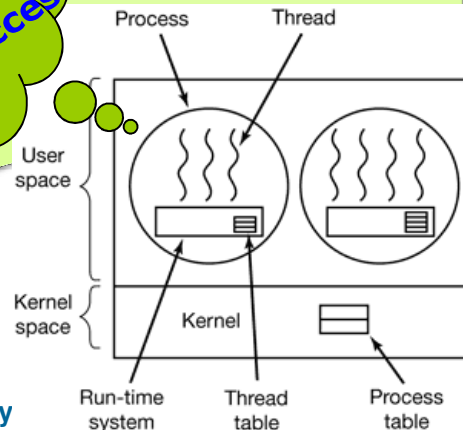


Implementing Threads

User Level

- Threads package entirely in user space
- Kernel knows nothing about threads
- Fast to create and switch-scheduler as local procedure

Blocking
System Call
CPU access

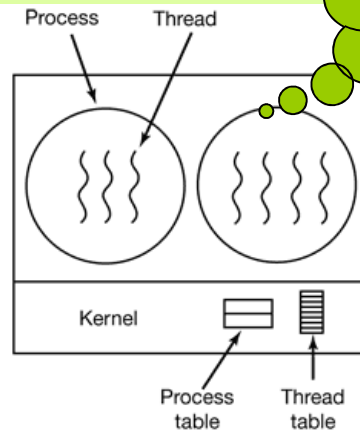


Kernel Level

- No runtime system
- Global Thread table, updated by kernel call
- Do not block process for syscall
- Thread switching: same or another process
- Not so fast as runtime system

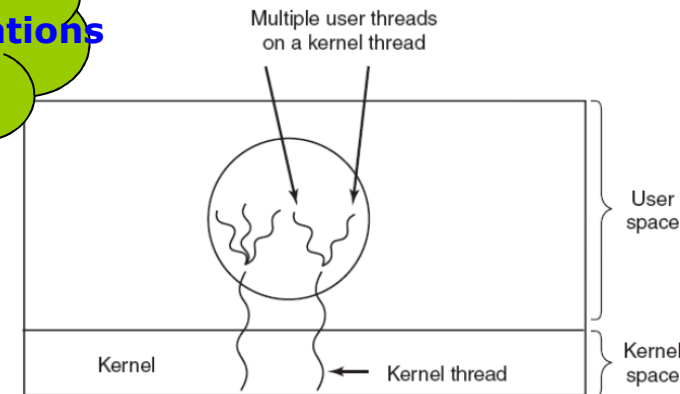
Thread
Recycling

Frequent
thread operations
Need many
System calls



Hybrid

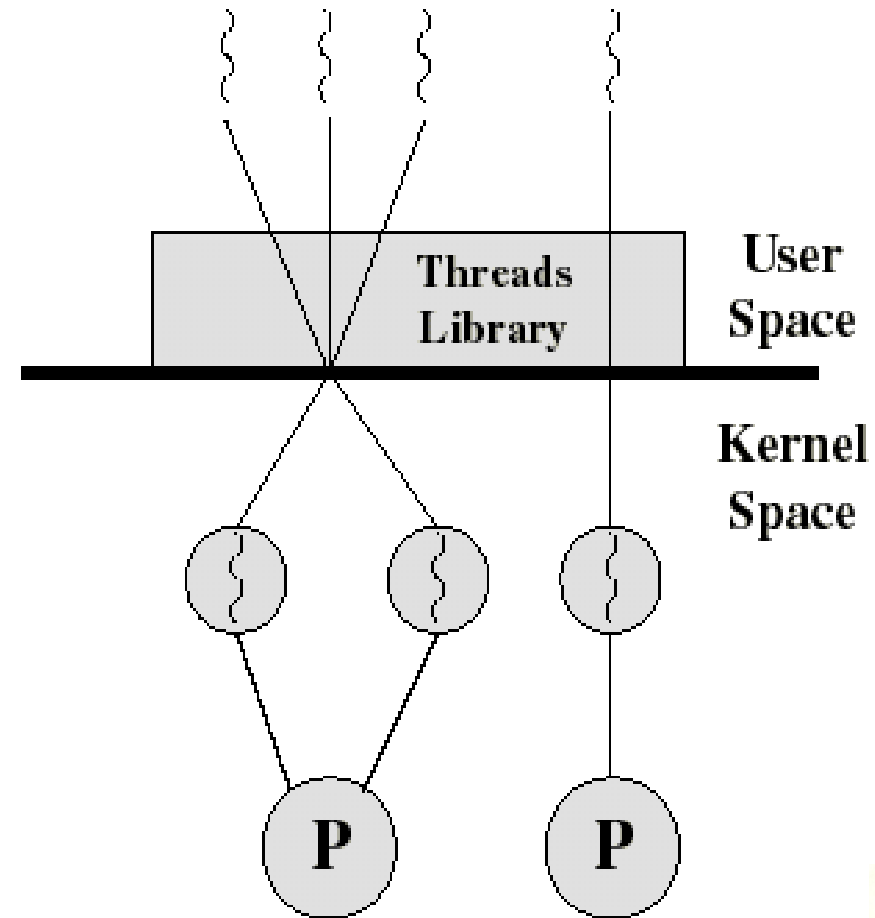
- Use kernel-level threads and then multiplex user-level threads onto same or all kernel threads.





Hybrid ULT/KLT Approaches

- Thread creation done in the user space.
- Scheduling and synchronization of threads done in the user space.
- The programmer may adjust the number of KLTs.
- May combine the best of both approaches.



From- A. Frank - P. Weisberg
3.20





POSIX Thread Program-1

```
#include <pthread.h>
#include <stdio.h>
#define NUM_THREADS      5

void *PrintHello(void *threadid)
{
    long tid;
    tid = (long)threadid;
    printf("Hello World! It's me, thread #%ld!\n", tid);
    pthread_exit(NULL);
}
```

```
int main (int argc, char *argv[])
{
    pthread_t threads[NUM_THREADS];
    int rc;
    long t;
    for(t=0; t<NUM_THREADS; t++){
        printf("In main: creating thread %ld\n", t);
        rc = pthread_create(&threads[t], NULL, PrintHello, (void *)t);
        if (rc){
            printf("ERROR: return code from pthread_create() is %d\n", rc);
            exit(-1);
        }
    }

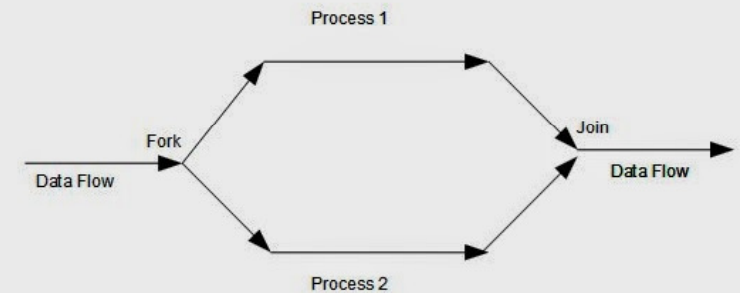
    /* Last thing that main() should do */
    pthread_exit(NULL);
}
```

Opaque Object



POSIX Thread Program-2

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <pthread.h>
4
5 /* Global variable: accessible to all threads */
6 int thread_count;
7
8 void* Hello(void* rank); /* Thread function */
9
10 int main(int argc, char* argv[]) {
11     long thread; /* Use long in case of a 64-bit system */
12     pthread_t* thread_handles;
13
14     /* Get number of threads from command line */
15     thread_count = strtol(argv[1], NULL, 10);
16
17     thread_handles = malloc (thread_count*sizeof(pthread_t));
18
19     for (thread = 0; thread < thread_count; thread++)
20         pthread_create(&thread_handles[thread], NULL,
21             Hello, (void*) thread);
22
23     printf("Hello from the main thread\n");
24
25     for (thread = 0; thread < thread_count; thread++)
26         pthread_join(thread_handles[thread], NULL);
27
28     free(thread_handles);
29     return 0;
30 } /* main */
31
32 void* Hello(void* rank) {
33     long my_rank = (long) rank
34         /* Use long in case of 64-bit system */
35
36     printf("Hello from thread %ld of %d\n", my_rank,
37         thread_count);
38
39     return NULL;
40 } /* Hello */
```



```
gcc -o pth_hello pth_hello.c -lpthread
./pth_hello 2
```



End of Lecture 3

