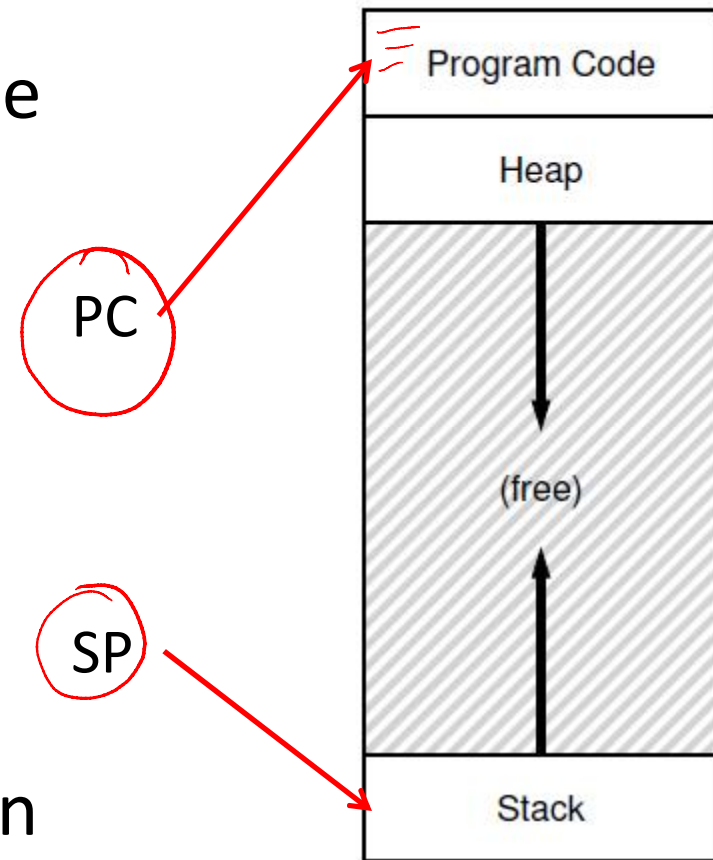


Lecture 12: Threads and Concurrency

Mythili Vutukuru
IIT Bombay

Single threaded process

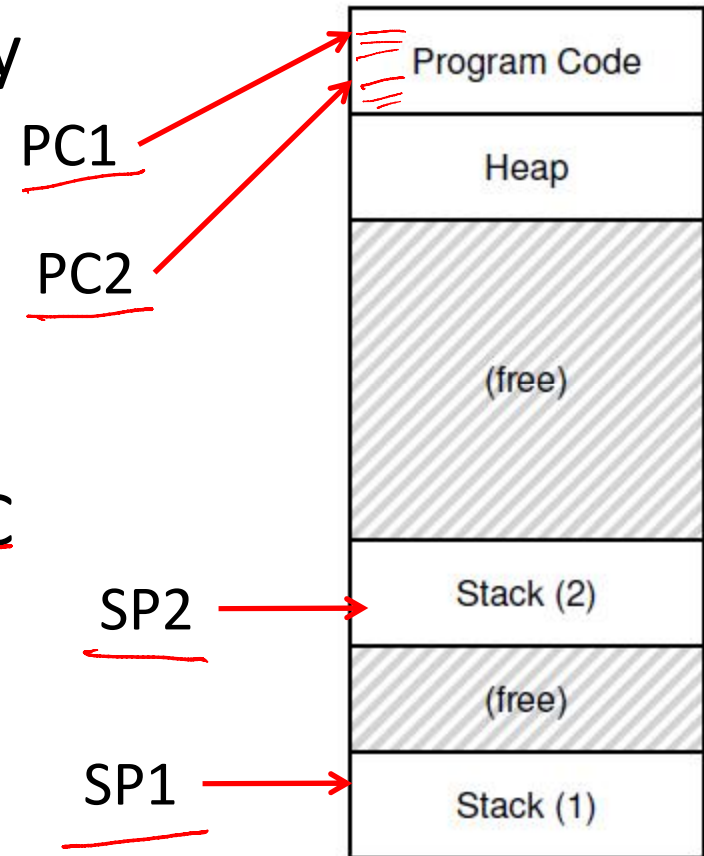
- So, far we have studied single threaded programs
- Recap: process execution
 - PC points to current instruction being run
 - SP points to stack frame of current function call
- A program can also have multiple threads of execution
- What is a thread?



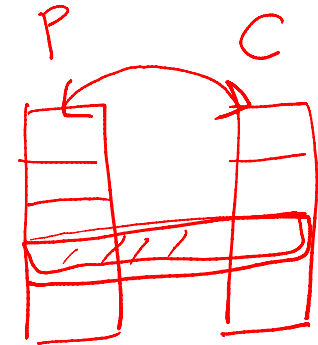
mem

Multi threaded process

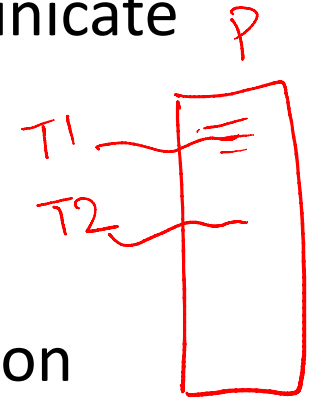
- A thread is like another copy of a process that executes independently
- Threads shares the same address space (code, heap)
- Each thread has separate PC
 - Each thread may run over different part of the program
- Each thread has separate stack for independent function calls



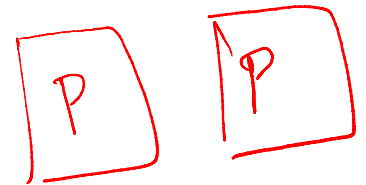
Process vs. threads



- Parent P forks a child C
 - P and C do not share any memory
 - Need complicated IPC mechanisms to communicate
 - Extra copies of code, data in memory
- Parent P executes two threads T1 and T2
 - T1 and T2 share parts of the address space
 - Global variables can be used for communication
 - Smaller memory footprint
- Threads are like separate processes, except they share the same address space



Why threads?

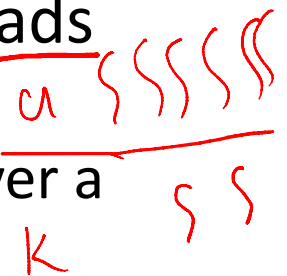
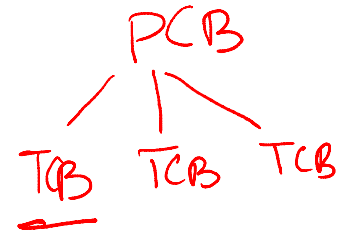


- Parallelism: a single process can effectively utilize multiple CPU cores
 - Understand the difference between concurrency and parallelism
 - Concurrency: running multiple threads/processes at the same time, even on single CPU core, by interleaving their executions
 - Parallelism: running multiple threads/processes in parallel over different CPU cores
- Even if no parallelism, concurrency of threads ensures effective use of CPU when one of the threads blocks (e.g., for I/O)



Scheduling threads

- OS schedules threads that are ready to run independently, much like processes
- The context of a thread (PC, registers) is saved into/restored from thread control block (TCB)
 - Every PCB has one or more linked TCBs
- Threads that are scheduled independently by kernel are called kernel threads
 - E.g., Linux pthreads are kernel threads
- In contrast, some libraries provide user-level threads
 - User program sees multiple threads
 - Library multiplexes larger number of user threads over a smaller number of kernel threads
 - Low overhead of switching between user threads (no expensive context switch)
 - But multiple user threads cannot run in parallel



Creating threads using pthreads API

```
1  #include <stdio.h>
2  #include <assert.h>
3  #include <pthread.h>
4
5  void *mythread(void *arg) {
6      printf("%s\n", (char *) arg);
7      return NULL;
8  }
9
10 int
11 main(int argc, char *argv[]) {
12     pthread_t p1, p2;
13     int rc;
14     printf("main: begin\n");
15     rc = pthread_create(&p1, NULL, mythread, "A"); assert(rc == 0);
16     rc = pthread_create(&p2, NULL, mythread, "B"); assert(rc == 0);
17     // join waits for the threads to finish
18     rc = pthread_join(p1, NULL); assert(rc == 0);
19     rc = pthread_join(p2, NULL); assert(rc == 0);
20     printf("main: end\n");
21     return 0;
22 }
```

Handwritten annotations:

- Red circles around `mythread` and `"A"` in line 15.
- Red circles around `mythread` and `"B"` in line 16.
- Red circles around `pthread_create` in lines 15 and 16.
- Red circles around `pthread_join` in lines 18 and 19.
- Red arrow pointing from the `mythread` parameter in line 15 to the `mythread` function definition in line 5.
- Handwritten text: `P`, `T1`, `T2`, `A`, `B`, `end`.

Figure 26.2: Simple Thread Creation Code (t0.c)

Example: threads with shared data

```
4
5  static volatile int counter = 0;
6
7  //
8  // mythread()
9  //
10 // Simply adds 1 to counter repeatedly, in a loop
11 // No, this is not how you would add 10,000,000 to
12 // a counter, but it shows the problem nicely.
13 //
14 void *
15 mythread(void *arg)
16 {
17     printf("%s: begin\n", (char *) arg);
18     int i;
19     for (i = 0; i < 1e7; i++) {
20         counter = counter + 1;
21     }
22     printf("%s: done\n", (char *) arg);
23     return NULL;
24 }
25
26 //
27 // main()
28 //
29 // Just launches two threads (pthread_create)
30 // and then waits for them (pthread_join)
31 //
32 int
33 main(int argc, char *argv[])
34 {
35     pthread_t p1, p2;
36     printf("main: begin (counter = %d)\n", counter);
37     pthread_create(&p1, NULL, mythread, "A");
38     pthread_create(&p2, NULL, mythread, "B");
39
40     // join waits for the threads to finish
41     pthread_join(p1, NULL);
42     pthread_join(p2, NULL);
43     printf("main: done with both (counter = %d)\n", counter);
44     return 0;
45 }
```

T1 T2
{ {
A B

Threads with shared data: what happens?

- What do we expect? Two threads, each increments counter by 10^7 , so 2×10^7

```
prompt> gcc -o main main.c -Wall -pthread
prompt> ./main
main: begin (counter = 0)
A: begin ✓
B: begin —
A: done ✓
B: done —
main: done with both (counter = 20000000)
```

- Sometimes, a lower value. Why?

```
prompt> ./main
main: begin (counter = 0)
A: begin
B: begin
A: done
B: done
main: done with both (counter = 19345221)
```

What is happening?

- Assembly code of
counter = counter + 1

100 mov 0x8049a1c, %eax
105 add \$0x1, %eax
108 mov %eax, 0x8049a1c

counter → *reg*

OS	Thread 1	Thread 2	(after instruction)		
			PC	%eax	counter
	<i>before critical section</i>		100	0	50
	mov 0x8049a1c, %eax		105	50	50
	add \$0x1, %eax		108	51	50
interrupt	save T1's state				
	restore T2's state				
		mov 0x8049a1c, %eax	100	0	50
		add \$0x1, %eax	105	50	50
		mov %eax, 0x8049a1c	108	51	50
			113	51	51
interrupt	save T2's state				
	restore T1's state				
	<u>mov %eax, 0x8049a1c</u>		108	51	51
			113	51	51

Handwritten notes:

- Red arrows point from the assembly code to the corresponding rows in the table.
- Red circles highlight the final counter values: 51 (after Thread 2's update) and 51 (after Thread 1's update).
- Red "52" is written next to the final counter value 51, indicating the expected result.

Figure 26.7: The Problem: Up Close and Personal

Race conditions and synchronization

- What just happened is called a race condition
 - Concurrent execution can lead to different results
- Critical section: portion of code that can lead to race conditions
- What we need: mutual exclusion
 - Only one thread should be executing critical section at any time
- What we need: atomicity of the critical section
 - The critical section should execute like one uninterruptible instruction
- How is it achieved? Locks (topic of next lecture)

