Lecture 11: Introduction to Process Scheduling

Vivek Kumar Computer Science and Engineering IIIT Delhi vivekk@iiitd.ac.in

Lecture 11: Introduction to process scheduling

Last Lecture



==> Processors communicate by sending messages via an interconnect



- Several machines (e.g., workstations) connected together in a cluster with some interconnect
- One big computation can be divided equally into N number of subcomputations that can run on N different processes
- Processes use MPI to communicate in distributed memory by exchanging messáges with each other
- Amdhal's law: if 50% of your application is parallel and 50% is serial, you can't get more than a factor of 2 speedup, no matter how many processors it runs on

© Vivek Kumar

Today's Class

- Introduction to process scheduling
 - Context switch
 - Process scheduling policies

In this course, we will not have any assignment involving carrying out changes inside any operating systems. Still, to understand some concepts, we would revisit well-known OS, such as xv6 OS – as in today's lecture!



Process State in these Scenarios?



- 1. Write on pipe that is currently full
- 2. scanf() waiting for user input
- 3. Recursive Fibonacci number calculation
- 4. Process executing sem_wait when
 - Semaphore > 0
 - Semaphore <= 0
- 5. Process executing sem_post
- 6. sleep()
- 7. wait()

Process State in these Scenarios?





© Vivek Kumar

How OS Can Juggle Between Processes?

- <u>Case-1</u>: Whenever a process is blocking for some resource
 - It does not make any sense to let a blocked process running on the CPU (e.g., waiting for IO, sem_wait, reading from an empty pipe, etc.)
 - Even if the CPUs are idle
 - OS can safely move the process from the running to waiting queue
 - Does not affect the performance as if CPUs are idle, the process can be brought back into the running queue as soon as it unblocks
- <u>Case-2</u>: If a process is not blocking for some resource (running until termination)
 - Timer interrupt is generated by the CPU
 - It is a hardware generated interrupt at a fixed interval
 - It is used to trigger the scheduling by moving the currently running process into the ready queue, and another process from the ready queue into the running queue



Case-1: Process is Going to Block



- The user process P1 calls a blocking API (e.g., **sleep** on an empty pipe)
- P1 switches its execution from the user stack to its kernel stack (each user process has its own kernel stack)
 - User context is saved on the bottom of the kernel stack and interrupt handler is invoked
 - The blocking system call will eventually call some sleep implementation
- The sleep will invoke a call to the CPU scheduler, which will: a) save the kernel context of the user process at the top of the kernel stack, b) load the context of the CPU scheduler from the corresponding stack, and c) jump to the CPU scheduler stack where it will call the schedule() routine
- Each CPU has its own CPU scheduler process that runs on a separate kernel stack than that of the process it was originally executing

o Why?



Case-2: Process is Running until Termination



- There could be several Fibonacci processes being displayed in running state from htop on a 4-core processor
 - However, only four of them are actually in running state any given time (1-1 mapping between a process and a core)
- Timer interrupt (IDT index 32 on x86) is generated by the hardware after every fixed intervals

Case-2: Process is Not Going to Block



- Timer interrupt moves the process execution from user stack to its kernel stack (each user process has its own kernel stack)
- User context is saved on the bottom of the kernel stack and interrupt handler is invoked
- The handler will call yield() method as it was a timer interrupt
- The **yield** will invoke a call to the CPU scheduler, which will: a) save the kernel context of the user process at the top of the kernel stack, b) load the context of the CPU scheduler from the corresponding stack, and c) jump to the CPU scheduler stack where it will call the schedule() routine
- Each CPU has its own CPU scheduler process that runs on a separate kernel stack than that of the process it was originally executing



© Vivek Kumar

The Context Switch – swtch(OldCtx, NewCtx)

- Typical steps for calling a context switch on x86
 - 1. Load the parameters passed to the context switching routine in registers, as it is currently on the kernel stack of the user process
 - Parameters would be the context of current process that has to be saved, and the context of the scheduler process that has to be restored
 - The context of the current process will no longer be accessible as the ESP will not be pointing to the old stack after context switch

CSE231: Operating Systems

- 2. Load context of the scheduler process
 - a) Load the callee save registers (EBX, ESI, and EDI), EBP, and ESP from the "top" of scheduler process stack
 - b) Load EIP from scheduler stack and jump to that instruction address



- Recall, there are six general purpose registers that can be used to store temporary data during a method execution
- If the caller is using these registers, it must save each of these 6 registers before transferring the call to the callee
- However, it could be possible that callee doesn't use all these registers. Hence, the saving and restoring of these 6 registers is divided across both caller and the callee
- If the <u>caller method</u> is using these data registers then it saves/restore only these 3 on its call stack before/after method call: EAX, EDX, and ECX

© Vivek Kuma

• If the <u>callee method</u> wants to use these 3 registers then it saves/restore the current value on its call stack at entry/exit: **EBX**, **ESI**, and **EDI**



© Vivek Kumar

The Juggling API: scheduler

• Basic questions about the design

- How many scheduler processes should be there in an OS?
- How many copies of process table should be there in the OS?
 - Process table is an array of PCB
- Can there be any race on the process table accesses?
 - How to resolve?

Who Calls the scheduler



 OS is booted with the help of another program called bootloader

- Reads the kernel ELF executable
- Loads the program segments
- Jump to the entry point (kernel_main)
- Call never returns to bootloader (similar to exec)

Kernel main

- Complete the basic setup (initializing CPUs, memory, disk, registers, etc.)
- Launch the first user process (?)
- Launch one scheduler process on each CPU that will eventually pick a process from the process table and start executing it until it gives up the CPU (see case-1 and case-2 in previous slides)





The scheduler should pay attention that the process picked up from the process table is in READY state (should not be in WAITING state)

- It should change the state of this process from READY to RUNNING
- When control is passed back to the scheduler from this process, its state would have changed to either READY or WAITING



- There is a single process table on a system irrespective of the number of CPUs
- Proper locking mechanism must be ensured to avoid any race condition on the process table by the scheduler processes running at other CPUs
- Modern OS may use scalable implementations of locking to improve performance

© Vivek Kumar

```
void scheduler() {
    while(true) {
        lock(process_table);
        foreach(Process p: scheduling_algorithm(process_table)) {
            if(p->state != READY) {
                continue;
            }
            p->state = RUNNING;
            unlock(process_table);
            swtch(scheduler_process, p);
            // p is done for now..
            lock(process_table);
        }
        unlock(process_table);
      }
      unlock(process_table);
    }
}
```

Process
 scheduling
 algorithm plays an
 important role in
 the design of
 operating system

 Different algorithms are chosen according to the need

Today's Class

- Introduction to process scheduling
 - o Context switch
 - \circ Process scheduling policies
 - First In First Out / First Come First Serve
 - Shortest Job First

Acknowledgements: Lecture content derived from the CS162 course offered at University of California, Berkeley



Context Switch: Cost?



- Context switch overhead measured on an AMD EPYC 32-core processor running Ubuntu 18.04.3 LTS
 - Data generated using Imbench benchmark (./lat_ctx –s 0 32 128 512 2048 8192)

Overheads

- Timer interrupt latency
- Saving restoring context
- Process scheduling
- Loss in cache locality

CPU and I/O Bursts



- Process execution model: programs alternate between bursts of CPU and I/O
 - Program typically uses the CPU for some period of time, then does I/O, then uses CPU again
 - Each scheduling algorithm is about which process to give to the CPU for use by its next CPU burst
 - With time slicing, process may be forced to give up CPU before finishing current CPU burst



Key Scheduling Goals

• Response Time

• What user sees: time to echo a keystroke in editor

• Throughput

- \circ Total operations per second
 - Affected by the overheads of context switching

• Fairness

 \circ Conflicts with the response time

Starvation

• Due to improper resource allocation policy

First-Come, First-Served (FCFS)

Run tasks in order of arrival.

Run task until completion (or blocks on IO). No preemption

Also called FIFO



First-Come, First-Served (FCFS)



What is the average completion time?	$\left(\frac{3+6+30}{3}=13\right)$
--------------------------------------	------------------------------------





First-Come, First-Served (FCFS)



What is the average completion time?

 $\left(\frac{24+27+30}{3}=27\right)$

What is the average waiting time?

 $\left(\frac{0+24+27}{3}=17\right)$



Can there be Convoy Effect in FCFS?

Convoy effect: Short process stuck behind long process



Today's Class

- Introduction to process scheduling
 - o Context switch
 - \circ Process scheduling policies
 - First In First Out / First Come First Serve
 - Shortest Job First

Acknowledgements: Lecture content derived from the CS162 course offered at University of California, Berkeley



Shortest Job First

How can we minimize average completion time?

By scheduling jobs in order of estimated completion time



CSE231: Operating Systems

Shortest Job First



What is the average completion time?





Can SJF Lead to Starvation?

Any scheduling policy that always favours a fixed property for scheduling leads starvation





Can there be Convoy Effect in SJF?

Any non-preemptible scheduling policy suffers from convoy effect







Any scheduling policy that always favors a fixed property for scheduling leads to starvation



CSE231: Operating Systems

Reference Material

- https://pdos.csail.mit.edu/6.828/2018/xv6/book-rev11.pdf
 - You may read the chapter-5 just to have a high-level understanding of scheduling and context switch in a simple OS
 - You can avoid referring to the xv6 codebase as we will not have any assignment on reading/modifying xv6 code

Next Lecture

• Process scheduling policies

