

CS 347M (Operating Systems Minor)

Spring 2022

# Lecture 5: Process Scheduling

Mythili Vutukuru  
CSE, IIT Bombay

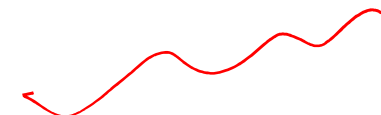
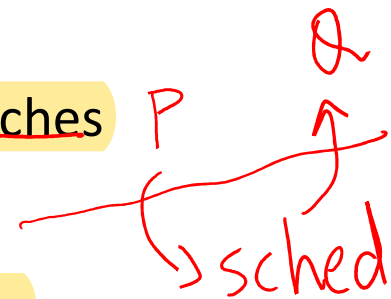
# OS scheduler



- OS scheduler schedules process on CPU cores
  - One process at a time per CPU core
  - Multiple processes can run in parallel on multiple cores
- **Scheduling policy**: which one of the ready/runnable processes should be run next on a given CPU core?
  - Mechanism of context switching (save context of old process in its kernel stack/PCB and restore context of new process) is independent of policy
- Simple scheduling policies have good theoretical guarantees, but not practical for real operating systems
  - Real-life schedulers are very complex, involve many heuristics

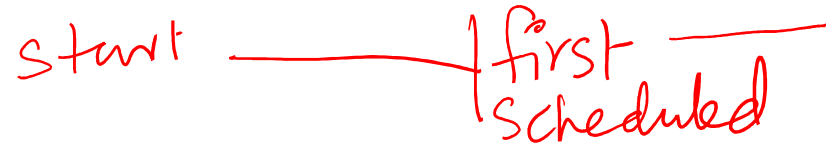
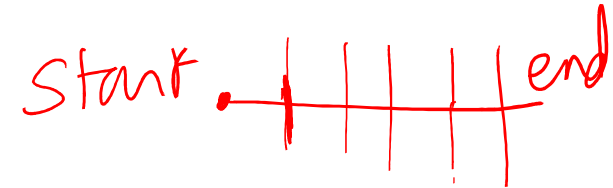
# Preemptive vs. non preemptive schedulers

- When is the OS scheduler invoked to trigger a context switch?
  - Only when a process is in kernel mode for a trap, but not on every trap
- Non-preemptive scheduler performs only voluntary context switches
  - Process makes blocking system call
  - Process has exited or has been terminated
- Preemptive scheduler performs involuntary context switches also
  - Process can be context switched out even if process is still runnable/ready
  - OS can ensure that no process runs for too long on CPU, starving others
- Timer interrupts: special interrupts that go off periodically to trap to OS
  - Used by OS to get back control, trigger involuntary context switches
- Modern systems use preemptive schedulers
  - Process can be context switched out any time in its execution



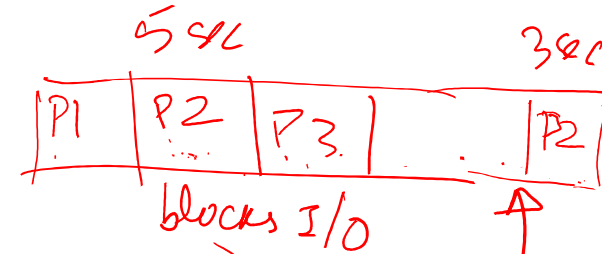
# Goals of CPU scheduling policy

- Maximize utilization: efficient use of CPU hardware
- Minimize completion time / turnaround time of a process (time from process creation to completion)
- Minimize response time of a process (time from process creation to first time it is run)
  - Important for interactive processes
- Fairness: all processes get a fair share of CPU
  - Can account for priorities also
- Low overhead of scheduling policy
  - Scheduler does not take too long to make a decision (even with large #processes)
  - Scheduler does not cause too many context switches ( $\sim 1$  microsecond to switch)



$10^{-6}$

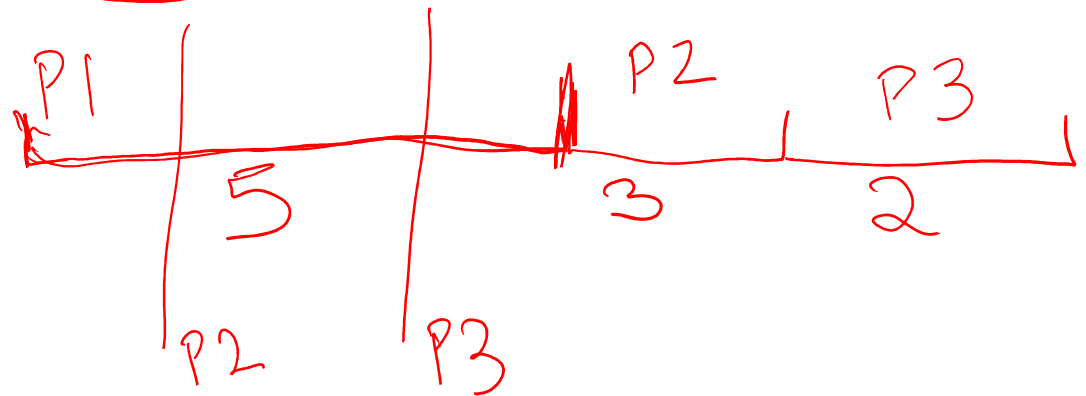
# Simplest policy: First In First Out



- Newly created processes are put in a **FIFO queue**, scheduler runs them one after another from queue
- Non-preemptive**: process allowed to run till it terminates or blocks
  - When process unblocks, the next run is separate “job”, added to queue again
  - That is, if a process comes back after I/O wait, it counts as a fresh CPU burst (**CPU burst** = the CPU time used by a process in a continuous stretch)
- Example schedule: P1 (1-5), P2 (6-8), P3 (9 to 10) *work over*

*Given*

Process	CPU time needed (units)	Arrives at end of time unit
P1	5	0
P2	3	1
P3	2	3



# Problem with FIFO

- Example: three processes arrive at  $t=0$  in the order A,B,C
- Problem: convoy effect (small processes get stuck behind long processes)
- Average turnaround times tend to be high

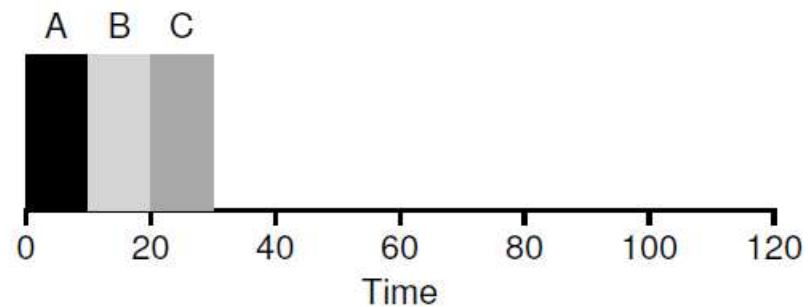


Figure 7.1: FIFO Simple Example

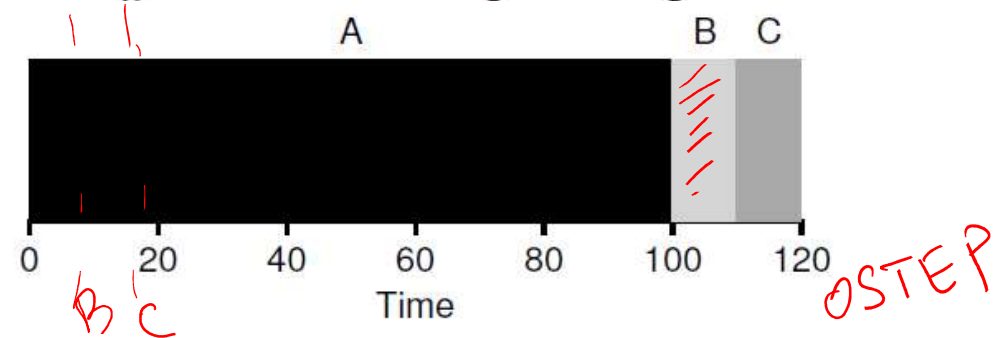
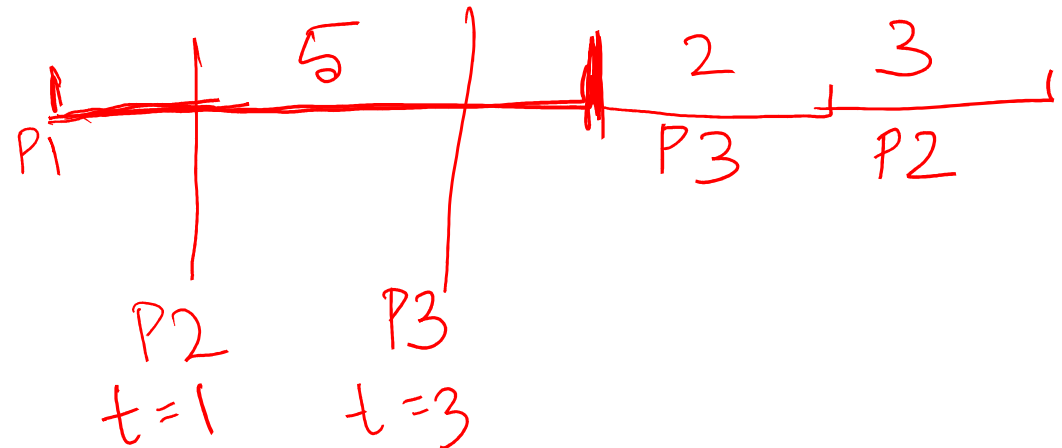


Figure 7.2: Why FIFO Is Not That Great

1 2 3 4 5 6 7 8 9 10

- unrealistic

A hand-drawn diagram illustrating a 1D lattice with three particles (P1, P2, P3) and their positions at two different times,  $t=1$  and  $t=3$ . The lattice is represented by a horizontal line with vertical tick marks. At  $t=1$ , P1 is at the first tick mark, P2 is at the second tick mark, and P3 is at the third tick mark. At  $t=3$ , P1 is at the fourth tick mark, P2 is at the fifth tick mark, and P3 is at the sixth tick mark. The distance between P1 and P2 at  $t=1$  is labeled '5', and the distance between P2 and P3 at  $t=3$  is labeled '2'.



# Problem with SJF

- Provably optimal when all processes arrive together
- Theoretically guaranteed to have the lowest average turnaround time across all policies (under certain assumptions)
- SJF is non-preemptive, so short jobs can still get stuck behind long ones.

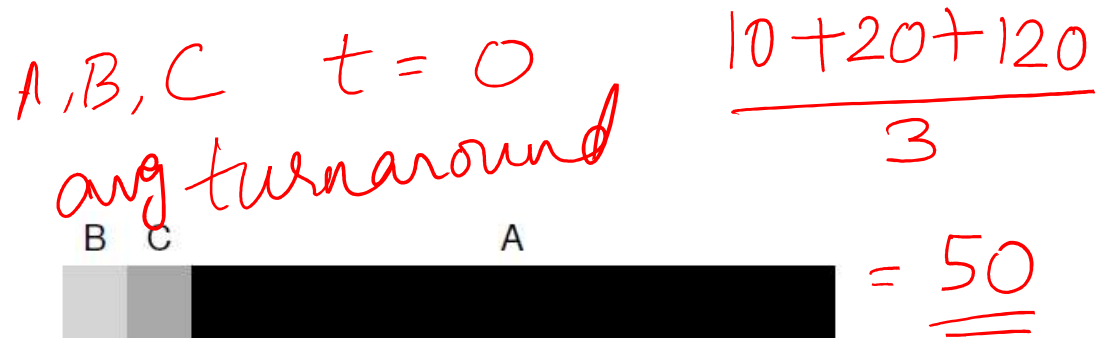


Figure 7.3: SJF Simple Example

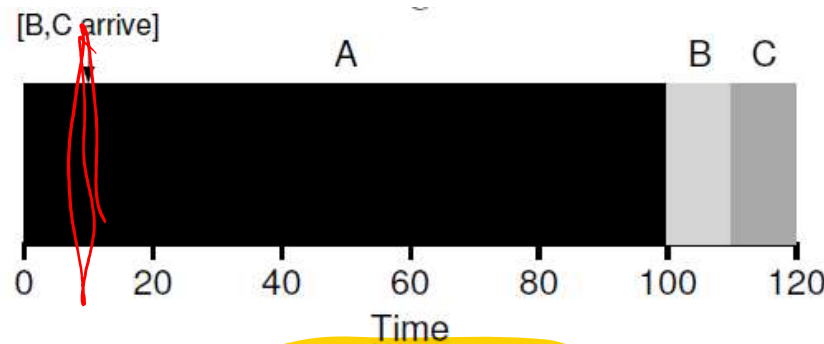


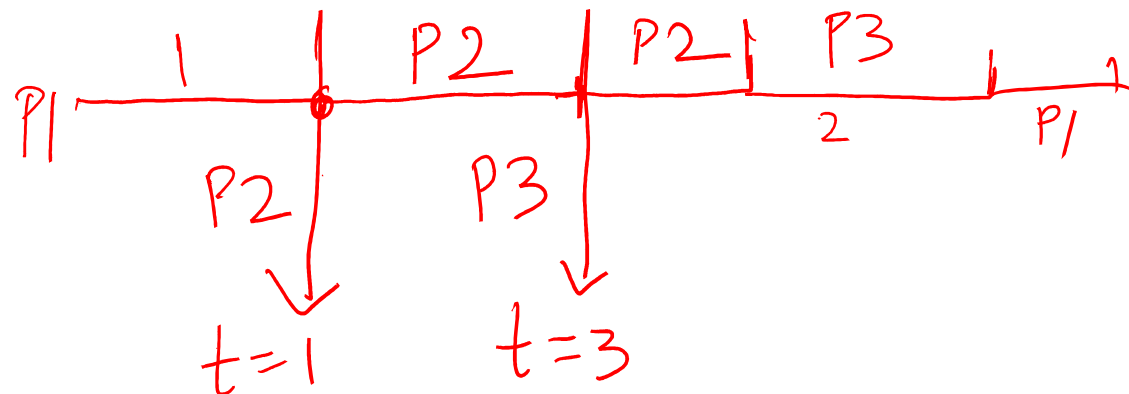
Figure 7.4: SJF With Late Arrivals From B and C



# Shortest Remaining Time First (SRTF)

- Preemptive version of SJF
- A newly arrived process can preempt a running process, if CPU burst of new process is shorter than remaining time of running process
  - Avoids problem of short process getting stuck behind long one
- Example schedule: P1 runs for 1 unit, P2 (2-4), P3 (5-6), P1 (7-10)

Process	CPU burst	Arrival time
P1	5 $- 1 = 4$	0
P2	3 $- 2 = 1$	1
P3	2	3



# Round Robin (RR)

- Every process executes for a fixed quantum slice
  - Slice not too small (to amortize cost of context switch) *1  $\mu$ s X*
  - Slice not too big (to provide good responsiveness) *10 sec X*
- Preemptive policy
  - Timer interrupt used to enforce periodic scheduling
- Good for response time and fairness
- Bad for turnaround time

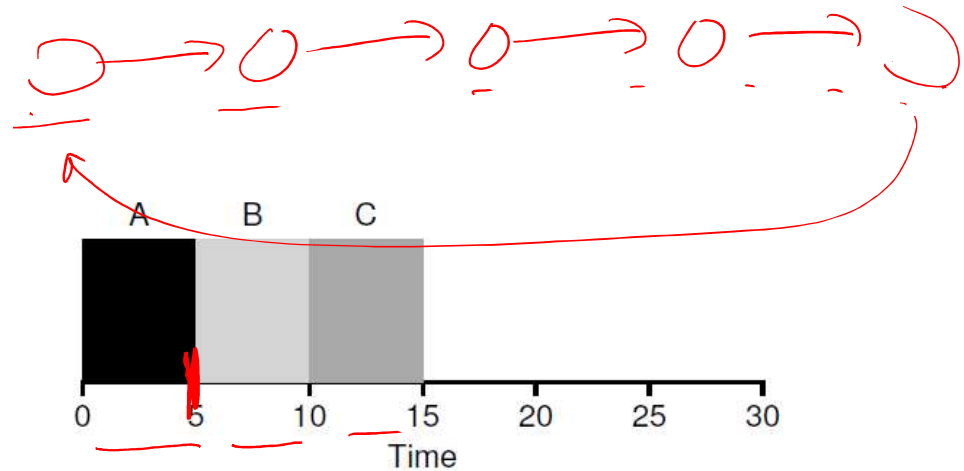


Figure 7.6: SJF Again (Bad for Response Time)

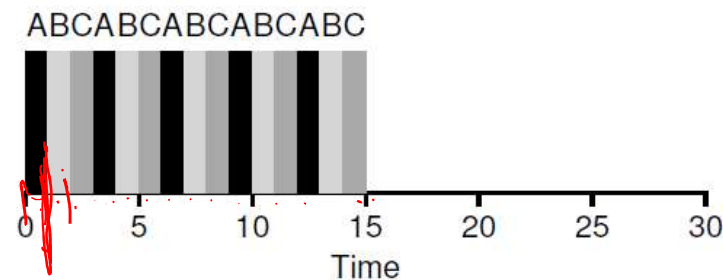


Figure 7.7: Round Robin (Good for Response Time)

# Weighted Fair Queueing (WFQ) WRR

- Round robin with different weights or priorities to processes

- Decided by scheduler or can be set by users
- Time slice will be in proportion to the weight or priority

- **Linux scheduler** is a variant of weighted fair queueing
- Real life schedulers may not be able to enforce time slice exactly

- What if timer interrupt is not exactly aligned with time slice?
- What if process blocks before its time slice?

- Practical modification: keep track of run time of process, schedule process that has used least fraction of its fair share

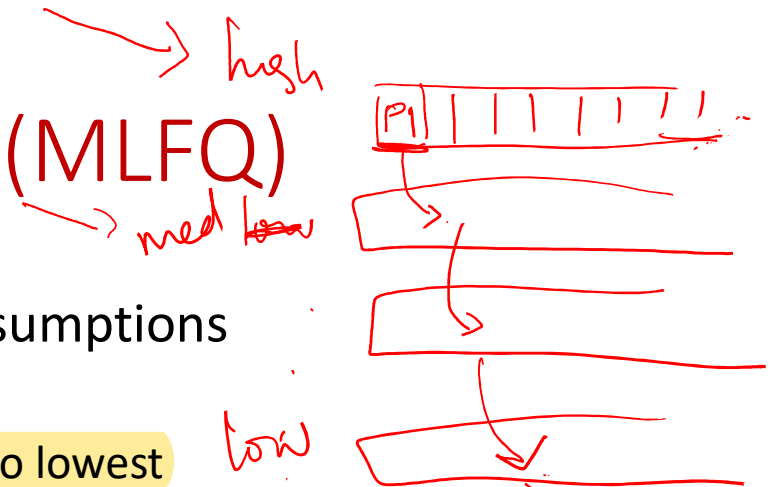
- Compensate excess/deficit running time in future time slices

P1	P2	P3
2	1	4
20ms	10ms	40ms
↓	↓	↓
16ms	11ms	33ms
<u>20</u>	<u>10</u>	<u>40</u>

work conserving

nice

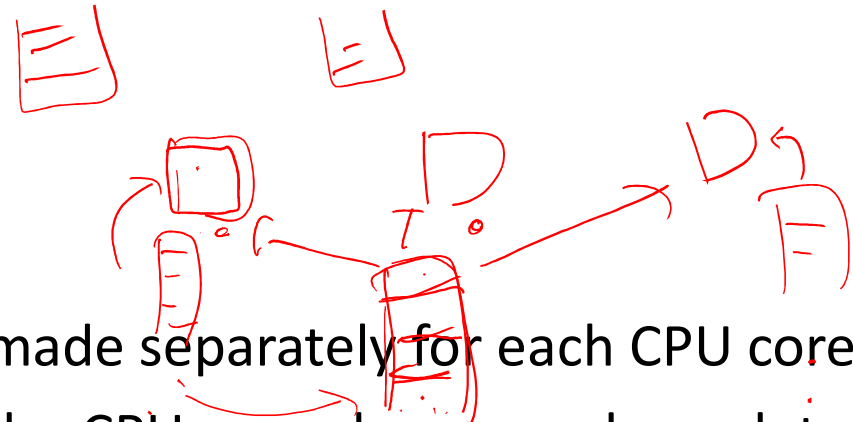
# Multi-level feedback queue (MLFQ)



- Another practical algorithm, with realistic assumptions
- Multiple queues, one for each priority level
  - Schedule processes from highest priority queue to lowest
  - Use round robin scheduling for processes within same priority level
- Priority set by user or OS, but decays with age
  - Job that uses up its time slice at a priority level goes to lower priority level
  - Why? Ensures short I/O-bound processes that don't use their full slice get priority STF over long CPU-bound processes that use their fair share all the time
  - Ensures short processes get priority over long processes, but without knowing CPU burst of process apriori
- Periodically reset all processes to highest priority level to avoid starvation of low priority or CPU-bound processes

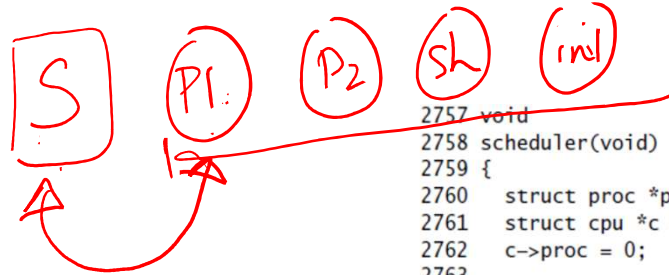
OSTEP

# Multicore scheduling



- Scheduling decision needs to be made separately for each CPU core
- Do we bind a process to a particular CPU core always, or do we let a process run on any CPU core that is free?
- Ensuring a process runs on the same core as far as possible is better
  - Avoids coordination overheads across cores, better CPU cache performance
- But, we must be flexible too
  - If CPU core overloaded, some of its processes must move to another core
  - Load balancing across cores to ensure uniform workload distribution

# Scheduler in xv6



- Every CPU has a scheduler thread (special process that runs scheduler code)
- Scheduler goes over list of processes and switches to one of the runnable ones (round robin)
- The special function “swtch” performs the actual context switch
  - Save context on kernel stack of old process
  - Restore context from kernel stack of new process

```

2757 void
2758 scheduler(void)
2759 {
2760     struct proc *p;
2761     struct cpu *c = mycpu();
2762     c->proc = 0;
2763
2764     for(;;){
2765         // Enable interrupts on this processor.
2766         sti();
2767
2768         // Loop over process table looking for process to run.
2769         acquire(&ptable.lock);
2770         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2771             if(p->state != RUNNABLE)
2772                 continue;
2773
2774             // Switch to chosen process. It is the process's job
2775             // to release ptable.lock and then reacquire it
2776             // before jumping back to us.
2777             c->proc = p;
2778             switchvm(p);
2779             p->state = RUNNING;
2780
2781             swtch(&(c->scheduler), p->context);
2782             switchkvm();
2783
2784             // Process is done running for now.
2785             // It should have changed its p->state before coming back.
2786             c->proc = 0;
2787         }
2788         release(&ptable.lock);
2789     }
2790 }
2791 }
    
```

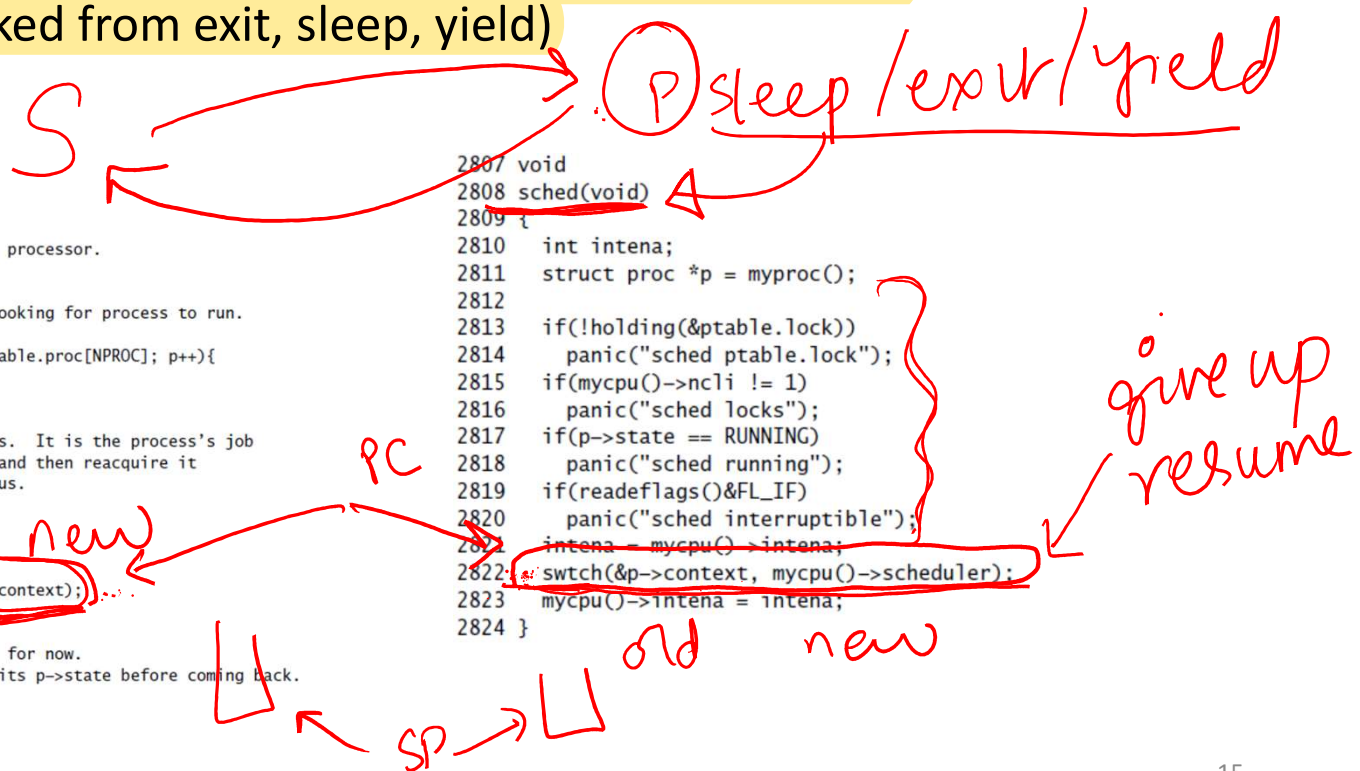
*Save context  
restore context*

# xv6 scheduler and sched functions

- Scheduler switches to user process in “scheduler” function
- User process switches to scheduler thread in the “sched” function (invoked from exit, sleep, yield)

```
2757 void
2758 scheduler(void)
2759 {
2760     struct proc *p;
2761     struct cpu *c = mycpu();
2762     c->proc = 0;
2763
2764     for(;;){
2765         // Enable interrupts on this processor.
2766         sti();
2767
2768         // Loop over process table looking for process to run.
2769         acquire(&ptable.lock);
2770         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2771             if(p->state != RUNNABLE)
2772                 continue;
2773
2774             // Switch to chosen process. It is the process's job
2775             // to release ptable.lock and then reacquire it
2776             // before jumping back to us.
2777             c->proc = p;
2778             switchvm(p);
2779             p->state = RUNNING;
2780
2781             switch(&(c->scheduler), p->context);
2782             switchvm(c);
2783
2784             // Process is done running for now.
2785             // It should have changed its p->state before coming back.
2786             c->proc = 0;
2787         }
2788         release(&ptable.lock);
2789     }
2790 }
2791 }
```

```
2807 void
2808 sched(void)
2809 {
2810     int intena;
2811     struct proc *p = myproc();
2812
2813     if(!holding(&ptable.lock))
2814         panic("sched ptable.lock");
2815     if(mycpu()->ncli != 1)
2816         panic("sched locks");
2817     if(p->state == RUNNING)
2818         panic("sched running");
2819     if(readeflags() & FL_IF)
2820         panic("sched interruptible");
2821     intena = mycpu()->intena;
2822     switch(&p->context, mycpu()->scheduler);
2823     mycpu()->intena = intena;
2824 }
```



# Who calls sched()?

- Yield: Timer interrupt occurs, process has run enough, gives up CPU
- Exit: Process has called exit, sets itself as zombie, gives up CPU
- Sleep: Process has performed a blocking action, sets itself to sleep, gives up CPU

give up CPU

timer interrupt

```
2826 // Give up the CPU for one scheduling round.
2827 void
2828 yield(void)
2829 {
2830     acquire(&ptable.lock);
2831     myproc()->state = RUNNABLE;
2832     sched();
2833     release(&ptable.lock);
2834 }
```

ready

exit

```
2662 // Jump into the scheduler, never to return.
2663 curproc->state = ZOMBIE;
2664 sched();
2665 panic("zombie exit");
2666 }
```

sleep

blocking syscall

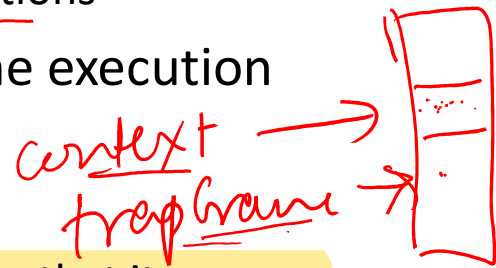
```
2894 // Go to sleep.
2895 p->chan = chan;
2896 p->state = SLEEPING;
2897
2898 sched();
2899
```

P → yield  
.. Sched()  
→ .. → yield  
Sched()



# What about new processes?

- The context switching code in xv6 restores context from kernel stack/PCB of a process and resumes execution where process stopped earlier
  - Recall: “context” structure in PCB stores context during context switch
  - Recall: “trap frame” in PCB stores context during user-kernel transitions
- What if a process has never run before? Where will it resume execution when it is switched in by scheduler?
- Kernel stack of new processes setup in such a way that
  - PC (EIP) of a function where it has to run is saved on kernel stack, so that it appears that process was switched out at the location where we want it to resume
  - Context structure and trap frame suitably created on kernel stack of new process
  - Process resumes execution in kernel mode, returns from trap to user space



# Recap: fork system call

- Parent allocates new process in ptable, copies parent state to child
- Child process set to runnable, scheduler runs it at a later time
- Return value in parent is PID of child, return value in child is set to 0

```
2579 int
2580 fork(void)
2581 {
2582     int i, pid;
2583     struct proc *np;
2584     struct proc *curproc = myproc();
2585     // Allocate process.
2586     if((np = allocproc()) == 0){
2587         return -1;
2588     }
2589 }
2590
2591 // Copy process state from proc.
2592 if((np->pgdir = copyvm(curproc->pgdir, curproc->sz)) == 0){
2593     kfree(np->kstack);
2594     np->kstack = 0;
2595     np->state = UNUSED;
2596     return -1;
2597 }
2598 np->sz = curproc->sz;
2599 np->parent = curproc;
```

*alloc new PID*

```
2600 *np->tf = *curproc->tf;
2601
2602 // Clear %eax so that fork returns 0 in the child.
2603 np->tf->eax = 0;
2604
2605 for(i = 0; i < NOFILE; i++)
2606     if(curproc->ofile[i])
2607         np->ofile[i] = filedup(curproc->ofile[i]);
2608 np->cwd = idup(curproc->cwd);
2609
2610 safestrcpy(np->name, curproc->name, sizeof(curproc->name));
2611
2612 pid = np->pid;
2613 acquire(&ptable.lock);
2614
2615 np->state = RUNNABLE;
2616 release(&ptable.lock);
2617
2618 return pid;
2620 }
2621 }
```

# allocproc

user mode

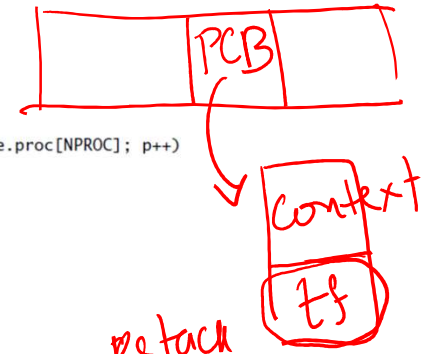
(K) forkret  
trapret

- Find unused entry in ptable, mark is as embryo
  - Marked as runnable after process creation completes
- New PID allocated
- New memory allocated for kernel stack
- Go to bottom of stack, leave space for trapframe (copied from parent trap frame in fork)
- Push return address of "trapret"
- Push context structure, with eip pointing to function "forkret"
- Why? When this new process is scheduled, it begins execution at forkret, then returns to trapret, then returns from trap to userspace
- Allocproc has created a hand-crafted kernel stack to make the process look like it had a trap and was context switched out in the past
  - Scheduler can switch this process in like any other

```

2468 // Look in the process table for an UNUSED proc.
2469 // If found, change state to EMBRYO and initialize
2470 // state required to run in the kernel.
2471 // Otherwise return 0.
2472 static struct proc*
2473 allocproc(void)
2474 {
2475     struct proc *p;
2476     char *sp;
2477
2478     acquire(&ptable.lock);
2479
2480     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++)
2481         if(p->state == UNUSED)
2482             goto found;
2483
2484     release(&ptable.lock);
2485     return 0;
2486
2487 found:
2488     p->state = EMBRYO;
2489     p->pid = nextpid++;
2490
2491     release(&ptable.lock);
2492
2493     // Allocate kernel stack.
2494     if((p->kstack = kalloc()) == 0){
2495         p->state = UNUSED;
2496         return 0;
2497     }
2498     sp = p->kstack + KSTACKSIZE;
2499
2500     // Leave room for trap frame.
2501     sp -= sizeof *p->tf;
2502     p->tf = (struct trapframe*)sp;
2503
2504     // Set up new context to start executing at forkret,
2505     // which returns to trapret.
2506     sp -= 4;
2507     *(uint*)sp = (uint)trapret;
2508
2509     sp -= sizeof *p->context;
2510     p->context = (struct context*)sp;
2511     memset(p->context, 0, sizeof *p->context);
2512     p->context->eip = (uint)forkret;
2513
2514     return p;
2515 }

```



mem kstack

Context  
forkret  
trapret

# Forking new process (revisited)

- Fork allocates new process via `allocproc`
- Parent memory and file descriptors copied (more later)
- Trapframe of child copied from that of parent
  - Result: child returns from trap to exact line of code as parent
  - Only return value in `eax` is changed, so parent and child have different return values from `fork`
- State of new child set to runnable, so scheduler thread will context switch to child process sometime in future
- Parent returns normally from trap/system call
- Child runs later when scheduled (`forkret`, `trapret`) and returns to user space like parent process

```

2579 int
2580 fork(void)
2581 {
2582     int i, pid;
2583     struct proc *np;
2584     struct proc *curproc = myproc();
2585
2586     // Allocate process.
2587     if((np = allocproc()) == 0){
2588         return -1;
2589     }
2590
2591     // Copy process state from proc.
2592     if((np->pgdir = copyvm(curproc->pgdir, curproc->sz)) == 0){
2593         kfree(np->kstack);
2594         np->kstack = 0;
2595         np->state = UNUSED;
2596         return -1;
2597     }
2598     np->sz = curproc->sz;
2599     np->parent = curproc;
2600     *np->tf = *curproc->tf;
2601
2602     // Clear %eax so that fork returns 0 in the child.
2603     np->tf->eax = 0;
2604
2605     for(i = 0; i < NOFILE; i++)
2606         if(curproc->ofile[i])
2607             np->ofile[i] = filedup(curproc->ofile[i]);
2608     np->cwd = idup(curproc->cwd);
2609
2610     safestrcpy(np->name, curproc->name, sizeof(curproc->name));
2611
2612     pid = np->pid;
2613     acquire(&table.lock);
2614     np->state = RUNNABLE;
2615     release(&table.lock);
2616
2617     return pid;
2621 }
2622

```

