

CS 347M (Operating Systems Minor)

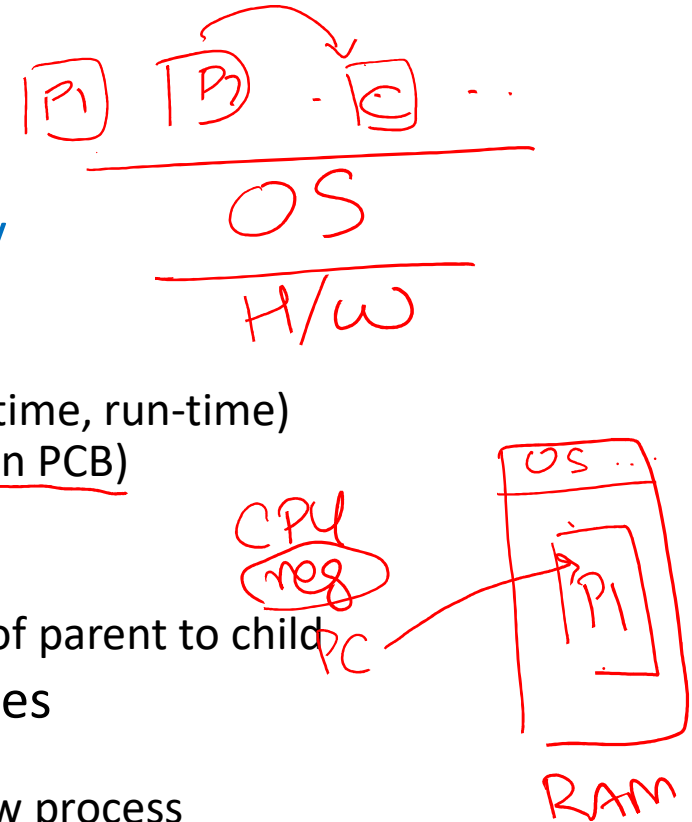
Spring 2022

Lecture 4: Kernel mode execution

Mythili Vutukuru
CSE, IIT Bombay

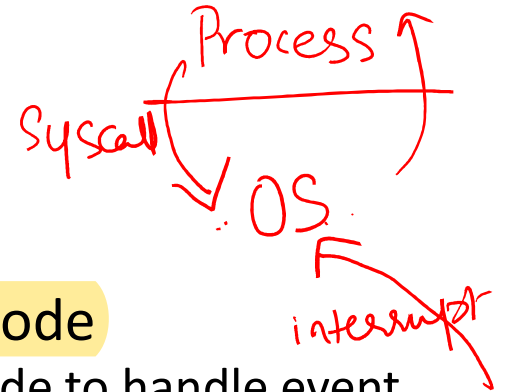
Recap: OS runs processes

- OS manages multiple active processes **concurrently**
 - Information on process in **PCB** (process control block)
- What is a process?
 - **Memory image** in RAM = compiled code, data (compile-time, run-time)
 - **CPU context** (in CPU registers when running, else saved in PCB)
 - Other things like I/O connections, ..
- Process created by **fork** from parent process
 - OS adds new process PCB to list, copies memory image of parent to child
- Periodically, **OS scheduler** loops over ready processes
 - Finds a suitable process to run next
 - Saves context of existing process, restores context of new process
- Once process is context switched in, OS is out of picture, CPU in user mode, runs user code directly
 - When does the OS run again?

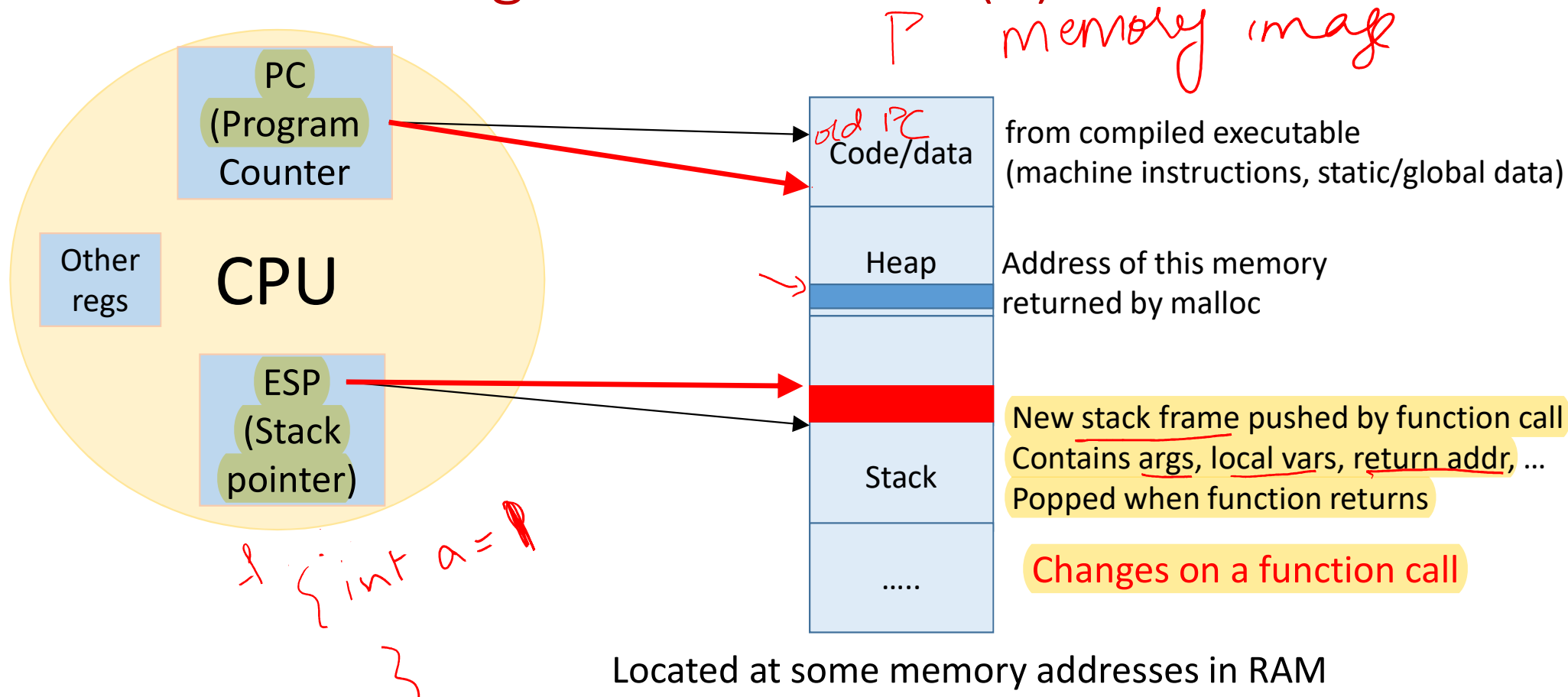


User mode vs. Kernel mode of a process

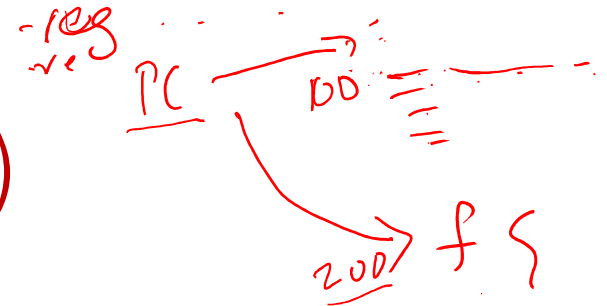
- CPU runs user code in user mode (low privilege) most of the time
- CPU switches to kernel mode execution when
 - Process makes system call, needs OS services
 - External device needs attention, raises interrupt
 - Some fault has happened during program execution
- All such events are called traps: CPU “traps” into OS code
 - CPU shifts to high privilege level (kernel mode), runs OS code to handle event
 - Later, CPU switches to low privilege level, back to user code in user mode
- Process P goes to kernel mode to run OS code, but it is still process P itself that is in running state
- OS not a separate process, runs in kernel mode of existing processes



Understanding a function call (1)



Understanding a function call (2)

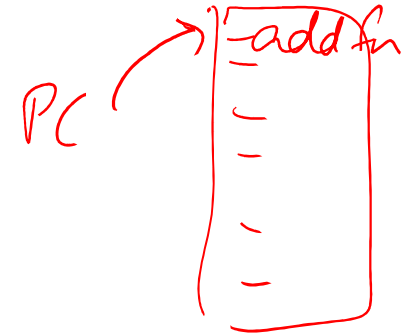


- What happens when a user program makes a function call?
 - Allocate memory on user stack for function arguments, local variables, ..
 - Push return address (PC where execution stopped), PC jumps to function code
 - Push register context (to resume execution when function returns)
 - Execute function old PC
 - When returning from function, pop return address, pop register context
- System call also must
 - Use a stack to push/pop register context
 - Save old PC, change PC to point to OS code to handle system call

compiler

privilege

System call vs. function call



- Changing PC in function call vs. system call

- In function call, address of function code known in executable, can jump to function code directly using a CPU instruction ("call" in x86) *call add fn*
- For system call, cannot trust user to jump to correct OS code (what if user jumps to inappropriate privileged code?)

- Saving register context on stack in function call vs. system call

- In function call, register context is saved and restored from user stack
- For system call, OS does not wish to use user stack (what if user has setup malicious values on the stack?)

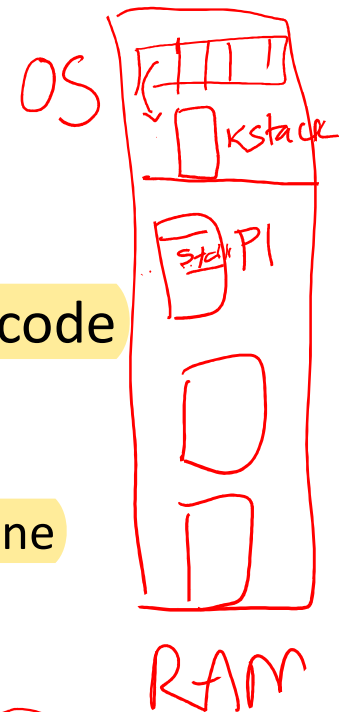


- We require: a secure stack, a secure way of jumping to OS code



Kernel stack and IDT

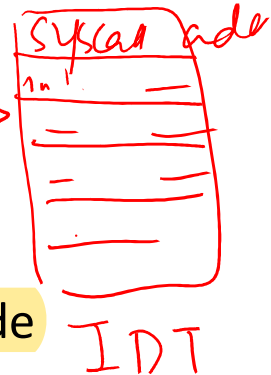
- Every process uses a separate kernel stack for running kernel code
 - Part of PCB of process, in OS memory, not accessible in user mode
 - Used like user stack, but for kernel mode execution
 - Context pushed on kernel stack during system call, popped when done
- To set PC, CPU accesses Interrupt Descriptor Table (IDT)
 - Data structure with addresses of kernel code to jump to for events
 - Setup by OS during bootup, not accessible in user mode
 - CPU uses IDT to locate address of OS code to jump to
- Together: secure way of locating OS code, secure stack for OS to run



PC

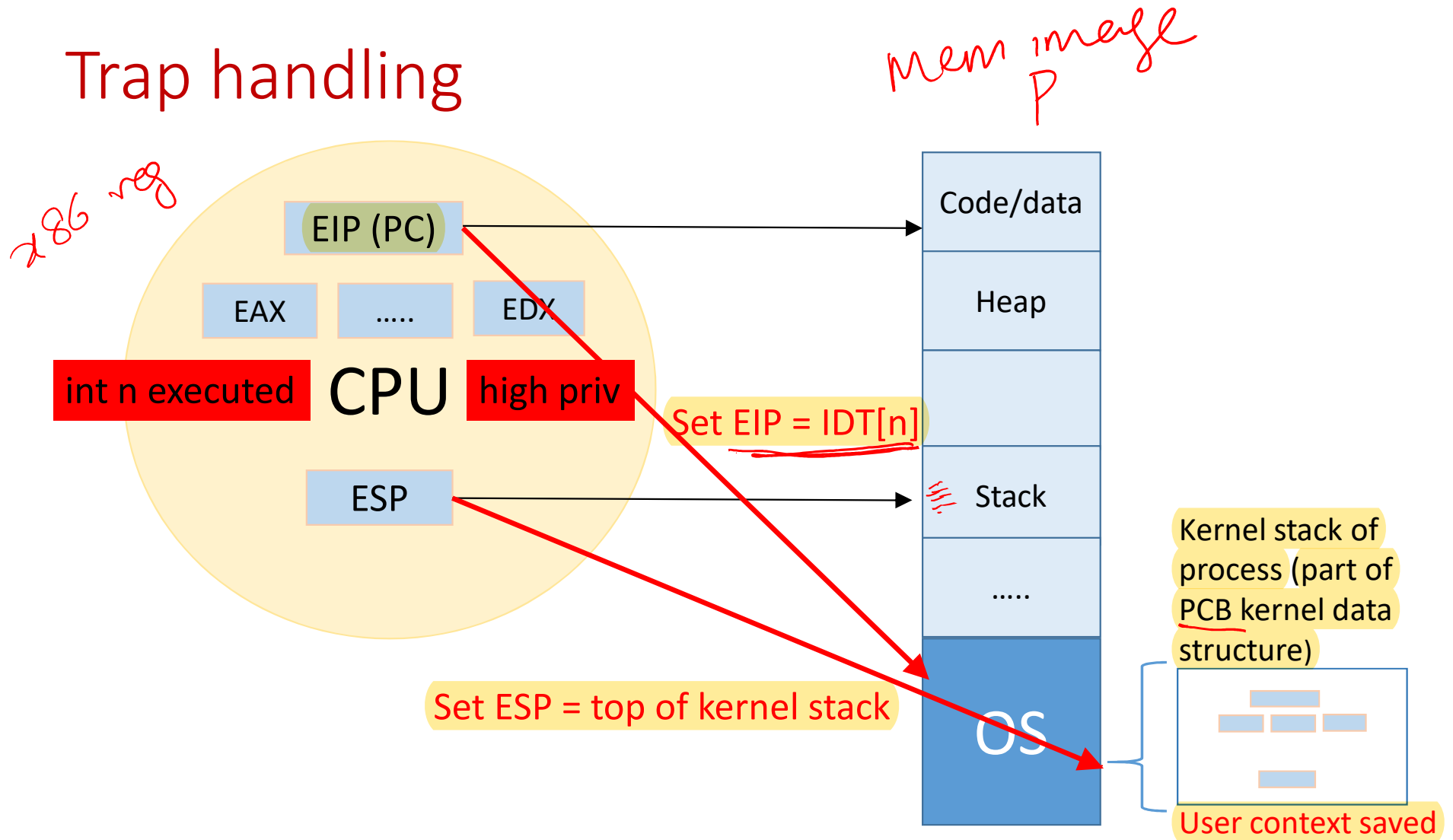
Hardware trap instruction

- When user code wants to make system call, it invokes special "trap instruction" with an argument
 - Example: "int n" in x86, argument "n" indicates type of trap (syscall, interrupt)
 - The value of "n" specifies index into IDT array, which OS function to jump to
- When CPU runs the trap instruction:
 - CPU moves to higher privilege level
 - CPU shifts stack pointer register to kernel stack of process
 - Register context is saved on kernel stack (part of PCB)
 - Address of OS code to jump to is obtained from IDT, PC points to OS code
 - OS code starts to run, on a secure stack



add
int
OS

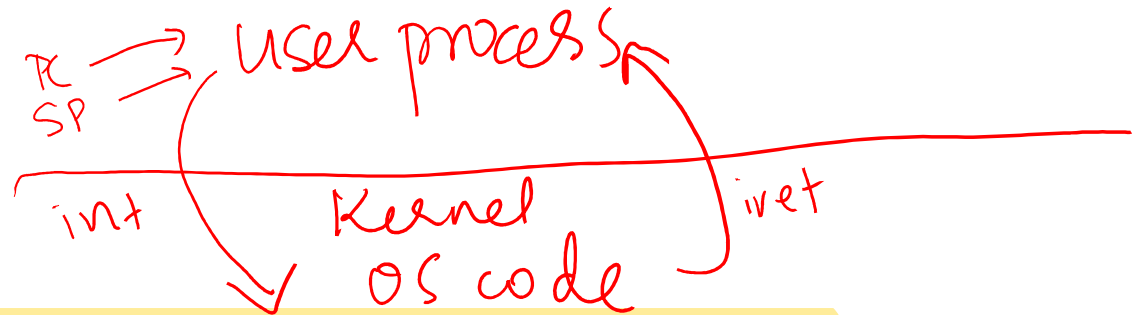
Trap handling



Why trap instruction?

- Need a secure way of jumping to OS code to handle traps
 - User code cannot be trusted to jump to correct OS code
 - Only CPU can be trusted to handover control from user to OS securely
- Who calls trap instruction?
 - System call code in a language library (printf invokes system call via int n)
 - External hardware raises interrupt, causes CPU to execute "int n"
 - Argument "n" indicates whether system call /IRQ number of hardware device
- Across all cases, the mechanism is: save context on kernel stack, switch to OS address in IDT, run OS code to handle trap

Return from trap



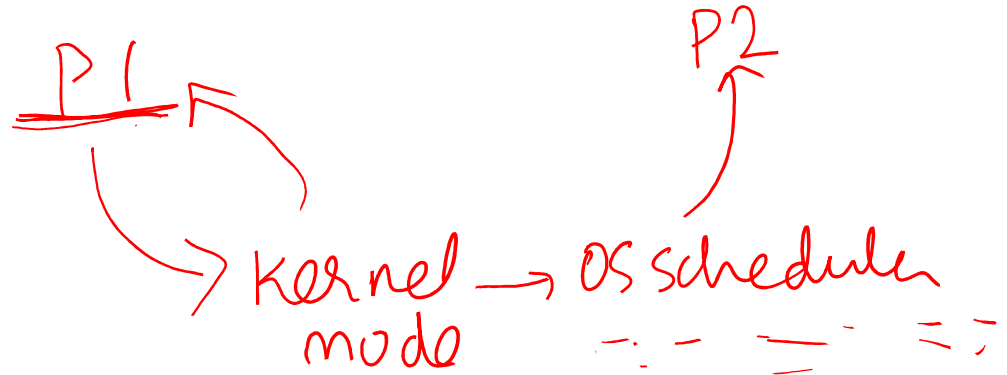
- When OS is done handling syscall or interrupt, it calls a special instruction return-from-trap
 - Restore context of CPU registers from kernel stack
 - Change CPU privilege from kernel mode to user mode
 - Restore PC and jump to user code after trap
- User process unaware that it was suspended, resumes execution at the point it stopped before
- Always return to the same user process from kernel mode? No
 - Before returning to user mode, OS checks if it must switch to another process

Why switch between processes?



- Sometimes when OS is in kernel mode, it cannot return back to the same process that was running in user mode before
 - Process has exited or must be terminated (e.g., segfault)
 - Process has made a blocking system call
- Sometimes, the OS does not want to return back to the same process
 - The process has run for too long
 - Must timeshare CPU with other processes
- In such cases, OS performs a context switch from one process to another

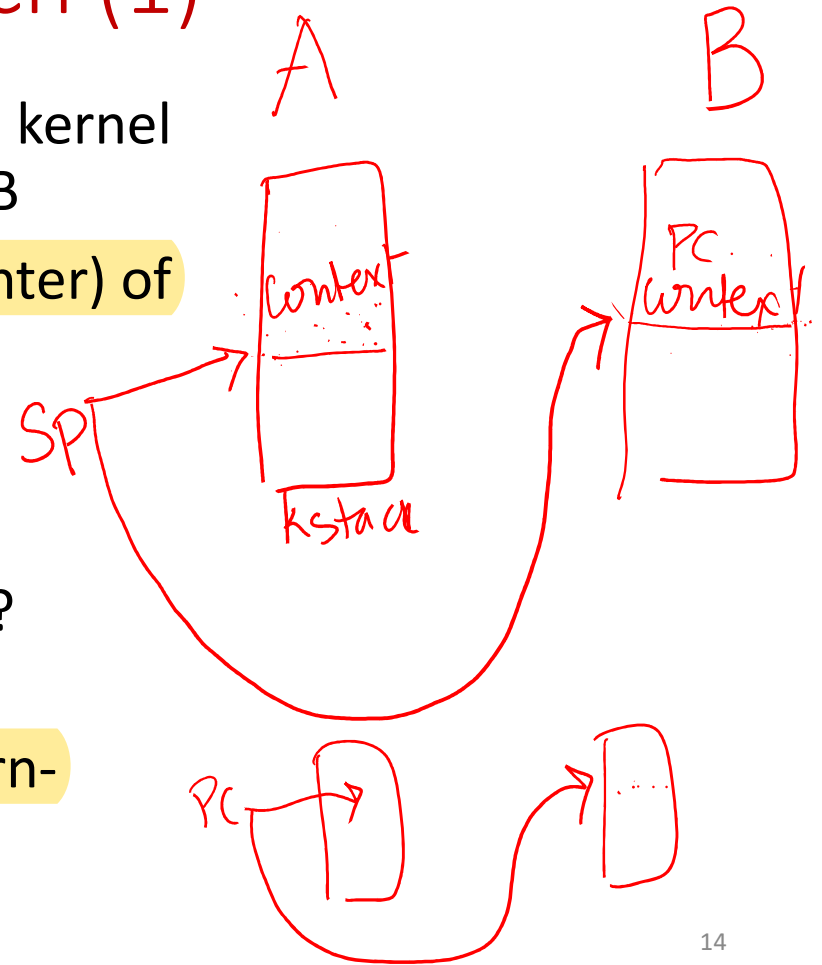
The OS scheduler



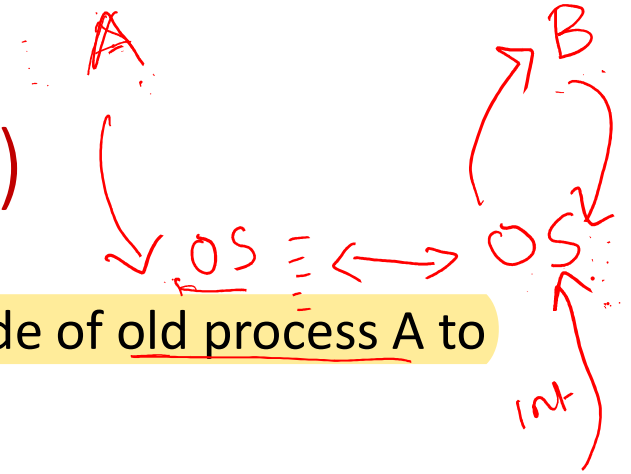
- OS scheduler has two parts
 - Policy to pick which process to run (next lecture)
 - Mechanism to switch to that process (this lecture)
- Non preemptive (cooperative) schedulers are polite
 - Switch only if process blocked or terminated
- Preemptive (non-cooperative) schedulers can switch even when process is ready to continue ✓
 - CPU generates periodic timer interrupt
 - After servicing interrupt, OS checks if the current process has run for too long

Mechanism of context switch (1)

- Example: process A has moved from user to kernel mode, OS decides it must switch from A to B
- Save context (PC, registers, kernel stack pointer) of A on kernel stack
- Switch SP to kernel stack of B
- Restore context from B's kernel stack
- Who has saved registers on B's kernel stack?
 - OS did, when it switched out B in the past
- Now, CPU is running B in kernel mode, return-from-trap to switch to user mode of B



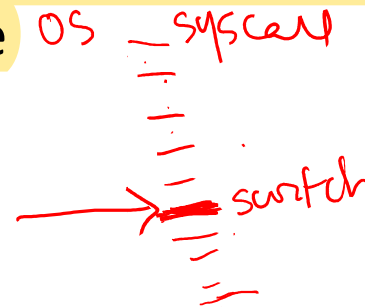
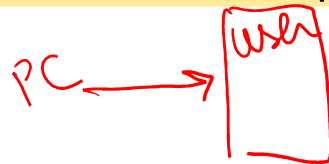
Mechanism of context switch (2)



- **Context switch:** switch CPU context from kernel mode of old process A to kernel mode of new process B
- **Before context switch**
 - A entered kernel mode, OS decides not to run A anymore (e.g., blocking system call)
 - CPU registers have context of A (stack pointer is pointing kernel stack of A, PC is pointing to some OS code being run by A)
- **Mechanism of context switch**
 - Save CPU context of A into kernel stack/PCB of A
 - Load CPU context from kernel stack/PCB of B into CPU registers
- **After context switch**
 - B resumes execution in kernel mode, stack pointer points to B's kernel stack
 - Where does B begin execution? At some point in the past, B went into kernel mode, and was switched out by OS. B resumes execution in same place.

Understand saving and restoring context (1)

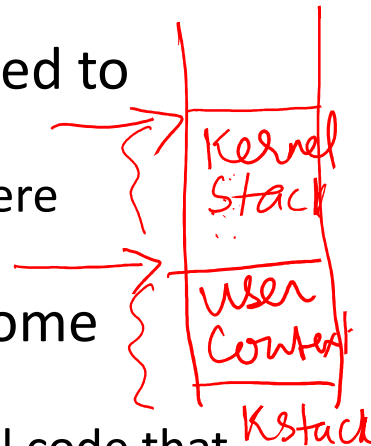
- Context (PC and other CPU registers) saved on the kernel stack in two different scenarios
- When going from user mode to kernel mode, user context (e.g., which instruction of user code you stopped at) is saved on kernel stack by the trap instruction
 - Restored by return-from-trap
- During a context switch, kernel context (e.g., where you stopped in the OS code) of process A is saved on the kernel stack of A by the context switching code
 - Restores kernel context of process B





Understand saving and restoring context (2)

- Suppose process A has made a blocking system call and moved to kernel mode
 - Context of user mode execution (e.g., PC pointing to user code where execution stopped) is saved on kernel stack/PCB
- After handling system call, OS decides to context switch to some other process, since A cannot continue now
 - Again, context of kernel mode execution (e.g., PC pointing to kernel code that has handled system call) is saved in kernel stack/PCB
- When A becomes ready and is run by scheduler again in future
 - Kernel context is restored from kernel stack/PCB into CPU registers, CPU resumes running in kernel mode of A
- A returns from trap into user mode
 - User context is restored, CPU resumes running user code of A



Trap handling in xv6

- The following events cause a user process to “trap” into the kernel (xv6 refers to all these events as traps)
 - System calls (requests by user for OS services)
 - Interrupts (external device wants attention)
 - Program fault (illegal action by program)
- When above events happen, CPU executes the special “int” instruction
 - Example seen in `usys.S`, “int” invoked to handle system calls
 - For hardware interrupts, device sends a signal to CPU, and CPU executes `int`
- Trap instruction has a parameter (int n), indicating type of interrupt
 - E.g., `syscall` has a different value of `n` from keyboard interrupt
 - The value of “n” is used to index into IDT, get address of kernel code to run

Trap frame on kernel stack

PC → user
inf IDT → OS code

- Trap frame: state is pushed on kernel stack during trap handling
 - CPU context of where execution stopped is saved, so that it can be resumed after trap
 - Some extra information needed by trap handler is also saved
- The “int n” instruction pushes a few entries (old PC, old SP etc.) and jumps to kernel code to handle trap
- The kernel code that is run next will push remaining registers on kernel stack, and then proceed to handle the trap

```
0600 // Layout of the trap frame built on the stack by the
0601 // hardware and by trapasm.S, and passed to trap().
0602 struct trapframe {
0603     // registers as pushed by pusha
0604     uint edi;
0605     uint esi;
0606     uint ebp;
0607     uint oesp; // useless & ignored
0608     uint ebx;
0609     uint edx;
0610     uint ecx;
0611     uint eax;
0612
0613     // rest of trap frame
0614     ushort gs;
0615     ushort padding1;
0616     ushort fs;
0617     ushort padding2;
0618     ushort es;
0619     ushort padding3;
0620     ushort ds;
0621     ushort padding4;
0622     uint trapno;
0623
0624     // below here defined by x86 hardware
0625     uint err;
0626     uint eip;
0627     ushort cs;
0628     ushort padding5;
0629     uint eflags;
0630
0631     // below here only when crossing rings, such as from user to kernel
0632     uint esp;
0633     ushort ss;
0634     ushort padding6;
0635 };
```

trap
frame
PC

Kstack

Trap handler function in xv6 (1)

- C trap handler performs different actions based on kind of trap
- If system call, "int n" is invoked with "n" equal to a value T_SYSCALL (in usys.S), indicating this trap is a system call
- Trap handler invokes common system call function
 - Looks at system call number stored in eax (whether fork or exec or)
 - Return value of syscall stored in eax

syscall
Kbd
disk

```
3400 void
3401 trap(struct trapframe *tf)
3402 {
3403     if(tf->trapno == T_SYSCALL){
3404         if(myproc()->killed)
3405             exit();
3406         myproc()->tf = tf;
3407         syscall();
3408         if(myproc()->killed)
3409             exit();
3410         return;
3411     }
```

```
3700 void
3701 syscall(void)
3702 {
3703     int num;
3704     struct proc *curproc = myproc();
3705     num = curproc->tf->eax;
3706     if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
3707         curproc->tf->eax = syscalls[num]();
3708     } else {
3709         cprintf("%d %s: unknown sys call %d\n",
3710             curproc->pid, curproc->name, num);
3711         curproc->tf->eax = -1;
3712     }
3713 }
3714 }
```

Trap handler function in xv6 (2)

- If interrupt from a device, corresponding device-related code is called
 - The trap number (value of “n” in “int n”) is different for different devices
- Timer is special hardware interrupt, and is generated periodically to trap to kernel

```
3413 switch(tf->trapno){
3414 case T_IRQ0 + IRQ_TIMER:
3415     if(cpuid() == 0){
3416         acquire(&tickslock);
3417         ticks++;
3418         wakeup(&ticks);
3419         release(&tickslock);
3420     }
3421     lapiceoi();
3422     break;
3423 case T_IRQ0 + IRQ_IDE:
3424     ideintr();
3425     lapiceoi();
3426     break;
3427 case T_IRQ0 + IRQ_IDE+1:
3428     // Bochs generates spurious IDE1 interrupts..
3429     break;
3430 case T_IRQ0 + IRQ_KBD:
3431     kbdintr();
3432     lapiceoi();
3433     break;
```

Trap handler function in xv6 (3)

- On timer interrupt, a process “yields” CPU to scheduler
- Ensures a process does not run for too long

```
3471 // Force process to give up CPU on clock tick.
3472 // If interrupts were on while locks held, would need to check nlock.
3473 if(myproc() && myproc()->state == RUNNING &&
3474     tf->trapno == T_IRQ0+IRQ_TIMER)
3475     yield();
3476
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 }
```



Context switching in xv6 (1)

- 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
- 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

swtch kstack

```
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             switchvm();
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 }
```


Context switching in xv6 (2)

```
graph LR
    P1 -- "sched()" --> Scheduler
    Scheduler --> P2
    P2 --> Scheduler
    Scheduler --> P1
```

- After running for some time, the process switches back to the scheduler thread, when:
 - Process has terminated (exit system call)
 - Process needs to sleep (e.g., blocking read system call)
 - Process yields after running for long (timer interrupt)
- Process calls “sched” which calls “switch” to switch to scheduler thread again
- Scheduler thread runs its loop and picks next process to run, and the story repeats

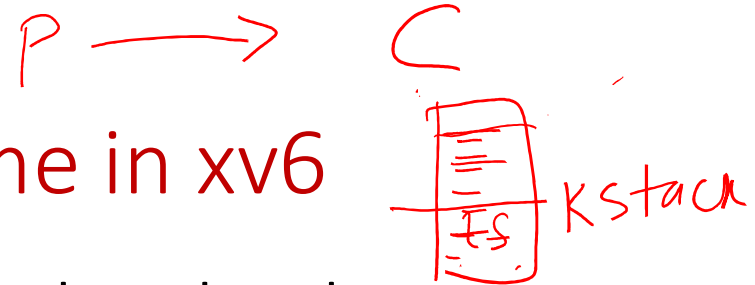
```
2662 // Jump into the scheduler, never to return.
2663 curproc->state = ZOMBIE;
2664 sched();
2665 panic("zombie exit");
2666 }
```

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

```
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 }
```

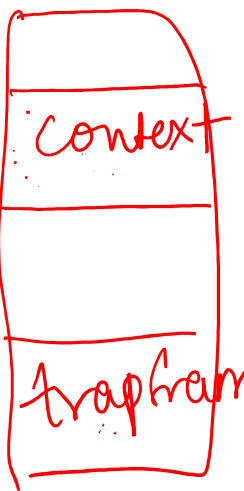
sched() switch

Context structure vs. trap frame in xv6



- Struct proc stores two different structures on kernel stack
 - Trapframe is saved when CPU switches to kernel mode (e.g., PC in trapframe is PC where syscall was made in user code)
 - Context structure is saved when process switches to another process (e.g., PC value when context switch is performed)
 - Both reside on kernel stack, struct proc has pointers to both
 - Example: Process has timer interrupt, saves trapframe on kstack, then context switch, saves context structure on kstack

kstack



```
2342 int pid; // Process ID
2343 struct proc *parent; // Parent process
2344 struct trapframe *tf; // Trap frame for current syscall
2345 struct context *context; // swch() here to run process
```

PCB
Struct
proc