

CS 347M (Operating Systems Minor)

Spring 2022

# Lecture 3: System calls for process management

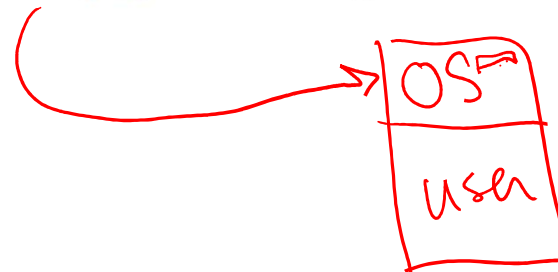
Mythili Vutukuru  
CSE, IIT Bombay

# Recap: Process and PCB



- OS manages multiple active processes
- Maintains one PCB for every active process (struct proc in xv6)
  - Process PID
  - Process state
  - Saved CPU context
  - Memory and I/O information

```
2334 enum procstate { UNUSED, EMBRYO, SLEEPING, RUNNABLE, RUNNING, ZOMBIE };
2335
2336 // Per-process state
2337 struct proc {
2338     uint sz; // Size of process memory (bytes)
2339     pde_t* pgdir; // Page table
2340     char *kstack; // Bottom of kernel stack for this process
2341     enum procstate state; // Process state
2342     int pid; // Process ID
2343     struct proc *parent; // Parent process
2344     struct trapframe *tf; // Trap frame for current syscall
2345     struct context *context; // switch() here to run process
2346     void *chan; // If non-zero, sleeping on chan
2347     int killed; // If non-zero, have been killed
2348     struct file *ofile[NOFILE]; // Open files
2349     struct inode *cwd; // Current directory
2350     char name[16]; // Process name (debugging)
2351 };
2352
```



## Recap: Process table (ptable) in xv6

```
2409 struct {
2410     struct spinlock lock;
2411     struct proc proc[NPROC];
2412 } ptable;
```

- ptable: Fixed-size array of all processes
  - Real kernels have dynamic-sized data structures
- CPU scheduler in the OS loops over all runnable processes, picks one, and sets it running on the CPU

```
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;
```

# This lecture: API for process management

*printf* → *syscall*

- What API does OS provide to user programs to manage processes?
  - How to create, run, terminate processes?
- API = Application Programming Interface  
= functions available to write user programs
- API provided by OS is a set of “system calls”
  - System call is a function call into OS code that runs at higher CPU privilege level
  - Sensitive operations (e.g., access to hardware) are allowed only at a higher privilege level
- Some “blocking” system calls cause the process to be blocked and descheduled (e.g., `read` from disk), while others can return immediately

*user mode*  
*low priv*

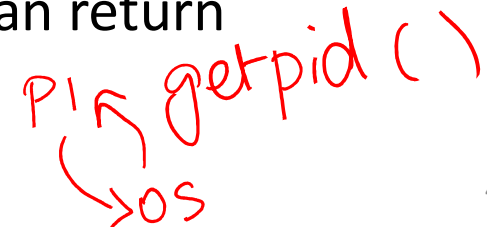
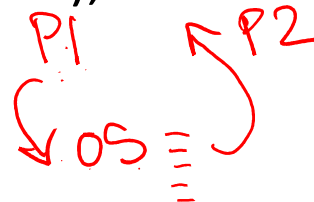
*user*

*syscall*

*OS*

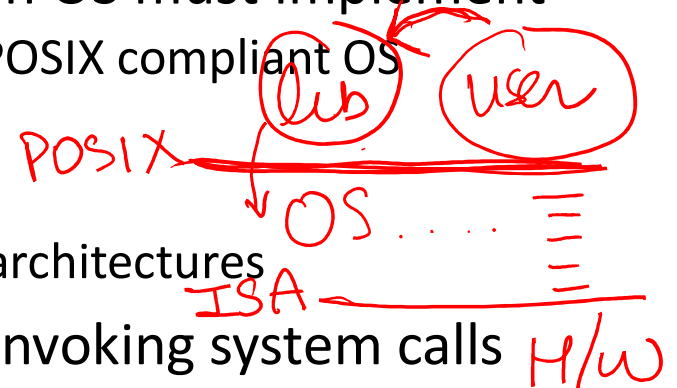
*high priv*

*kernel mode*



# So, should we rewrite programs for each OS?

- POSIX API: a standard set of system calls that an OS must implement
  - Programs written using POSIX API can run on any POSIX compliant OS
  - Most modern OSes are POSIX compliant
  - Ensures program portability
  - Program may need to be recompiled for different architectures
- Program language libraries hide the details of invoking system calls
  - The printf function in the C library calls the write system call to write to screen
  - User programs usually do not need to worry about invoking system calls



# What happens on a system call? (1)

- System calls are usually made by user library functions, e.g., C library
  - User code invokes library function only
- Example in xv6: system calls available to user programs are defined in user library header “user.h”
  - Equivalent to C library headers (xv6 doesn't use standard C library)

```
struct stat;  
struct rtcdate;  
  
// system calls  
int fork(void);  
int exit(void) __attribute__((noreturn));  
int wait(void);  
int pipe(int*);  
int write(int, const void*, int);  
int read(int, void*, int);  
int close(int);  
int kill(int);  
int exec(char*, char**);  
int open(const char*, int);  
int mknod(const char*, short, short);  
int unlink(const char*);  
int fstat(int fd, struct stat*);  
int link(const char*, const char*);  
int mkdir(const char*);  
int chdir(const char*);  
int dup(int);  
int getpid(void);  
char* sbrk(int);  
int sleep(int);  
int uptime(void);
```

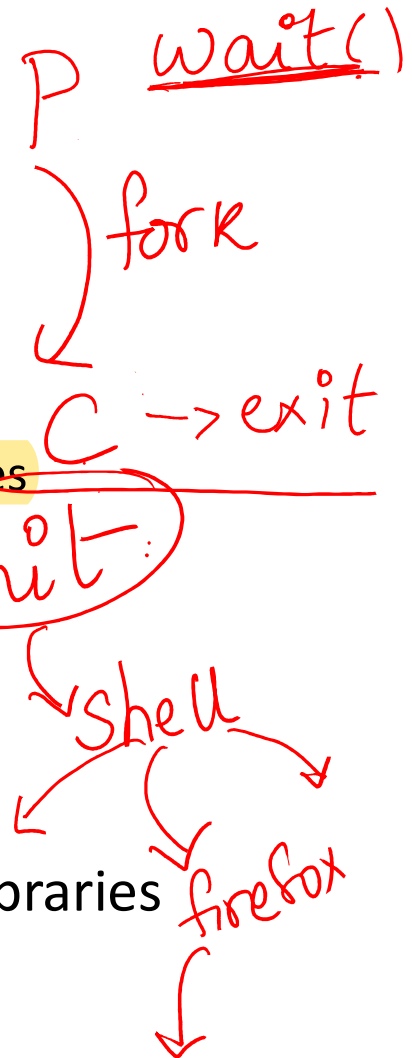


→ syscall



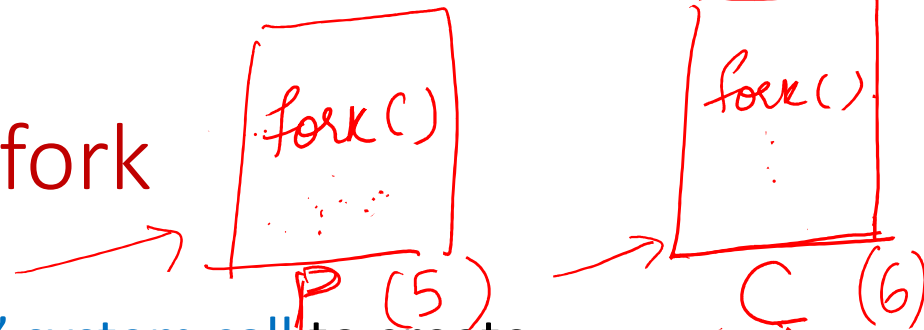
# Process related system calls (in Unix)

- fork() creates a new child process
  - All processes are created by forking from a parent
  - OS starts init process after boot up, which forks other processes
  - The init process is ancestor of all processes
- exec() makes a process execute a given executable
- exit() terminates a process
- wait() causes a parent to block until child terminates
- Many variants of the above system calls exist in language libraries with different arguments





# Process creation: fork

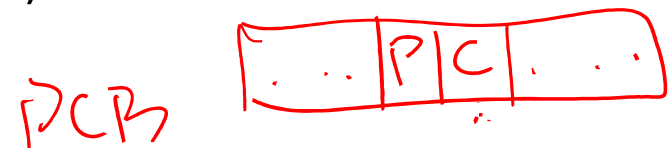


- Parent process calls “fork” system call to create (spawn) a new process
  - New child process created with new PID
  - Memory image of parent is copied into that of child
  - Parent and child run different copies of same code
  - Parent and child resume execution in the code after “fork”
  - Child starts executing with a return value of 0 from fork
  - Parent resumes executing with a return value of child PID
  - After fork, parent and child run independently (any changes in parent’s data after fork does not impact child)

```
... a = 1
int ret = fork();
if (ret == 0) {
    print "I am child"
}
else if (ret > 0) {
    print "I am parent"
}
... a = 2
```

Handwritten annotations on the code: 'a = 1' is circled in red. A red arrow points from the 'if (ret == 0)' block to the text 'a = 3'. A red bracket groups the 'else if' block with the text 'ret = pid'. Another red bracket groups the 'else if' block with the text 'a = 2'.

test.c



```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  cnt cnt a
5  int
6  main(int argc, char *argv[])
7  {
8      printf("hello world (pid:%d)\n", (int) getpid());
9      int rc = fork();
10     if (rc < 0) {          // fork failed; exit
11         fprintf(stderr, "fork failed\n");
12         exit(1);
13     } else if (rc == 0) { // child (new process)
14         printf("hello, I am child (pid:%d)\n", (int) getpid());
15     } else {              // parent goes down this path (main)
16         printf("hello, I am parent of %d (pid:%d)\n",
17               rc, (int) getpid());
18     }
19     return 0;
20 }

```

OSTEP

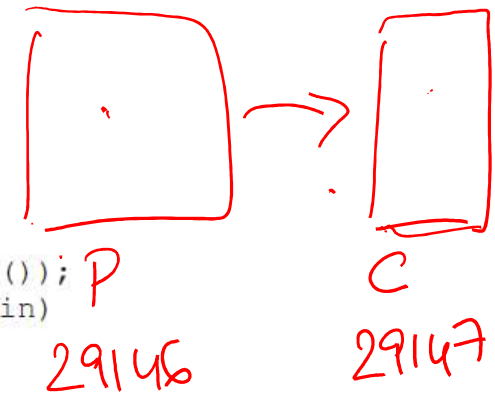


Figure 5.1: Calling fork () (p1.c)

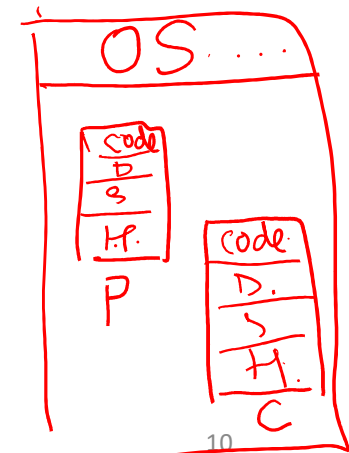
When you run this program (called p1.c), you'll see the following:

```

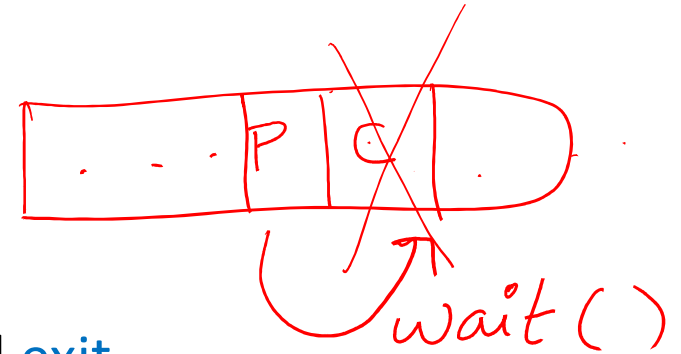
prompt> ./p1
hello world (pid:29146).
hello, I am parent of 29147 (pid:29146)
hello, I am child (pid:29147)
prompt>

```

RAM



# Exit and wait system calls



- When a process finishes execution, it called exit system call to terminate
  - OS switches the process out and never runs it again
  - Exit is automatically called at end of main
  - Process does not disappear, only becomes zombie
- Parent calls "wait" system call to reap (clean up memory of) a zombie child
  - Wait system call blocks parent until child exits
  - After child exit, wait cleans up memory of child and returns in parent process

```
...  
int ret = fork()  
if(ret == 0) {  
    print "I am child"  
    exit()  
}  
else if(ret > 0) {  
    print "I am parent"  
    wait()  
}  
...
```

x P      C ✓  
wait() ← exit

exit ↘

P

- fork()
- fork()
- wait()
- wait()

- GP



```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <sys/wait.h>
5
6  int
7  main(int argc, char *argv[])
8  {
9      printf("hello world (pid:%d)\n", (int) getpid());
10     int rc = fork();
11     if (rc < 0) {          // fork failed; exit
12         fprintf(stderr, "fork failed\n");
13         exit(1);
14     } else if (rc == 0) { // child (new process)
15         printf("hello, I am child (pid:%d)\n", (int) getpid());
16     } else {              // parent goes down this path (main)
17         int wc = wait(NULL);
18         printf("hello, I am parent of %d (wc:%d) (pid:%d)\n",
19               rc, wc, (int) getpid());
20     }
21     return 0;
22 }

```

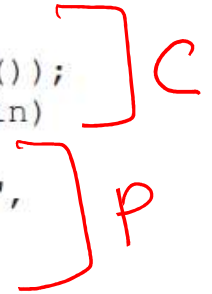
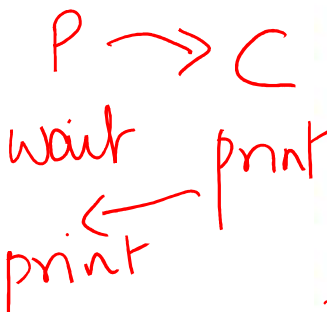
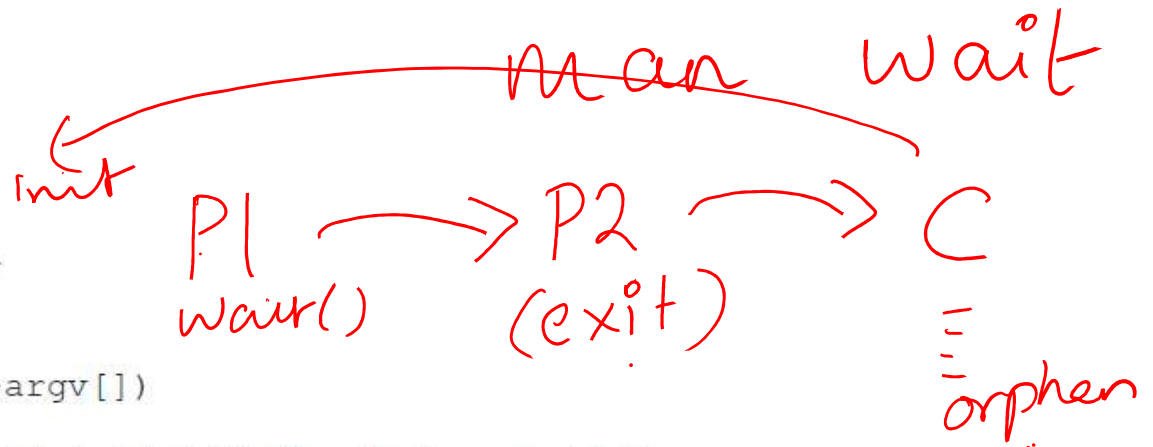
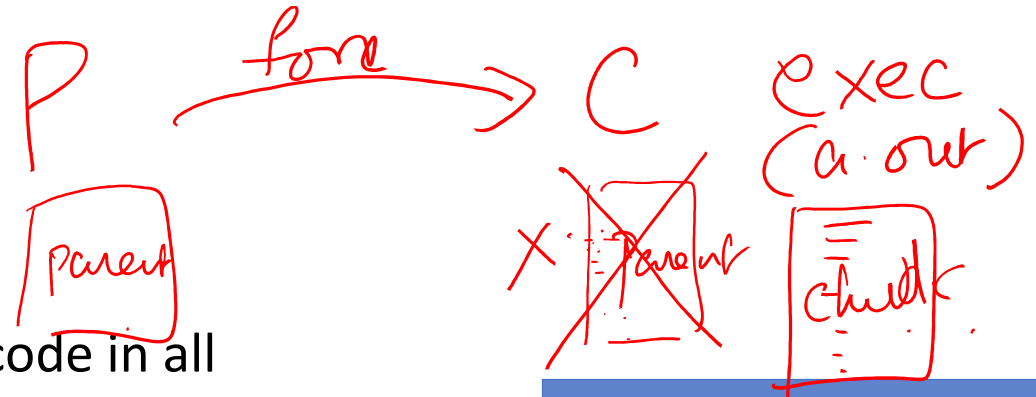


Figure 5.2: Calling `fork()` And `wait()` (p2.c)

# Exec system call



- Isn't it impractical to run the same code in all processes?
  - Sometimes parent creates child to do similar work..
  - .. but other times, child may want to run different code
- Child process uses "exec" system call to get a new "memory image"
  - Allows a process to switch to running different code
  - Exec system call takes another executable as argument
  - Memory image is reinitialized with new executable, new code, data, stack, heap, ...
  - Child process does not return to old parent program (unless exec fails)
  - Print statement after exec never prints unless exec fails

```
...
int ret = fork();
if(ret == 0) {
    exec("some_executable")
    print "error: exec failed"
}
else if(ret > 0) {
    print "I am parent"
}
...
```

parent.c

(child.c → a.out)

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <string.h>
5  #include <sys/wait.h>
6
7  int
8  main(int argc, char *argv[])
9  {
10     printf("hello world (pid:%d)\n", (int) getpid());
11     int rc = fork();
12     if (rc < 0) {           // fork failed; exit
13         fprintf(stderr, "fork failed\n");
14         exit(1);
15     } else if (rc == 0) { // child (new process)
16         printf("hello, I am child (pid:%d)\n", (int) getpid());
17         char *myargs[3];
18         myargs[0] = strdup("wc"); // program: "wc" (word count)
19         myargs[1] = strdup("p3.c"); // argument: file to count
20         myargs[2] = NULL;          // marks end of array
21         execvp(myargs[0], myargs); // runs word count
22         printf("this shouldn't print out");
23     } else {                // parent goes down this path (main)
24         int wc = wait(NULL);
25         printf("hello, I am parent of %d (wc:%d) (pid:%d)\n",
26               rc, wc, (int) getpid());
27     }
28     return 0;
29 }

```

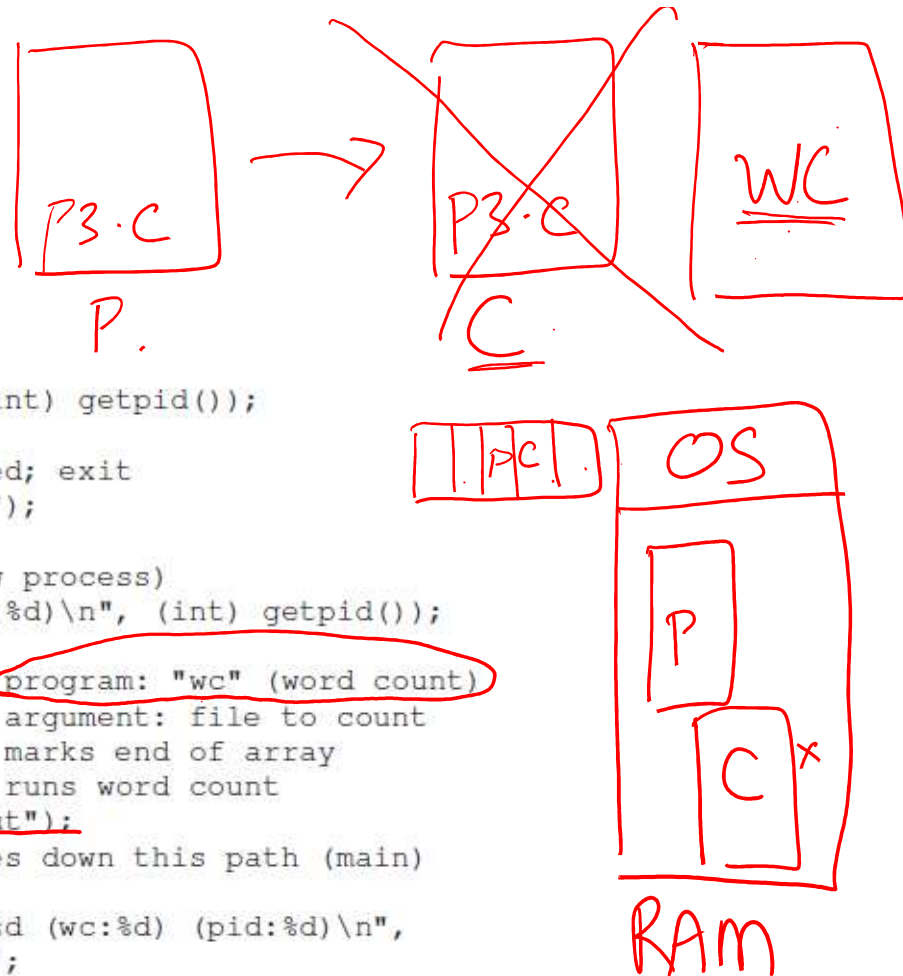
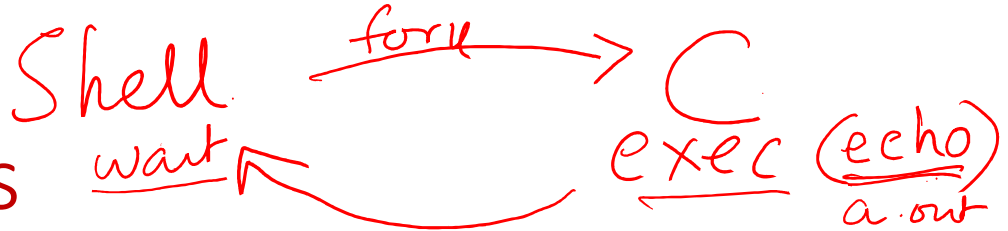


Figure 5.3: Calling `fork()`, `wait()`, And `exec()` (p3.c)



~~sleep 5 &~~

# How the shell works



- OS exposes a terminal/shell to run user programs
  - Can be created by first “init” process on boot up
- What happens when you type a command in the shell?
  - Shell runs command, returns back to command prompt again
- How does the shell work?
  - Shell reads input from user
  - Shell process **forks** a child process
  - Child process runs **exec** with “echo” program executable as argument (most Linux commands are programs written already for your convenience)
  - Child runs “echo” command, calls **exit** at end of program
  - Parent shell calls **wait**, blocks till child terminates, reaps it
  - Once child is done, reads next input command from user
- Think: why doesn't shell exec command directly?
  - Do we want the shell program code to be rewritten fully?

`$ echo hello`  
hello  
`$`

Shell.c

```
do forever {
    input(command)
    exec(...)
    int ret = fork()
    if(ret == 0) {
        exec(command)
    }
    else {
        wait()
    }
}
```



## More on Shell

STD IN ←  
out → Screen  
err → Screen

Shell → C  
(stdout)  
↓  
file  
exec

- Shell can manipulate the child in strange ways
- Suppose you want to redirect output from a command to a file

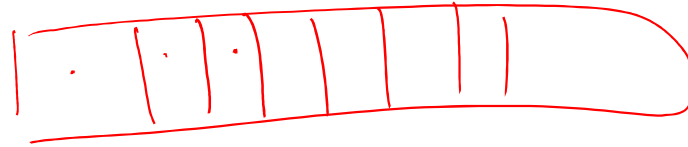
`$echo hello > foo.txt`

- Shell spawns a child, rewires its standard output to a file, then calls exec on the child

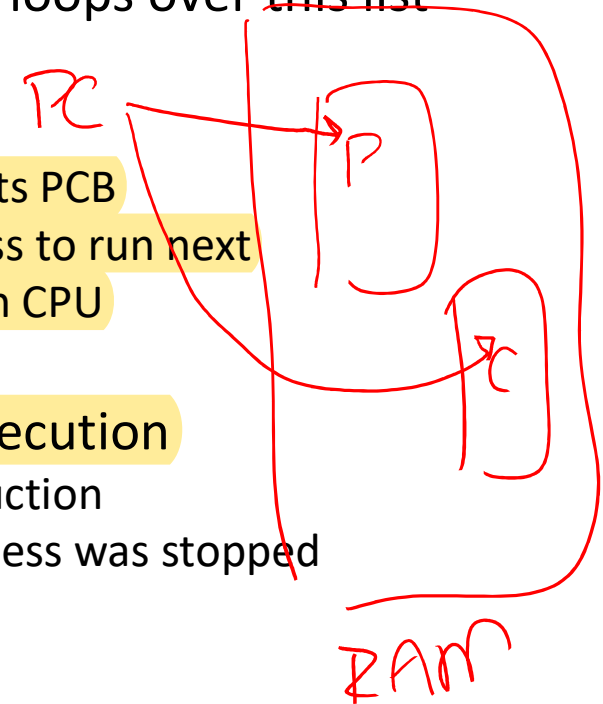
```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <string.h>
5  #include <fcntl.h>
6  #include <sys/wait.h>
7
8  int
9  main(int argc, char *argv[])
10 {
11     int rc = fork();
12     if (rc < 0) { // fork failed; exit
13         fprintf(stderr, "fork failed\n");
14         exit(1);
15     } else if (rc == 0) { // child: redirect standard output to a file
16         close(STDOUT_FILENO);
17         open("./p4.output", O_CREAT|O_WRONLY|O_TRUNC, S_IRWXU);
18         // now exec "wc"...
19         char *myargs[3];
20         myargs[0] = strdup("wc"); // program: "wc" (word count)
21         myargs[1] = strdup("p4.c"); // argument: file to count
22         myargs[2] = NULL; // marks end of array
23         execvp(myargs[0], myargs); // runs word count
24     } else { // parent goes down this path (main)
25         int wc = wait(NULL);
26     }
27     return 0;
28 }
29 }
```

Figure 5.4: All Of The Above With Redirection (p4.c)

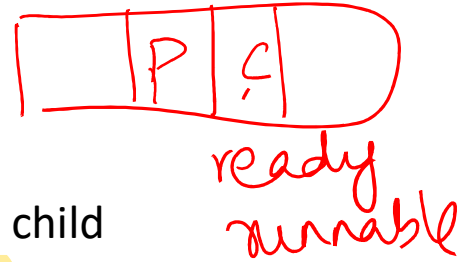
# OS scheduler



- OS maintains list of all active processes (PCBs) in a data structure
  - Processes added during fork, removed after clean up in wait
- OS **scheduler** is special code in the OS that periodically loops over this list and picks processes to run
- Basic outline of scheduler code
  - When invoked, save context of currently running process in its PCB
  - Loop over all ready/runnable processes and identify a process to run next
  - Restore context of new process from PCB and get it to run on CPU
  - Repeat this process as long as system is running
- Note that restoring context of a process resumes its execution
  - PC points to instruction in process code, starts running instruction
  - Other registers are filled with values that existed before process was stopped
  - Process continues execution without realizing it was paused



# xv6: fork system call implementation



- 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

fork  
int-  
OS

P → fork()

PCB entry

```

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;
    
```

```

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
2614 acquire(&ptable.lock);
2615
2616 np->state = RUNNABLE;
2617
2618 release(&ptable.lock);
2619
2620 return pid;
2621 }
    
```

ready

# xv6: exit system call implementation

- Exiting process cleans up state (e.g., close files)
- Pass abandoned children (orphans) to init
- Mark itself as zombie and invoke scheduler



```
2626 void
2627 exit(void)
2628 {
2629     struct proc *curproc = myproc();
2630     struct proc *p;
2631     int fd;
2632
2633     if(curproc == initproc)
2634         panic("init exiting");
2635
2636     // Close all open files.
2637     for(fd = 0; fd < NOFILE; fd++){
2638         if(curproc->ofile[fd]){
2639             fileclose(curproc->ofile[fd]);
2640             curproc->ofile[fd] = 0;
2641         }
2642     }
2643
2644     begin_op();
2645     iput(curproc->cwd);
2646     end_op();
2647     curproc->cwd = 0;
2648
2649     acquire(&ptable.lock);
```

```
2650     // Parent might be sleeping in wait().
2651     wakeup1(curproc->parent);
2652
2653     // Pass abandoned children to init.
2654     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2655         if(p->parent == curproc){
2656             p->parent = initproc;
2657             if(p->state == ZOMBIE)
2658                 wakeup1(initproc);
2659         }
2660     }
2661
2662     // Jump into the scheduler, never to return.
2663     curproc->state = ZOMBIE;
2664     sched();
2665     panic("zombie exit");
2666 }
```

# xv6: wait system call implementation

```
2670 int
2671 wait(void)
2672 {
2673     struct proc *p;
2674     int havekids, pid;
2675     struct proc *curproc = myproc();
2676     acquire(&ptable.lock);
2677     for(;;){
2678         // Scan through table looking for exited children.
2679         havekids = 0;
2680         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2681             if(p->parent != curproc)
2682                 continue;
2683             havekids = 1;
2684             if(p->state == ZOMBIE){
2685                 // Found one.
2686                 pid = p->pid;
2687                 kfree(p->kstack);
2688                 p->kstack = 0;
2689                 freevm(p->pgdir);
2690                 p->pid = 0;
2691                 p->parent = 0;
2692                 p->name[0] = 0;
2693                 p->killed = 0;
2694                 p->state = UNUSED;
2695                 release(&ptable.lock);
2696                 return pid;
2697             }
2698         }
2699     }
```

```
2700 // No point waiting if we don't have any children.
2701 if(!havekids || curproc->killed){
2702     release(&ptable.lock);
2703     return -1;
2704 }
2705 // Wait for children to exit. (See wakeup1 call in proc_exit.)
2706 sleep(curproc, &ptable.lock);
2707 }
2708 }
2709 }
```

*clean up*

- Search for dead children in process table
- If dead child found, clean up memory of zombie, return PID of dead child
- If no dead child, sleep until one dies

## xv6: main function of shell

```
8700 int
8701 main(void)
8702 {
8703     static char buf[100];
8704     int fd;
8705
8706     // Ensure that three file descriptors are open.
8707     while((fd = open("console", O_RDWR)) >= 0){
8708         if(fd >= 3){
8709             close(fd);
8710             break;
8711         }
8712     }
8713
8714     // Read and run input commands.
8715     while(getcmd(buf, sizeof(buf)) >= 0){
8716         if(buf[0] == 'c' && buf[1] == 'd' && buf[2] == ' '){
8717             // Chdir must be called by the parent, not the child.
8718             buf[strlen(buf)-1] = 0; // chop \n
8719             if(chdir(buf+3) < 0)
8720                 printf(2, "cannot cd %s\n", buf+3);
8721             continue;
8722         }
8723         if(fork1() == 0)
8724             runcmd(parsecmd(buf)); exec
8725         wait();
8726     }
8727     exit();
8728 }
```