

Lecture 26: Process creation in xv6

Mythili Vutukuru

IIT Bombay

<https://www.cse.iitb.ac.in/~mythili/os/>

New process creation in xv6

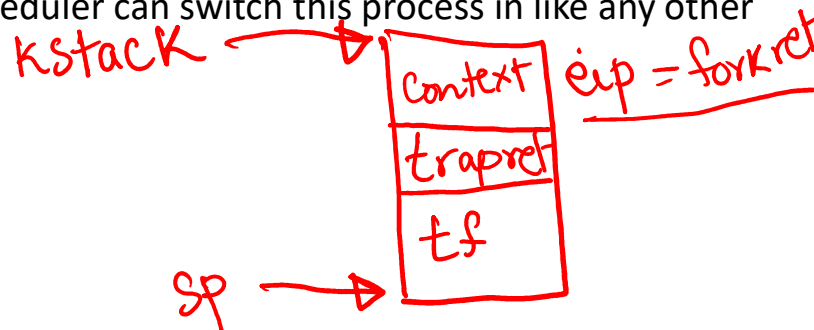
init \longrightarrow *Shell* \longrightarrow *user process*

- Init process: first process created by xv6 after boot up
 - This init process forks shell process, which in turn forks other processes to run user commands
 - The init process is the ancestor of all processes in Unix-like systems
- After init, every other process is created by the fork system call, where a parent forks/spawns a child process
- The function “allocproc” called during both init process creation and in fork system call
 - Allocates new process structure, PID etc
 - Sets up the kernel stack of process so that it is ready to be context switched in by scheduler

allocproc

- 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 (more later)
- 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

make it look like it had a trap and it was context switched out in the past.



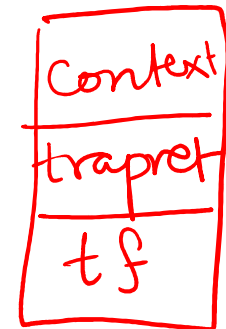
```

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 }

```

even though the one-time initialization runs only once, every new process must pass through forkret to safely transition from the scheduler to user execution

Init process creation



- Alloc proc has created new process
 - When scheduled, it runs function forkret, then trapret
- Trapframe of process set to make process return to first instruction of init code (initcode.S) in userspace
- The code “initcode.S” simply performs “exec” system call to run the init program

```
2518 // Set up first user process.
2519 void
2520 userinit(void)
2521 {
2522     struct proc *p;
2523     extern char _binary_initcode_start[], _binary_initcode_size[];
2524
2525     p = allocproc();
2526
2527     initproc = p;
2528     if((p->pgdir = setupkvm()) == 0)
2529         panic("userinit: out of memory?");
2530     inituvm(p->pgdir, _binary_initcode_start, (int)_binary_initcode_size);
2531     p->sz = PGSIZE;
2532     memset(p->tf, 0, sizeof(*p->tf));
2533     p->tf->cs = (SEG_UCODE << 3) | DPL_USER;
2534     p->tf->ds = (SEG_UDATA << 3) | DPL_USER;
2535     p->tf->es = p->tf->ds;
2536     p->tf->ss = p->tf->ds;
2537     p->tf->eflags = FL_IF;
2538     p->tf->esp = PGSIZE;
2539     p->tf->eip = 0; // beginning of initcode.S
2540
2541     safestrncpy(p->name, "initcode", sizeof(p->name));
2542     p->cwd = namei("/");
2543
2544     // this assignment to p->state lets other cores
2545     // run this process. the acquire forces the above
2546     // writes to be visible, and the lock is also needed
2547     // because the assignment might not be atomic.
2548     acquire(&ptable.lock);
2549
2550     p->state = RUNNABLE;
2551
2552     release(&ptable.lock);
2553 }
```

Init process

- Init program opens STDIN, STDOUT, STDERR files
 - Inherited by all subsequent processes as child inherits parent's files
- Forks a child, **execs shell executable in the child, waits for child to die**
- **Reaps dead children (its own or other orphan descendants)**

```
8500 // init: The initial user-level program
8501
8502 #include "types.h"
8503 #include "stat.h"
8504 #include "user.h"
8505 #include "fcntl.h"
8506
8507 char *argv[] = { "sh", 0 };
8508
8509 int
8510 main(void)
8511 {
8512     int pid, wpid;
8513
8514     if(open("console", O_RDWR) < 0){
8515         mknod("console", 1, 1);
8516         open("console", O_RDWR);
8517     }
8518     dup(0); // stdout
8519     dup(0); // stderr
8520
8521     for(;;){
8522         printf(1, "init: starting sh\n");
8523         pid = fork();
8524         if(pid < 0){
8525             printf(1, "init: fork failed\n");
8526             exit();
8527         }
8528         if(pid == 0){
8529             exec("sh", argv);
8530             printf(1, "init: exec sh failed\n");
8531             exit();
8532         }
8533         while((wpid=wait()) >= 0 && wpid != pid)
8534             printf(1, "zombie!\n");
8535     }
8536 }
```

in xv6 init process forks the shell which is used to fork the other processes.

Forking new process

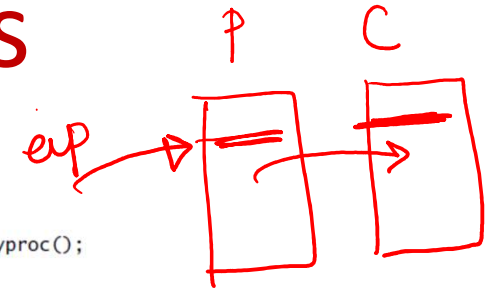
- Fork allocates new process via `allocproc`
- Parent memory and file descriptors copied (more later)
- Trapframe of child copied from that of parent
 - Result: child returns from trap to exact line of code as parent
 - Different physical memory but same virtual address (location in code)
 - Only return 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 from `trapret`.

in Linux even the physical memory is same until some write call is made.

```

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(&ptable.lock);
2614     np->state = RUNNABLE;
2615     release(&ptable.lock);
2616
2617     return pid;
2621 }
2622

```



Summary of new process creation

- New process created by marking a new entry in ptable as RUNNABLE, after configuring the kernel stack, memory image etc of new process
- Neat hack: kernel stack of new process made to look like that of a process that had been context switched out in the past, so that scheduler can context switch it in like any other process
 - No special treatment for newly forked process during “swtch”