

Архитектура операционной системы

Process scheduling

Что такое процесс?



Что такое процесс?

- Иллюзия персональной абстрактной машины
 - Адресное пространство
 - Страницы, таблицы
 - CPU
 - Файлы, и другие абстракции ОС
 - Состояние?
 - Стек?
- `proc.h:52`, `sched.h`

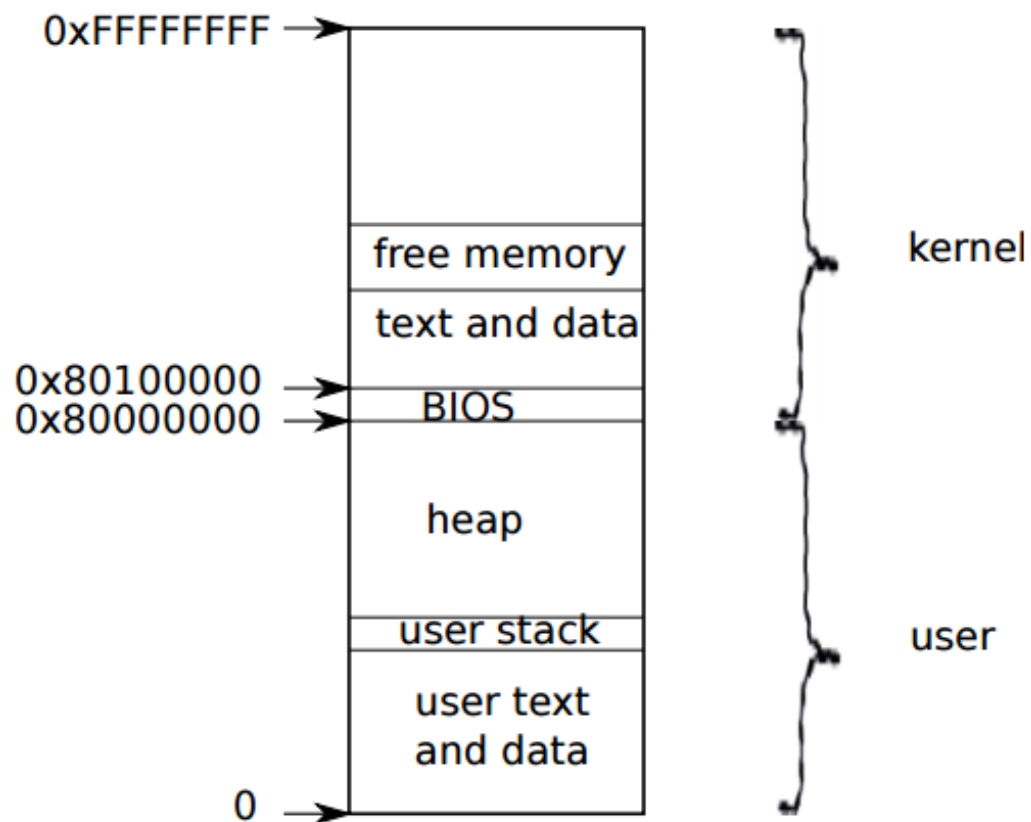
Process context (xv6)

```
27 struct context {
28     uint edi;
29     uint esi;
30     uint ebx;
31     uint ebp;
32     uint eip;
33};
```

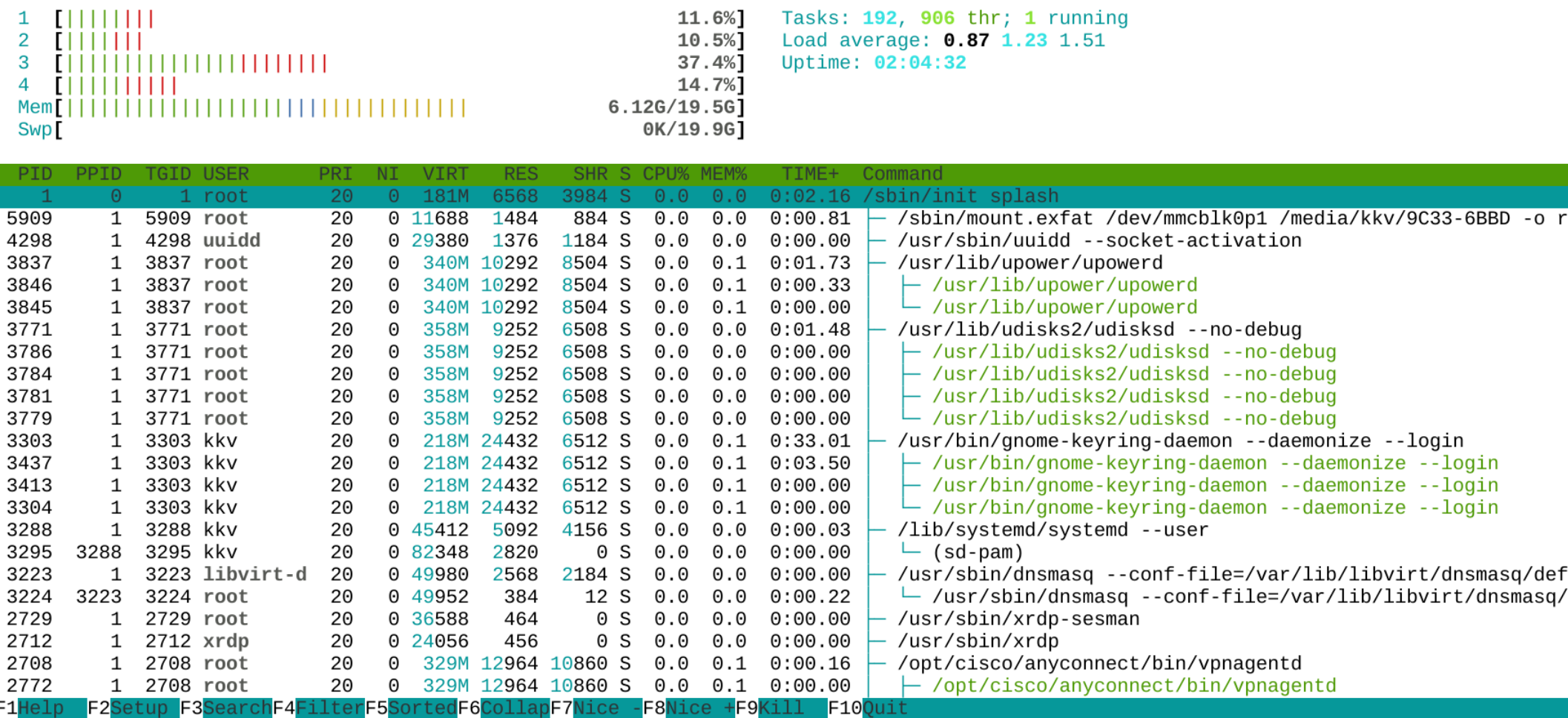
```
enum procstate { UNUSED, EMBRYO, SLEEPING, RUNNABLE, RUNNING, ZOMBIE };
```

```
37 // Per-process state
38 struct proc {
39     uint sz;                // Size of process memory (bytes)
40     pde_t* pgdir;           // Page table
41     char *kstack;           // Bottom of kernel stack for this process
42     enum procstate state;    // Process state
43     int pid;                 // Process ID
44     struct proc *parent;     // Parent process
45     struct trapframe *tf;    // Trap frame for current syscall
46     struct context *context; // switch() here to run process
47     void *chan;              // If non-zero, sleeping on chan
48     int killed;               // If non-zero, have been killed
49     struct file *ofile[NOFILE]; // Open files
50     struct inode *cwd;        // Current directory
51     char name[16];           // Process name (debugging)
52};
```

Память процесса (layout)



Что такое поток?



Fork/exec/clone

Code analysis

- main.c
- bootasm.S
- proc.c
- vm.c
- main.c +54

OS Main

```
18 main(void)
19 {
20     kinit1(end, P2V(4*1024*1024)); // phys page allocator
21     kvmalloc(); // kernel page table
22     mpinit(); // detect other processors
23     lapicinit(); // interrupt controller
24     seginit(); // segment descriptors
25     picinit(); // disable pic
26     ioapicinit(); // another interrupt controller
27     consoleinit(); // console hardware
28     uartinit(); // serial port
29     pinit(); // process table
30     tvinit(); // trap vectors
31     binit(); // buffer cache
32     fileinit(); // file table
33     ideinit(); // disk
34     startothers(); // start other processors
35     kinit2(P2V(4*1024*1024), P2V(PHYSTOP)); // must come after startothers()
36     userinit(); // first user process
37     mpmain(); // finish this processor's setup
38 }
```

```

121 userinit(void)
122 {
123     struct proc *p;
124     extern char _binary_initcode_start[], _binary_initcode_size[];
125
126     p = allocproc();
127
128     initproc = p;
129     if((p->pgdir = setupkvm()) == 0)
130         panic("userinit: out of memory?");
131     inituvm(p->pgdir, _binary_initcode_start, (int)_binary_initcode_size);
132     p->sz = PGSIZE;
133     memset(p->tf, 0, sizeof(*p->tf));
134     p->tf->cs = (SEG_UCODE << 3) | DPL_USER;
135     p->tf->ds = (SEG_UDATA << 3) | DPL_USER;
136     p->tf->es = p->tf->ds;
137     p->tf->ss = p->tf->ds;
138     p->tf->eflags = FL_IF;
139     p->tf->esp = PGSIZE;
140     p->tf->eip = 0; // beginning of initcode.S
141
142     safestrcpy(p->name, "initcode", sizeof(p->name));
143     p->cwd = namei("/");
144
145     // this assignment to p->state lets other cores
146     // run this process. the acquire forces the above
147     // writes to be visible, and the lock is also needed
148     // because the assignment might not be atomic.
149     acquire(&ptable.lock);
150
151     p->state = RUNNABLE;
152
153     release(&ptable.lock);
154 }

```

Создание контекста процесса

Создание памяти процесса

```

181 // Load the initcode into address 0 of pgdir.
182 // sz must be less than a page.
183 void
184 inituvm(pde_t *pgdir, char *init, uint sz)
185 {
186     char *mem;
187
188     if(sz >= PGSIZE)
189         panic("inituvm: more than a page");
190     mem = kalloc();
191     memset(mem, 0, PGSIZE);
192     mappages(pgdir, 0, PGSIZE, V2P(mem), PTE_W|PTE_U);
193     memmove(mem, init, sz);
194 }

```

Запуск

```

10
11 // Interrupt descriptor table (shared by all CPUs).
12 struct gatedesc idt[256];
13 extern uint vectors[]; // in vectors.S: array of 256 entry pointers
14 struct spinlock tickslock;
15 uint ticks;

```

trap.c

```

29 void
30 idtinit(void)
31 {
32     lidt(idt, sizeof(idt));
33 }

```

```

50// Common CPU setup code.
51static void
52mpmain(void)
53{
54     cprintf("cpu%d: starting %d\n", cpuid(), cpuid());
55     idtinit(); // load idt register
56     xchg(&(mycpu()->started), 1); // tell startothers() we're up
57     scheduler(); // start running processes
58}

```

```

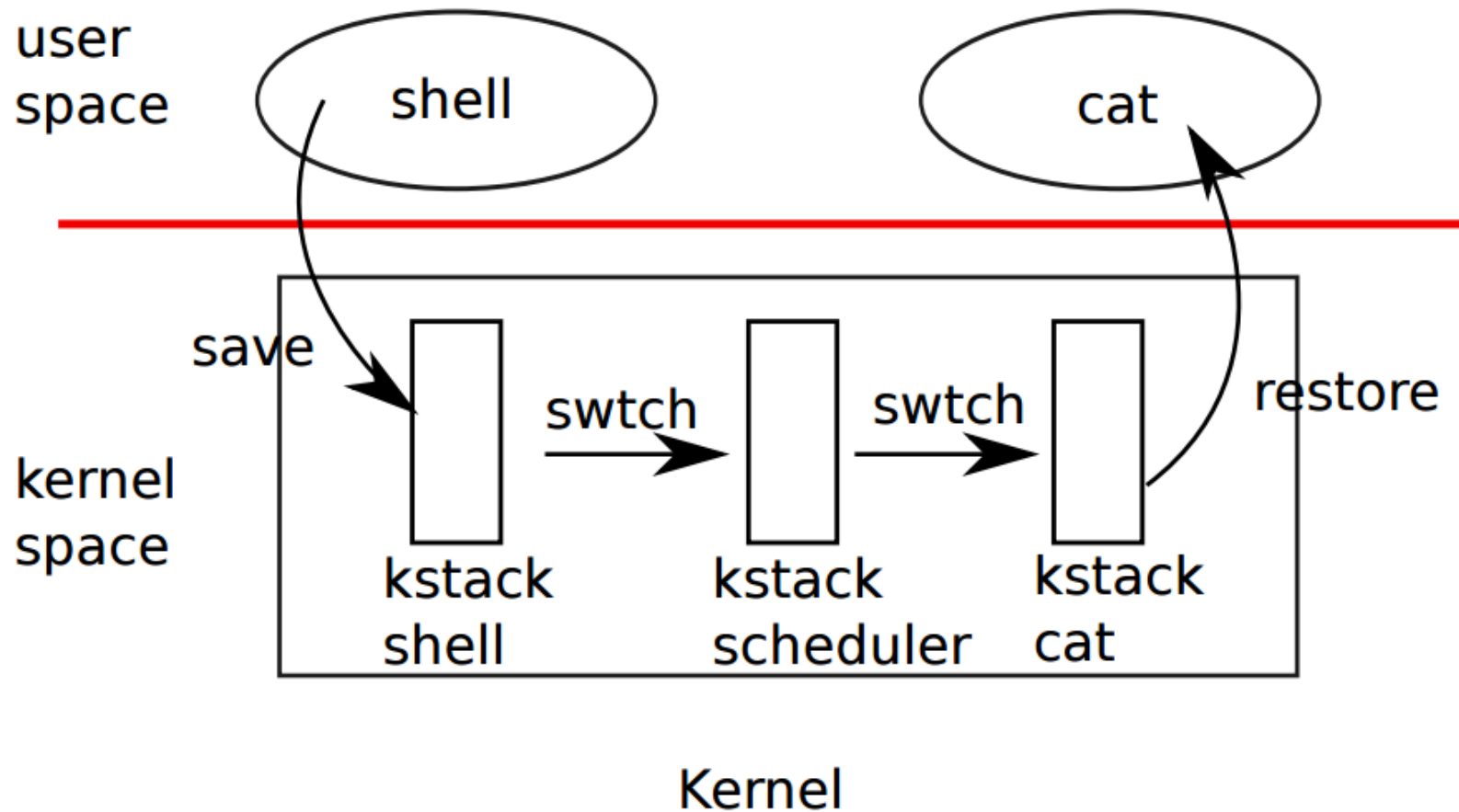
280 // Per-CPU process scheduler.
281 // Each CPU calls scheduler() after setting itself up.
282 // Scheduler never returns. It loops, doing:
283 // - choose a process to run
284 // - switch to start running that process
285 // - eventually that process transfers control
286 //   via switch back to the scheduler.
287 void
288 scheduler(void)
289 {
290     struct proc *p;
291
292     for(;;){
293         // Enable interrupts on this processor.
294         sti();
295
296         // Loop over process table looking for process to run.
297         acquire(&ptable.lock);
298         for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
299             if(p->state != RUNNABLE)
300                 continue;
301
302             // Switch to chosen process. It is the process's job
303             // to release ptable.lock and then reacquire it
304             // before jumping back to us.
305             proc = p;
306             switchvm(p);
307             p->state = RUNNING;
308             swtch(&cpu->scheduler, p->context);
309             switchkvm();
310
311             // Process is done running for now.
312             // It should have changed its p->state before coming back.
313             proc = 0;
314         }
315         release(&ptable.lock);
316     }
317 }
318 }

```

Много[задачность]

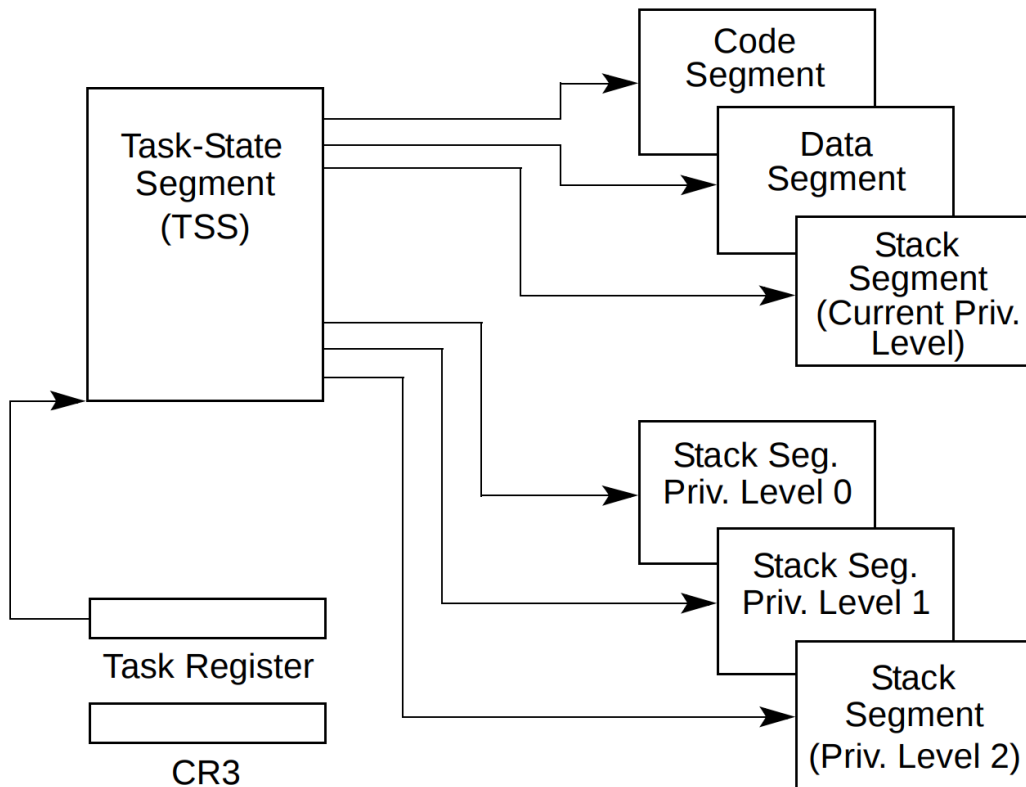
- Прерывание таймера
- Смена контекста
- План блокировок, при наличии нескольких CPU
- Освобождение ресурсов при завершении процесса

Переключение контекста



switchkvm/switchuvmm

```
147// Switch h/w page table register to the kernel-only page table,
148// for when no process is running.
149void
150switchkvm(void)
151{
152    lcr3(V2P(kpgdir)); // switch to the kernel page table
153}
```



```
155// Switch TSS and h/w page table to correspond to process p.
156void
157switchuvmm(struct proc *p)
158{
159    if(p == 0)
160        panic("switchuvmm: no process");
161    if(p->kstack == 0)
162        panic("switchuvmm: no kstack");
163    if(p->pgdir == 0)
164        panic("switchuvmm: no pgdir");
165
166    pushcli();
167    mycpu()->gdt[SEG_TSS] = SEG16(STS_T32A, &mycpu()->ts,
168                                sizeof(mycpu()->ts)-1, 0);
169    mycpu()->gdt[SEG_TSS].s = 0;
170    mycpu()->ts.ss0 = SEG_KDATA << 3;
171    mycpu()->ts.esp0 = (uint)p->kstack + KSTACKSIZE;
172    // setting IOPL=0 in eflags *and* iomb beyond the tss segment limit
173    // forbids I/O instructions (e.g., inb and outb) from user space
174    mycpu()->ts.iomb = (ushort) 0xFFFF;
175    ltr(SEG_TSS << 3);
176    lcr3(V2P(p->pgdir)); // switch to process's address space
177    popcli();
178}
```

CPU state

```
2053 // Per-CPU state
2054 struct cpu {
2055     uchar id;                    // Local APIC ID; index into cpus[] below
2056     struct context *scheduler;   // swtch() here to enter scheduler
2057     struct taskstate ts;         // Used by x86 to find stack for interrupt
2058     struct segdesc gdt[NSEGS];   // x86 global descriptor table
2059     volatile uint started;       // Has the CPU started?
2060     int ncli;                    // Depth of pushcli nesting.
2061     int intena;                  // Were interrupts enabled before pushcli?
2062
2063     // Cpu-local storage variables; see below
2064     struct cpu *cpu;
2065     struct proc *proc;           // The currently-running process.
2066 };
```

Переключение контекста (+1)

- Замечания:
 - Процесс:
 - Свой набор регистров
- ```
2093 struct context {
2094 uint edi;
2095 uint esi;
2096 uint ebx;
2097 uint ebp;
2098 uint eip;
2099 };
```
- Свой стек в ядре
- **Каждый CPU имеет свой scheduler thread**

```
365 void
366 sched(void)
367 {
368 int intena;
369 struct proc *p = myproc();
370
371 if(!holding(&ptable.lock))
372 panic("sched ptable.lock");
373 if(mycpu()->ncli != 1)
374 panic("sched locks");
375 if(p->state == RUNNING)
376 panic("sched running");
377 if(readeflags() & FL_IF)
378 panic("sched interruptible");
379 intena = mycpu()->intena;
380 swtch(&p->context, mycpu()->scheduler);
381 mycpu()->intena = intena;
382 }
```



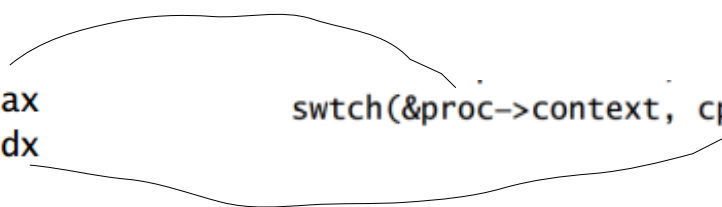
```

2457 void
2458 scheduler(void)
2459 {
2460 struct proc *p;
2461
2462 for(;;){
2463 // Enable interrupts on this processor.
2464 sti();
2465
2466 // Loop over process table looking for process to run.
2467 acquire(&ptable.lock);
2468 for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2469 if(p->state != RUNNABLE)
2470 continue;
2471
2472 // Switch to chosen process. It is the process's job
2473 // to release ptable.lock and then reacquire it
2474 // before jumping back to us.
2475 proc = p;
2476 switchvm(p);
2477 p->state = RUNNING;
2478 swtch(&cpu->scheduler, proc->context);
2479 switchkvm();
2480
2481 // Process is done running for now.
2482 // It should have changed its p->state before coming back.
2483 proc = 0;
2484 }
2485 release(&ptable.lock);

```

# swtch

```
2700 # Context switch
2701 #
2702 # void swtch(struct context **old, struct context *new);
2703 #
2704 # Save current register context in old
2705 # and then load register context from new.
2706
2707 .globl swtch
2708 swtch:
2709 movl 4(%esp), %eax swtch(&proc->context, cpu->scheduler);
2710 movl 8(%esp), %edx
2711
2712 # Save old callee-save registers
2713 pushl %ebp
2714 pushl %ebx
2715 pushl %esi
2716 pushl %edi
2717
2718 # Switch stacks
2719 movl %esp, (%eax)
2720 movl %edx, %esp
2721
2722 # Load new callee-save registers
2723 popl %edi
2724 popl %esi
2725 popl %ebx
2726 popl %ebp
2727 ret
```



## syscall.h

```
2#define SYS_fork 1
3#define SYS_exit 2
4#define SYS_wait 3
5#define SYS_pipe 4
6#define SYS_read 5
7#define SYS_kill 6
8#define SYS_exec 7
9#define SYS_fstat 8
10#define SYS_chdir 9
11#define SYS_dup 10
12#define SYS_getpid 11
13#define SYS_sbrk 12
14#define SYS_sleep 13
15#define SYS_uptime 14
16#define SYS_open 15
17#define SYS_write 16
18#define SYS_mknod 17
19#define SYS_unlink 18
20#define SYS_link 19
21#define SYS_mkdir 20
22#define SYS_close 21
```

## syscall.c

```
85 extern int sys_chdir(void);
86 extern int sys_close(void);
87 extern int sys_dup(void);
88 extern int sys_exec(void);
89 extern int sys_exit(void);
90 extern int sys_fork(void);
91 extern int sys_fstat(void);
92 extern int sys_getpid(void);
93 extern int sys_kill(void);
94 extern int sys_link(void);
95 extern int sys_mkdir(void);
96 extern int sys_mknod(void);
97 extern int sys_open(void);
98 extern int sys_pipe(void);
99 extern int sys_read(void);
100 extern int sys_sbrk(void);
101 extern int sys_sleep(void);
102 extern int sys_unlink(void);
103 extern int sys_wait(void);
104 extern int sys_write(void);
105 extern int sys_uptime(void);
106
107 static int (*syscalls[])(void) = {
108 [SYS_fork] sys_fork,
109 [SYS_exit] sys_exit,
110 [SYS_wait] sys_wait,
111 [SYS_pipe] sys_pipe,
112 [SYS_read] sys_read,
113 [SYS_kill] sys_kill,
114 [SYS_exec] sys_exec,
115 [SYS_fstat] sys_fstat,
116 [SYS_chdir] sys_chdir,
117 [SYS_dup] sys_dup,
118 [SYS_getpid] sys_getpid,
119 [SYS_sbrk] sys_sbrk,
120 [SYS_sleep] sys_sleep,
121 [SYS_uptime] sys_uptime,
122 [SYS_open] sys_open,
123 [SYS_write] sys_write,
124 [SYS_mknod] sys_mknod,
125 [SYS_unlink] sys_unlink,
126 [SYS_link] sys_link,
127 [SYS_mkdir] sys_mkdir,
128 [SYS_close] sys_close,
129 };
```

## sysproc.c

```
10 int
11 sys_fork(void)
12 {
13 return fork();
14 }
```

```
177 // Create a new process copying p as the parent.
178 // Sets up stack to return as if from system call.
179 // Caller must set state of returned proc to RUNNABLE.
180 int
181 fork(void)
182 {
183 int i, pid;
184 struct proc *np;
185 struct proc *curproc = myproc();
186
187 // Allocate process.
188 if((np = allocproc()) == 0){
189 return -1;
190 }
191
192 // Copy process state from curproc.
193 if((np->pgdir = copyvm(curproc->pgdir, curproc->sz)) == 0){
194 kfree(np->kstack);
195 np->kstack = 0;
196 np->state = UNUSED;
197 return -1;
198 }
199 np->sz = curproc->sz;
200 np->parent = curproc;
201 *np->tf = *curproc->tf;
202
203 // Clear %eax so that fork returns 0 in the child.
204 np->tf->eax = 0;
205
206 for(i = 0; i < NOFILE; i++)
207 if(curproc->ofile[i])
208 np->ofile[i] = filedup(curproc->ofile[i]);
209 np->cwd = idup(curproc->cwd);
210
211 safestrcpy(np->name, curproc->name, sizeof(curproc->name));
212
213 pid = np->pid;
214
215 acquire(&table.lock);
216
217 np->state = RUNNABLE;
218
219 release(&table.lock);
220
221 return pid;
222 }
```

## usys.S

```
1 #include "syscall.h"
2 #include "traps.h"
3
4 #define SYSCALL(name) \
5 .globl name; \
6 name: \
7 movl $SYS_ ## name, %eax; \
8 int $T_SYSCALL; \
9 ret
10
11 SYSCALL(fork)
12 SYSCALL(exit)
13 SYSCALL(wait)
```