

Chapter 3a: xv6 Process Internals

Gene Cooperman

Contents

1	xv6 Process Internals: Data Structures and Operations in <code>proc.h</code> and <code>proc.c</code>	1
1.1	1. Data Structures in <code>proc.h</code>	1
1.2	2. The Process Table Lock	4
1.3	3. Operations in <code>proc.c</code>	4
1.4	4. Context Switching: The Full Picture	8
1.5	5. Summary of Functions	9
1.6	6. The Mechanics of <code>swtch.S</code>	10

1 xv6 Process Internals: Data Structures and Operations in `proc.h` and `proc.c`

This is an explanation of how xv6 represents processes and CPUs, and how the core process-management functions work.

1.1 1. Data Structures in `proc.h`

1.1.1 `struct cpu` – per-CPU state

Each physical CPU in the system has a `struct cpu`. In modern Linux, we might have multiple CPU cores, but not multiple CPUs. The xv6 `struct cpu` contains information that might be associated with both a CPU core and also a thread, rather than a process. In our analysis, we will assume `NCPU = 1`.

For completeness, we show the full `struct cpu` below. But we will be interested only in:

- `struct context *scheduler;` // saved context for the scheduler thread
- `struct proc *proc;` // the process currently running on this CPU, or NULL

the `proc` field is especially important if we are assuming one CPU core. The `proc` is a pointer to the entry in the process table that is the currently running process. There will always be exactly one running process, since if no user process is running, then the scheduler process will be running.

```
struct cpu {
    uchar apicid;           // Local APIC ID (hardware CPU identifier)
    struct context *scheduler; // saved context for the scheduler thread
    struct taskstate ts;    // x86 task state segment (for privilege transitions)
    struct segdesc gdt[NSEGS]; // x86 global descriptor table
    volatile uint started; // has this CPU been initialized?
```

```

int ncli;           // depth of pushcli nesting (interrupt disable count)
int intena;        // were interrupts enabled before the first pushcli?
struct proc *proc; // the process currently running on this CPU, or NULL
};

```

Key points:

- **scheduler** stores the saved register state for the per-CPU scheduler loop. When a process yields the CPU, `swtch()` saves the process's context and restores this one, returning control to `scheduler()`.
- **proc** points to whichever process is currently **RUNNING** on this CPU. Kernel code throughout xv6 accesses the current process through this pointer (historically accessed via the `%gs` segment register for efficiency).

1.1.2 Accessing the current CPU and process

`proc.h` defines two special variables using GCC's `asm` register-binding syntax:

```

extern struct cpu *cpu asm("%gs:0"); // current CPU (Intel assembly, only)
extern struct proc *proc asm("%gs:4"); // current process (Intel assembly, only)

```

Each CPU has its own `%gs` segment pointing to its `struct cpu`. This allows kernel code on any CPU to access `cpu` and `proc` without locks or array lookups – the hardware resolves the correct structure automatically.

1.1.3 `struct context` – saved registers for context switching

```

struct context {
    uint edi;
    uint esi;
    uint ebx;
    uint ebp;
    uint eip;
};

```

This stores the **callee-saved registers** plus the instruction pointer. It is the minimal state needed to suspend and resume a kernel thread. Note that this is *not* the full user-space register state – that is saved in the trapframe. `struct context` is used only for switching between kernel execution contexts: from a process's kernel thread to the scheduler, or vice versa.

The layout matches what `swtch()` (in `swtch.S`) pushes and pops on the kernel stack.

1.1.4 `enum procstate` – process lifecycle states

```

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

```

State	Meaning
UNUSED	Slot in process table is free
EMBRYO	Process is being initialized
RUNNABLE	Ready to run, waiting for scheduler
RUNNING	Currently executing on a CPU
SLEEPING	Blocked, waiting for an event
ZOMBIE	Exited, waiting for parent to call <code>wait()</code>

1.1.5 struct proc – the process control block

```
struct proc {
    uint sz;                // size of process memory (bytes)
    pde_t *pgdir;          // page directory (virtual address space)
    char *kstack;          // bottom of kernel stack for this process
    enum procstate state;  // process state
    int pid;               // process ID
    struct proc *parent;   // parent process
    struct trapframe *tf;  // trapframe for current syscall/interrupt
    struct context *context; // saved kernel context (for swtch())
    void *chan;           // if SLEEPING: the channel being waited on
    int killed;          // if non-zero, process has been killed
    struct file *ofile[NOFILE]; // open files
    struct inode *cwd;    // current working directory
    char name[16];        // process name (for debugging)
};
```

Key fields grouped by purpose:

Memory management: - `sz` – the size of the user address space. Grows via `sbrk()` / `growproc()`. - `pgdir` – the page directory. Each process has its own, giving it a private virtual address space. - `kstack` – each process has a dedicated kernel stack, used when it traps into the kernel.

Scheduling and context switching: - `state` – the current lifecycle state. - `context` – pointer to saved callee-saved registers on the kernel stack, used by `swtch()`. - `tf` – pointer to the trapframe on the kernel stack, used by `trapret` to return to user space. - `chan` – when a process is SLEEPING, this records *what* it's sleeping on (an arbitrary address used as a key by `sleep()`/`wakeup()`).

Process relationships: - `pid` – unique process identifier. - `parent` – pointer to the parent process. Used by `exit()` and `wait()` for process cleanup. - `killed` – a flag set by `kill()`. The process isn't immediately destroyed; instead, it is marked and will check this flag at the next safe point (e.g., returning to user space from a trap) and call `exit()`.

File system: - `ofile[NOFILE]` – array of pointers to open `struct file` entries. The file descriptor `fd` is an index into this `ofile[]` array. - `cwd` – pointer to the inode of the current working directory.

Debugging: - `name` – a human-readable name, typically set to the name of the executable by `exec()`. (In modern Linux, they use `execvp()` or a variant, instead of `exec()`).

1.1.6 ptable – the global process table

Defined in `proc.c`:

```
struct {
    struct spinlock lock;
    struct proc proc[NPROC];
} ptable;
```

All processes in the system live in this fixed-size array (typically `NPROC = 64`). So, this early UNIX allowed at most 64 processes at one time. The spinlock protects concurrent access – any code that reads or modifies process states must hold `ptable.lock`. The process table is the central data structure for process management: `allocproc()` scans it for free slots, the `scheduler()` scans it for runnable processes, `wait()` scans it for zombie children, and so on.

1.2 2. The Process Table Lock

Nearly every function in `proc.c` acquires `ptable.lock` at some point. Because we are assuming `NCPU=1`, there can be at most one simultaneous process in the operating system kernel of xv6. In this special case, we would not need locks. So, feel free to skip this subsection.

This single lock protects:

- Process state transitions (e.g., `RUNNABLE -> RUNNING`)
- The `pid` assignment in `allocproc()`
- The parent/child relationship during `fork()`, `exit()`, and `wait()`
- The `chan` field during `sleep()` and `wakeup()`

This coarse-grained locking is simple but means that process operations are serialized across all CPUs. Real operating systems use finer-grained locking, but xv6 prioritizes clarity over scalability.

1.3 3. Operations in `proc.c`

1.3.1 `pinit()` – initialize the process table (line 24)

Called once at boot (from `main.c`).
Initializes `ptable.lock`.

This is straightforward: just `initlock(&ptable.lock, "ptable")`. After this, the process table is ready for use.

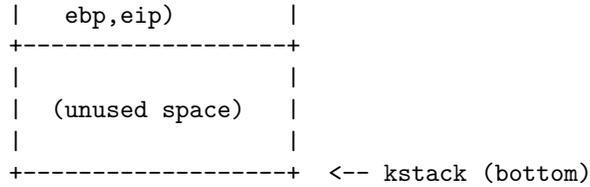
1.3.2 `allocproc()` – allocate a new process slot (called internally)

`allocproc()` is not listed in the file reference because it is `static` (internal to `proc.c`), but it is the foundation for both `userinit()` and `fork()`. It:

1. Acquires `ptable.lock`
2. Scans `ptable.proc[]` for a slot with `state == UNUSED`
3. If found, sets `state = EMBRYO` and assigns `pid = nextpid++`
4. Releases `ptable.lock`
5. Allocates a kernel stack (`KSTACKSIZE` bytes via `kalloc()`)
6. Sets up the bottom of the kernel stack with:
 - A **trapframe** – space for saving user registers on trap entry
 - A **context** – initialized so that `swtch()` will jump to `forkret()`, which then falls into `trapret`

The kernel stack layout after `allocproc()`:

```
+-----+ <-- kstack + KSTACKSIZE (top of stack)
|  trapframe      | <-- proc->tf points here
+-----+
|  trapret addr   | <-- return address for forkret()
+-----+
|  context        | <-- proc->context points here
|  (edi,esi,ebx,  |      eip = forkret
```



This synthetic stack setup is the mechanism that allows a newly created process to be “switched to” by the scheduler as if it were a process that had previously been switched away from. The first `swtch()` into this process restores the context, “returns” to `forkret()`, which “returns” to `trapret`, which pops the trapframe and enters user space.

1.3.3 `userinit()` – create the first user process (line 79)

Called once at boot (from `main.c`), after `pinit()`.
Creates the `init` process (pid 1).

Since there is no parent process to `fork()` from, `userinit()` manually constructs the first process:

1. Calls `allocproc()` to get a new process slot and kernel stack
2. Calls `setupkvm()` to create a page table with kernel mappings
3. Copies a small piece of code (`initcode.S`) into the process’s user memory at address 0 – this code will call `exec("/init", ...)`
4. Sets up the trapframe: `eip` points to address 0 (the `initcode`), `esp` points to `PGSIZE` (one page of stack), segment registers set for user mode
5. Sets `cwd` to the root directory inode
6. Sets `state = RUNNABLE`

When the scheduler first picks this process, it enters user space running `initcode.S`, which immediately `exec()`s `/init`, which in turn `fork()`s a shell.

1.3.4 `fork()` – create a child process (line 129)

This is described in detail in the [“xv6 System Call Flow: `fork\(\)`” document](#). In brief (assuming that we have already entered kernel mode):

1. `allocproc()` – get a new slot and kernel stack
2. `copyvm()` – duplicate the parent’s entire address space
3. Copy the trapframe, but set `eax = 0` (child’s return value)
4. `filedup()` each open file descriptor (shared, not copied)
5. `idup()` the current working directory
6. Copy the process name
7. Set `state = RUNNABLE` and return the child’s pid

1.3.5 `growproc(n)` – grow or shrink the address space (line 108)

Called by `sys_sbrk()`.

Adjusts `proc->sz` by `n` bytes (positive = grow, negative = shrink).

Calls `allocvm()` (to grow) or `deallocvm()` (to shrink) to adjust the page table, then reloads the page table into the hardware via `switchvm()`. Returns 0 on success, -1 on failure (e.g., out of memory).

1.3.6 `exit()` – terminate the current process (line 173)

When a process calls `exit()`:

1. Closes all open files (`fclose()` on each `ofile[]` entry)
2. Releases the current working directory inode (`iput()`)
3. Acquires `ptable.lock`
4. Wakes up the parent (who may be blocked in `wait()`)
5. **Reparents orphaned children:** any children of the exiting process are reassigned to the `init` process (pid 1). If any of those children are zombies, `init` is woken up so it can reap them.
6. Sets `state = ZOMBIE`
7. Calls `sched()` to switch back to the scheduler – **this process never runs again**

The process's kernel stack, page table, and `ptable` entry are *not* freed here. That is the parent's responsibility, in `wait()`.

1.3.7 `wait()` – wait for a child to exit (line 217)

1. Acquires `ptable.lock`
2. Scans the entire `ptable` for a child process (any process whose `parent == proc`) in state `ZOMBIE`
3. If a zombie child is found:
 - Frees the child's kernel stack (`kfree()`)
 - Frees the child's page table (`freevm()`)
 - Sets the child's state to `UNUSED` (the slot is now free for reuse)
 - Returns the child's pid
4. If children exist but none are zombies yet, calls `sleep()` to block until a child exits
5. If no children exist at all, returns -1

This is the other half of process cleanup. `exit()` makes the process a zombie; `wait()` reaps it and frees its resources.

1.3.8 `scheduler()` – the per-CPU scheduling loop (line 266)

Each CPU runs an independent `scheduler()` loop (entered from `main.c` after boot, and never returns). The algorithm is simple round-robin:

```
loop forever:
    enable interrupts (so this CPU isn't deaf to devices)
    acquire ptable.lock
    for each process p in ptable:
        if p->state == RUNNABLE:
            switch to p:
                set cpu->proc = p
                switchvm(p)           // load p's page table
                p->state = RUNNING
                swtch(&cpu->scheduler, p->context) // context switch!
                // ... execution returns here when p yields ...
                switchkvm()           // switch back to kernel page table
                cpu->proc = NULL
    release ptable.lock
```

The `swtch()` call saves the scheduler's registers and loads the process's saved context. The scheduler doesn't run again until the process calls `sched()` (via `yield()`, `sleep()`, or `exit()`).

1.3.9 sched() – switch from a process back to the scheduler (line 301)

sched() is the counterpart to the swtch() call in scheduler(). It performs sanity checks and then calls:
swtch(&proc->context, cpu->scheduler);

This saves the current process's kernel registers and restores the scheduler's context, resuming the scheduler() loop right after the swtch() call that originally entered this process.

sched() requires that ptable.lock is held (to protect the state transition) and that interrupts are disabled.

1.3.10 yield() – voluntarily give up the CPU (line 320)

```
void yield(void) {
    acquire(&ptable.lock);
    proc->state = RUNNABLE;
    sched();
    release(&ptable.lock);
}
```

Called from trap() on every timer interrupt. The process is still runnable – it just gives other processes a chance to execute. This is how xv6 implements **preemptive multitasking**: even if a process never voluntarily sleeps, the timer interrupt forces it to yield periodically.

1.3.11 sleep(chan, lk) – block until woken (line 352)

sleep() puts the current process to sleep on a **channel** – an arbitrary pointer value used as a key to match sleepers with wakeups:

1. Acquires ptable.lock (needed to change state atomically)
2. Releases the caller's lock lk (this is critical – see below)
3. Sets proc->chan = chan and proc->state = SLEEPING
4. Calls sched() to switch to the scheduler
5. When eventually woken up: clears chan, re-acquires lk, and returns

Why release lk before sleeping? The caller always holds some lock lk that protects the condition it's waiting on (e.g., the pipe lock when waiting for data). If the process slept while still holding lk, no other process could ever modify the condition and call wakeup() – deadlock. But simply releasing lk and then sleeping would create a race: the condition might change between the release and the sleep. The solution is to hold ptable.lock across both the release of lk and the state change to SLEEPING. Since wakeup() also requires ptable.lock, no wakeup can be lost.

1.3.12 wakeup(chan) – wake all processes sleeping on a channel (line 401)

```
Acquires ptable.lock.
Scans the entire ptable.
For each process with state == SLEEPING and chan == chan:
    set state = RUNNABLE
Releases ptable.lock.
```

Note that wakeup() wakes *all* sleepers on the channel. This is simple but means that woken processes must re-check their condition (typically in a while loop), because the condition may no longer be true by the time they actually run – another process may have consumed the event first.

1.3.13 kill(pid) – mark a process for termination (line 412)

kill() does **not** immediately destroy the target process. It simply:

1. Scans the ptable for a process with the given pid
2. Sets p->killed = 1
3. If the process is SLEEPING, sets it to RUNNABLE (so it wakes up and can notice it's been killed)

The actual termination happens later: when the process returns from a system call or interrupt, trap() checks proc->killed and calls exit(). This deferred approach is necessary because killing a process while it's in the middle of a file system operation (holding locks, partially updating disk structures) could leave the system in an inconsistent state.

1.3.14 procdump() – print process table for debugging (line 436)

Prints the state, pid, and name of each active process, along with a kernel stack backtrace. Triggered by pressing Ctrl+P on the console. Useful for diagnosing deadlocks and hangs.

1.4 4. Context Switching: The Full Picture

The most subtle aspect of proc.c is how context switching works. Here is the complete path when a running process gives up the CPU:



```

|      (continues scanning ptable for next RUNNABLE process Q)
|      |
|      v
|      acquire ptable.lock
|      Q->state = RUNNING
|      switchvm(Q)
|      |
|      v
|      swtch(&cpu->scheduler, Q->context)
|      |
|      saves scheduler's {edi,esi,ebx,ebp,eip}
|      restores Q's {edi,esi,ebx,ebp,eip}
|      |
|      v
|      Q resumes in sched(), returns to yield()
|      release ptable.lock
|      returns to trap(), returns to user space
|
|
v
(P sits in ptable with state RUNNABLE
until the scheduler picks it again)

```

The key insight is that **every process switch goes through the scheduler**. A process never switches directly to another process. The path is always: process -> scheduler -> (possibly different) process. The `swtch()` function is called twice per transition – once to enter the scheduler, once to leave it.

1.5 5. Summary of Functions

Function	Line	Purpose
<code>pinit()</code>	24	Initialize <code>ptable.lock</code> at boot
<code>allocproc()</code>	(static)	Find unused slot, set up kernel stack
<code>userinit()</code>	79	Create the first user process (init)
<code>growproc(n)</code>	108	Grow or shrink process address space
<code>fork()</code>	129	Create child process (copy of parent)
<code>exit()</code>	173	Terminate current process -> ZOMBIE
<code>wait()</code>	217	Reap zombie child -> UNUSED, free resources
<code>scheduler()</code>	266	Per-CPU loop: pick RUNNABLE, switch to it
<code>sched()</code>	301	Switch from process back to scheduler
<code>yield()</code>	320	Set RUNNABLE, call <code>sched()</code> (preemption)
<code>sleep(chan, lk)</code>	352	Block on channel, release lock atomically
<code>wakeup(chan)</code>	401	Wake all processes sleeping on channel
<code>kill(pid)</code>	412	Mark process for deferred termination
<code>procdump()</code>	436	Debug: print process table to console

Variable	Line	Purpose
<code>ptable</code>	10-13	Global process table (lock + array of <code>NPROC</code> procs)
<code>nextpid</code>	17	Next pid to assign (incremented by <code>allocproc()</code>)

1.6 6. The Mechanics of `swtch.S`

The `swtch()` function in `swtch.S` is the lowest-level building block of context switching. It is a short piece of assembly – only about 10 instructions – but it is the mechanism that makes multitasking possible.

1.6.1 The C prototype

From `defs.h`, the declaration is:

```
void swtch(struct context **old, struct context *new);
```

- `old` – a pointer to a `struct context *` where the current register state will be saved. This is typically `&proc->context` (saving into the process table) or `&cpu->scheduler` (saving the scheduler's state).
- `new` – a pointer to a previously saved `struct context` to restore. This is typically `cpu->scheduler` (switching to the scheduler) or `proc->context` (switching to a process).

1.6.2 What the assembly does

`swtch.S` performs the following steps:

```
# Step 1: Get arguments from the stack
movl 4(%esp), %eax    # eax = old (pointer to where to save current context)
movl 8(%esp), %edx    # edx = new (pointer to context to restore)

# Step 2: Save current (callee-saved) registers by pushing onto current stack
pushl %ebp
pushl %ebx
pushl %esi
pushl %edi

# Step 3: Save the current stack pointer into *old
movl %esp, (%eax)    # *old = esp (i.e., proc->context = esp)

# Step 4: Load the new stack pointer from new
movl %edx, %esp     # esp = new (i.e., esp = cpu->scheduler, or vice versa)

# Step 5: Restore (callee-saved) registers by popping from the new stack
popl %edi
popl %esi
popl %ebx
popl %ebp

# Step 6: Return (pops eip from the new stack)
ret
```

1.6.3 Where the context is saved

Step 3 is the critical save operation: `movl %esp, (%eax)`. Since `old` is `&proc->context`, this writes the current stack pointer into `proc->context` – a field in the process’s entry in the `ptable`. The four registers just pushed (`edi`, `esi`, `ebx`, `ebp`) are now sitting on the kernel stack, and `proc->context` points to them. Together with the return address (`eip`) that was pushed by the `call swtch` instruction, these five values constitute the `struct context`:

```
kernel stack (at the moment of save)
+-----+
|   ...   |
+-----+
|  eip    | <-- pushed by "call swtch" (return address)
+-----+
|  ebp    | <-- pushed by swtch
+-----+
|  ebx    | <-- pushed by swtch
+-----+
|  esi    | <-- pushed by swtch
+-----+
proc->context ---> |  edi    | <-- pushed by swtch (esp saved here)
+-----+
```

This is why `struct context` in `proc.h` has exactly these five fields in this order:

```
struct context {
    uint edi;    // top of stack (lowest address) -- popped last on restore
    uint esi;
    uint ebx;
    uint ebp;
    uint eip;    // bottom of saved block -- the return address from "call swtch"
};
```

The layout of `struct context` matches the push/pop order in `swtch.S` exactly. The `proc->context` pointer doesn’t point to a separately allocated structure – it points directly into the process’s kernel stack, at the location where the registers were pushed.

1.6.4 Why only callee-saved registers?

`swtch()` is called as a normal C function (from `sched()` or `scheduler()`). The x86 C calling convention requires that `swtch` preserve `ebp`, `ebx`, `esi`, and `edi` – these are the callee-saved registers. The caller-saved registers (`eax`, `ecx`, `edx`) are the caller’s responsibility, and the C compiler has already saved them if needed before calling `swtch()`. So `swtch()` only needs to save and restore the callee-saved registers plus the return address.

This means `swtch()` is not saving the *entire* CPU state – just enough to resume a C function that was suspended at a `call swtch` instruction. The full user-space state (all general-purpose registers, segment registers, flags, etc.) is saved separately in the `trapframe` when the process traps into the kernel.

1.6.5 The two levels of saved state

This is worth emphasizing, because it is a common source of confusion:

What is saved	Where it's saved	When it's saved	What saves it
User-space registers (all of them)	<code>struct trapframe</code> on the kernel stack	On trap entry (syscall, interrupt)	<code>trapasm.S</code> (hardware + software)
Kernel-thread registers (callee-saved only)	<code>struct context</code> on the kernel stack, pointed to by <code>proc->context</code>	On context switch	<code>swtch.S</code>

When a process is not running, its full state is captured in these two layers on its kernel stack: the trapframe (holding the user state) and the context (holding the kernel state at the point of the switch). Restoring the process means: `swtch()` restores the kernel context -> the process resumes in `sched()` -> returns through `yield()` or `sleep()` -> returns through `trap()` -> `trapret` restores the trapframe -> `iret` returns to user space.

1.6.6 A concrete example

When the scheduler switches from process P to process Q:

In `sched()` [process P's kernel thread]:

```
swtch(&proc->context, cpu->scheduler)
  -> saves P's {edi,esi,ebx,ebp} on P's kernel stack
  -> stores esp into P's ptable entry: P->context = esp
  -> loads esp from cpu->scheduler
  -> restores scheduler's {edi,esi,ebx,ebp} from scheduler's stack
  -> ret jumps to scheduler's eip (resumes scheduler loop)
```

In `scheduler()` [scheduler thread]:

```
swtch(&cpu->scheduler, Q->context)
  -> saves scheduler's {edi,esi,ebx,ebp} on scheduler's stack
  -> stores esp into cpu->scheduler
  -> loads esp from Q's ptable entry: esp = Q->context
  -> restores Q's {edi,esi,ebx,ebp} from Q's kernel stack
  -> ret jumps to Q's eip (resumes Q in sched(), which returns to yield/sleep)
```

Each `swtch()` call is both an ending and a beginning: it saves one execution context and restores another. The beauty of the design is that `swtch()` itself has no idea what it's switching between – it just saves registers, swaps stack pointers, and restores registers. All the policy (which process to run, when to switch) lives in `scheduler()`, `yield()`, and `sleep()`.

Based on xv6 rev8 (September 1, 2015) – MIT 6.828 teaching operating system.