

CS3210: Processes and switching

Taesoo Kim edited by Kyle Harrigan

Administrivia

- (Nov 1) Team Proposal Day (just slides, target 3-5 min/team)
 - Problem statement
 - Idea
 - Demo plan (aka evaluation)
 - Timeline
- **DUE** : submit slides (as a team) by 10 pm, Oct 31
- **NOTE** : if you submit early, we can give you feedback

Summary of last lectures

- Power-on → BIOS → bootloader → kernel → user programs
- OS: abstraction, multiplexing, isolation, sharing
- Design: monolithic (xv6) vs. micro kernels (jos)
- Abstraction: process, system calls
- Isolation mechanisms: CPL, segmentation, paging

Today's plan

- Aside: dirtyc0w
- A few more notes on locking in xv6
- About process
 - For multiplexing (e.g., more processes than CPUs)
 - In particular, switching and scheduling

dirtycow (CVE-2016-5195)

- What is it?
 - Race condition in kernel memory manager
 - 11 years old! Only recently reported (Oct 19, 2016 by Phil Oester) -- perhaps only recently exploitable?
- Why do we care for purposes of this class?
 - Extremely relevant to in-class topics (memory management, race conditions, paging, copy-on-write, etc.)
 - As always, a simple bug in kernel can have drastic consequences
- [Let us learn some more](#)

Locks

- **Mutual exclusion** : only one core can hold a given lock
 - concurrent access to the same memory location, at least one write
 - example: `acquire(l); x = x + 1; release(l);`

Example: why do we need a lock?

```
00     struct file* filealloc(void) {
01         struct file *f;
02
03         acquire(&ftable.lock);
04         for(f = ftable.file; f < ftable.file + NFILE; f++){
05             if(f->ref == 0){
06                 f->ref = 1;
07                 release(&ftable.lock);
08                 return f;
09             }
10         }
11         release(&ftable.lock);
12         return 0;
13     }
```

Locks

- **Mutual exclusion** : only one core can hold a given lock
 - concurrent access to the same memory location, at least one write
 - example: `acquire(l); x = x + 1; release(l);`
- **Atomic execution** : hide intermediate state
 - another example: transfer money from account A to B
 - `put(a + 100)` and `put(b - 100)` must be both effective, or neither

A different way to think about locks

- Locks help operations maintain **invariants** on a data structure
 - assume the *invariants are true* at start of operation
 - operation uses locks to hide *temporary violation* of invariants
 - operation *restores invariants* before releasing locks
- Q: `put(a + 100)` and `put(b - 100)` ?

Strawman: locking

```
01  struct lock { int locked; };
02
03  void acquire(struct lock *l) {
04      for (;;) {
05          if (l->locked == 0) { // A: test
06              l->locked = 1;    // B: set
07              return;
08          }
09      }
10  }
11
12  void release(struct lock *l) {
13      l->locked = 0;
14  }
```

Problem: concurrent executions on line 05

```
// process A
if (l->locked == 0)
    l->locked = 1;
```

```
// process B
if (l->locked == 0)
    l->locked = 1;
```

- Recall:

```
$ while true; do ./count 2 10 | grep 10 ; done
cpu = 2, count = 10
...
```

Relying on an atomic operation

```
01  struct lock { int locked; };
02
03  void acquire(struct lock *l) {
04      for (;;) {
05          if (xchg(&l->locked, 1) == 0)
06              return;
07      }
08  }
09
10  void release(struct lock *l) {
11      // Q?
12      xchg(&l->locked, 0);
13  }
```


acquire() in xv6

```
01 void acquire(struct spinlock *lk) {
02     // Q1?
03     pushcli();
04     // Q2?
05     if (holding(lk))
06         panic("acquire");
07
08     while (xchg(&lk->locked, 1) != 0)
09         ;
10
11     lk->cpu = cpu;
12     getcallerpcs(&lk, lk->pcs);
13 }
```

release() in xv6

```
01 void release(struct spinlock *lk) {
02     // Q1?
03     if (!holding(lk))
04         panic("release");
05
06     // Q2?
07     lk->pcs[0] = 0;
08     lk->cpu = 0;
09
10     xchg(&lk->locked, 0);
11
12     // Q3?
13     popcli();
14 }
```

Why spinlocks?

- Q: don't they waste CPU while waiting?
- Q: why not give up the CPU and switch to another process, let it run?
- Q: what if holding thread needs to run; shouldn't you yield CPU?

Spinlock guidelines

- hold for very short times
- don't yield CPU while holding lock
- (un)fairness issues: FIFO ordering?
- **NOTE** "blocking" locks for longer critical sections
 - waiting threads yield the CPU
 - but overheads are typically higher (later)

Problem 1: deadlock (e.g., double acquire)

- Q: what happens in xv6?

```
01  struct spinlock lk;  
02  initlock(&lk, "test lock");  
03  acquire(&lk);  
04  acquire(&lk);
```

Problem 2: interrupt (preemption)

- Race in `iderw()` (`ide.c`)
 - `sti()` after `acquire()`
 - `cli()` before `release()`

Q: iderw()

- Q: what goes wrong with adding sti/cli in iderw?
- Q: what ensures atomicity between processors
- Q: what ensures atomicity within a single processor?

What about racing in file.c

- Race in `filealloc()` (`file.c`)
- Q: `ftable.lock` ?
 - `sti()` after `acquire()`
 - `cli()` before `release()`

Q: filealloc()

- Q: could the disk interrupt handler run while interrupts are enabled?
- Q: does any any interrupt handler grab the `ftable.lock`?
- Q: what interrupt could cause trouble?

Scheduling

- Which process to run?
 - Pick one from a set of RUNNABLE processes (or env in jos)
 - Q: what have you seen from lab?
- (next lecture) Switching/scheduling in detail

Scheduling: design space

- Q: Preemptive vs. cooperative?
- Q: Global queue vs. per-CPU queue?

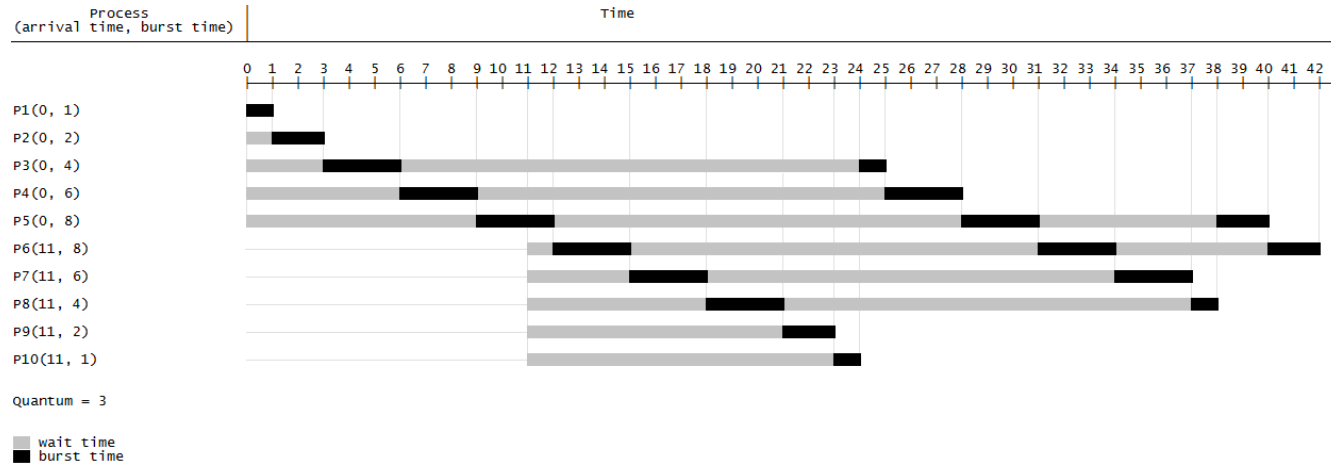
Scheduling: design space

- Scalability: w/ many runnable processes?
- Granularity (timeslice, quantum): 10ms vs 100ms? (dynamic? tickless?)
- Fairness: time quota, epoch (inversion? group?)
- QoS: priority? (e.g., nice)
- Constraints: realtime, deadlines (e.g., airplane)
- etc: resource starvation, performance consolidation (e.g., cloud)

Scheduling: difficult in practice

- No perfect/universal solution/policy
- Contradicting goals:
 - maximizing throughput vs. minimizing latency
 - minimizing response time vs. maximizing scalability
 - maximizing fairness vs. maximizing scalability

Example: round-robin scheduling



- Simple: assign fixed time unit per process
- Starvation-free (no priority)

Complexity in real scheduling algorithms

- Linux?

Complexity in real scheduling algorithms

- Linux
 - `kernel/sched/*.c` : 17k LoC with 7k lines of comments
 - vs. your RR in jos? 10 LoC?

```
01     for (j = 1; j <- NENV; j++) {
02         k = (j + i) % NENV;
03         if (envs[k].env_status == ENV_RUNNABLE)
04             env_run(&envs[k]);
05     }
```

Summary (Wikipedia)

<i>Operating System</i>	<i>Preemption</i>	<i>Algorithm</i>
Amiga OS	Yes	Prioritized round-robin scheduling
FreeBSD	Yes	Multilevel feedback queue
Linux kernel before 2.6.0	Yes	Multilevel feedback queue
Linux kernel 2.6.0–2.6.23	Yes	O(1) scheduler
Linux kernel after 2.6.23	Yes	Completely Fair Scheduler
Mac OS pre-9	None	Cooperative scheduler
Mac OS 9	Some	Preemptive scheduler for MP tasks, and cooperative for processes and threads
Mac OS X	Yes	Multilevel feedback queue
NetBSD	Yes	Multilevel feedback queue
Solaris	Yes	Multilevel feedback queue
Windows 3.1x	None	Cooperative scheduler
Windows 95, 98, Me	Half	Preemptive scheduler for 32-bit processes, and cooperative for 16-bit processes
Windows NT (including 2000, XP, Vista, 7, and Server)	Yes	Multilevel feedback queue

Example: available options in Linux

```
$ sudo sysctl -A | grep "sched" | grep -v "domain"
kernel.sched_child_runs_first = 0
kernel.sched_latency_ns = 18000000
kernel.sched_migration_cost_ns = 500000
kernel.sched_min_granularity_ns = 2250000
kernel.sched_rr_timeslice_ms = 30
kernel.sched_rt_period_us = 1000000
kernel.sched_rt_runtime_us = 950000
kernel.sched_shares_window_ns = 10000000
kernel.sched_time_avg_ms = 1000
kernel.sched_wakeup_granularity_ns = 3000000
...

$ less /proc/sched_debug
$ less /proc/[pid]/sched
```

Characterizing processes

- CPU-bound vs IO-bound
- Interactive processes (e.g., vim, emacs)
- Batch processes (e.g., cronjob)
- Real-time processes (e.g., audio/video players)

Scheduling policies in Linux

- `SCHED_FIFO` : first in, first out, real time processes
- `SCHED_RR` : round robin real time processes
- `SCHED_OTHER` : normal time/schedule sharing (default)
- `SCHED_BATCH` : CPU intensive processes
- `SCHED_IDLE` : Very low prioritized processes

Example

- Q: `count.c` ?

```
$ sudo ./count 3 1000000000  
8522: runs  
8524: runs  
8523: runs  
8523: 2.05 sec  
8522: 2.34 sec  
8524: 2.49 sec
```

Example: available policies

```
$ chrt -m
SCHED_OTHER min/max priority : 0/0
SCHED_FIFO min/max priority : 1/99
SCHED_RR min/max priority   : 1/99
SCHED_BATCH min/max priority : 0/0
SCHED_IDLE min/max priority : 0/0
```

Example: FIFO (real time scheduling)

```
$ sudo ./count 10 1000000000 "chrt -f -p 99"  
...
```

References

- [Intel Manual](#)
- [UW CSE 451](#)
- [OSPP](#)
- [MIT 6.828](#)
- Wikipedia
- The Internet