

# 第二部分

## 一.条件变量

相关文件：proc.c 重点关注问题：

1. 请查阅资料，给出 pthread 提供的条件变量操作，并与 xv6 提供的 sleep & wakeup 操作比较。
2. 为什么 sleep 需要一个 lk 作为参数？（可以参照 xv6 配套讲义中 sleep and wakeup 一节中的 代码改进过程进行说明）
3. 为什么 sleep （以及 wakeup ）要使用 ptable.lock ？ xv6 如何通过锁的使用来解决 lost wake-up 问题？
4. sleep 函数中，2890 行 if 语句的作用是什么？
5. sleep 函数中，2891 行和 2892 行能不能交换顺序？为什么？
6. 阐述 sleep 函数执行时，进程是如何转入睡眠态，又转入就绪态和运行态，并继续执行 sleep 的。
7. xv6 的 wakeup 操作，为什么要拆分成 wakeup 和 wakeup1 两个函数，请举例说明。
8. 假设 wakeup 操作唤醒了多个等待相同 channel 的进程，此时这多个进程会如何执行？ xv6 的 wakeup 是否符合 Mesa semantics？
9. wakeup 时如果没有 sleeping 的进程， wakeup 会阻塞吗？

- 1) 为什么 sleep 要获取 ptable.lock ？

因为在真正把进程设置为 SLEEP 状态之前，子进程可能就已经成为僵死进程并在 exit() 函数中调用了 wakeup1() 来唤醒父进程，但此时父进程还未 SLEEP，这会使得父进程接收不到 wakeup 信号从而进入死锁状态。这种现象的原因是睡眠操作与检测睡眠条件不是一个原子操作，xv6 让进程在 SLEEP 之前获取了 ptable.lock，而 exit() 调用 wakeup 又必须要获取 ptable.lock，所以不会发生 lose wakeup。

- 2) sleep 一直占着 ptable.lock，哪里有释放才能使其他进程获取锁？

在 xv6 中，ptable.lock 总是由旧进程获得，并将锁的控制权转移给切换代码，由新进程释放。锁的这种使用方式很少见，通常来说，持有锁的线程应该负责释放该锁，这样更容易让我们理解其正确性。但对于上下文切换来说，我们必须使用这种方式，因为

ptable.lock 会保证进程的 state 和 context 在运行 swtch 时保持不变。如果在 swtch 中没有持有 ptable.lock，可能引发这样的问题：在 yield 将某个进程状态设置为 RUNNABLE 之后，但又在 swtch 让它停止在其内核栈上运行之前，有另一个 CPU 的 scheduler 检测到该进程位 RUNNABLE，然后调用 swtch，结果将是两个 CPU 都运行在同一个栈上，这显然是不该发生的。这里也可以解释为什么在第一个进程中要把 context 的 eip 设为 forkret，就是为了按照惯例释放 ptable.lock，否则这个新进程是可以直接从 trapret 就开始执行的。所以在进程被置为 SLEEPING 之后要 give up CPU，并不由当前进程释放锁，而是调用 sched()，sched() 会通过下面这句切换到 scheduler：

```
1 | swtch(&p->context, mycpu()->scheduler);
```

而前面说过所有内核线程被切换掉时都会被卡在 `swtch` , `scheduler` 也不例外, 所以 `scheduler` 会从 `swtch` 开始继续往下执行。而下面正有 `release(&ptable.lock);`

3) 为什么有 `lk` 锁, 为什么可以释放?

因为调用 `sleep` 的不一定是 `wait` , 比如发送者和接收者也需要调用 `sleep` 和 `wakeup` (详细的例子参考[xv6调度[睡眠与唤醒]](<http://www.mamicode.com/info-detail-2514856.html>))。 `sleep` 的调用者必须持有锁, 这个锁不是 `ptable.lock`, 是保证二者比如发送者和接收者操作原子性的锁, 作用是防止在调用者在 `test` 和 `sleep` 之间, 对手进程在另一个 `cpu` 上运行, 更新了共享量, 调用了 `wakeup` 却发现没有需要唤醒的进程, 导致对手进程被阻塞在 `test` 上, 而调用者进程也依然 `sleep`, 这就造成死锁。接着 `sleep` 要求持有 `ptable.lock`。于是该进程就会同时持有锁 `ptable.lock` 和 `lk` 了。而如今 `sleep` 已经持有了 `ptable.lock`, 因为 `wakeup` 一定要持有 `ptable.lock`, 所以即便对手进程调用了 `wakeup`, `wakeup` 也不可能在持有 `ptable.lock` 的情况下运行直至 `sleep` 让进程睡眠后, 所以此时 `sleep` 完全可以释放 `lk`, 让对手进程能够修改共享量防止另一种死锁, 同时这样一来, `wakeup` 也不会错过 `sleep`。

在这里是由 `wait` 调用 `sleep`, `lk` 就是 `ptable.lock` 的时候, 这样 `sleep` 就会直接跳过 `lk` 和 `ptable.lock` 这两个步骤。

4) 事实上 `xv6` 实现了 `wakeup` 和 `wakeup1`, 前者获取 `ptable.lock` 之后调用后者。为什么要这样做?

因为有时调度器既可能在持有 `ptable.lock` 的情况下唤醒进程, 也可能在不持有的情况下唤醒。

## • 0. xv6 sleep & wakeup 功能与源码

`Sleep(chan)` sleeps on the arbitrary value `chan`, called the wait channel. `Sleep` puts the calling process to sleep, releasing the CPU for other work.

`Wakeup(chan)` wakes all processes sleeping on `chan` (if any), causing their sleep calls to return. If no processes are waiting on `chan`, `wakeup` does nothing.

```
1 | // Per-process state
2 | struct proc {
3 |     uint sz;                      // Size of process memory (bytes)
4 |     pde_t* pgdir;                 // Page table
5 |     char *kstack;                 // Bottom of kernel stack for this
    process
6 |     enum procstate state;         // Process state
7 |     int pid;                      // Process ID
8 |     struct proc *parent;          // Parent process
9 |     struct trapframe *tf;         // Trap frame for current syscall
10 |    struct context *context;       // swtch() here to run process
11 |    void *chan;                   // If non-zero, sleeping on chan
12 |    int killed;                   // If non-zero, have been killed
13 |    struct file *ofile[NOFILE];   // Open files
```

```

14     struct inode *cwd;           // Current directory
15     char name[16];              // Process name (debugging)
16 };
17

```

## - 0.1 sleep源码

```

1 // Atomically release lock and sleep on chan.
2 // Reacquires lock when awakened.
3 void
4 sleep(void *chan, struct spinlock *lk)
5 {
6     struct proc *p = myproc();
7
8     if(p == 0)
9         panic("sleep");
10
11     if(lk == 0)
12         panic("sleep without lk");
13
14     // Must acquire ptable.lock in order to
15     // change p->state and then call sched.
16     // Once we hold ptable.lock, we can be
17     // guaranteed that we won't miss any wakeup
18     // (wakeup runs with ptable.lock locked),
19     // so it's okay to release lk.
20     if(lk != &ptable.lock){ //DOC: sleeplock0
21         acquire(&ptable.lock); //DOC: sleeplock1
22         release(lk);
23     }
24     // Go to sleep.
25     p->chan = chan;
26     p->state = SLEEPING;
27
28     sched();
29
30     // Tidy up.
31     p->chan = 0;
32
33     // Reacquire original lock.
34     if(lk != &ptable.lock){ //DOC: sleeplock2
35         release(&ptable.lock);
36         acquire(lk);
37     }
38 }

```

## - 0.2 wakeup源码

```
1 //PAGEBREAK!
2 // Wake up all processes sleeping on chan.
3 // The ptable lock must be held.
4 static void
5 wakeup1(void *chan)
6 {
7     struct proc *p;
8
9     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++)
10         if(p->state == SLEEPING && p->chan == chan)
11             p->state = RUNNABLE;
12 }
13
14 // Wake up all processes sleeping on chan.
15 void
16 wakeup(void *chan)
17 {
18     acquire(&ptable.lock);
19     wakeup1(chan);
20     release(&ptable.lock);
21 }
```

## • 1. pthread vs xv6

### - 1.1 pthread中条件阻塞&唤醒函数

#### - 1.1.1 pthread\_cond\_wait

对应于sleep，pthread\_cond\_wait()函数等待条件变量变为真的。它需要两个参数，第一个参数就是条件变量，而第二个参数mutex是保护条件变量的互斥量。也就是说这个函数在使用的时候需要配合pthread\_mutex\_lock()一起使用。

```
1 pthread_mutex_lock(&mutex);
2 pthread_cond_wait(&cond, &mutex);
```

该函数功能：

1. 等待条件变量满足；
2. 把获得的锁释放掉；（注意：1，2两步是一个**原子操作**）

⋮ 释放锁这一步和等待条件满足一定是一起执行

pthread\_cond\_timedwait()函数和 pthread\_cond\_wait()函数比起来多一个时间参数, 这个参数可以指定等待这个条件多长时间, 通过timespec结构指定。

• [线程同步之条件变量 \(pthread cond wait\) - 腾讯云开发者社区-腾讯云 \(tencent.com\)](#)

对比于xv6 的sleep, 其返回前对锁的操作如下:

```
1  if(lk != &ptable.lock){ //DOC: sleeplock2
2      release(&ptable.lock);
3      acquire(lk);
4  }
```

倘若当前locked的锁就是lk, 啥也不干; 倘若不是, 释放当前锁, 获得lk, 返回。

也即是说, wait 返回后会放弃当前线程对 lock的持有, 而 sleep 还保持着锁的lock 状态。

## - 1.1.2 pthread的条件“wakeup”

pthread\_cond\_wait只针对于当前调用该函数的进程, 而xv6的wakeup (参数chan) 唤醒了堵在chan (某个特定的channel) 的所有进程, 其更类似 pthread\_cond\_broadcast:

```
1  int pthread_cond_signal(pthread_cond_t *cond);
2  int pthread_cond_broadcast(pthread_cond_t *cond);
```

pthread\_cond\_signal()激活一个等待该条件的线程, 存在多个等待线程时按入队顺序激活其中一个; 而 pthread\_cond\_broadcast 释放了所有堵在cond上的进程。

补充资料:

在 wakeup 中遍历整个进程表来寻找对应 chan 的进程是非常低效的。更好的办法是用另一个结构体代替 sleep 和 wakeup 中的 chan, 该结构体中维护了一个睡眠进程的链表。Plan 9 的 sleep 和 wakeup 把该结构体称为集合点 (rendezvous point) 或者 Rendez。许多线程库都把相同的一个结构体作为一个状态变量; 如果是这样的话, sleep 和 wakeup 操作则被称为 wait 和 signal。所有此类机制都有同一个思想: 使得睡眠状态可以被某种执行原子操作的锁保护。

在 wakeup 的实现中, 它唤醒了特定队列中的所有进程, 而有可能很多进程都在同一个队列中等待被唤醒。操作系统会调度这里的所有进程, 它们会互相竞争以检查其睡眠状态。这样的一堆进程被称作惊群 (thundering herd), 而我们最好是避免这种情况的发生。大多数的状态变量都有两种不同的 wakeup, 一种唤醒一个进程, 即 signal; 另一种唤醒所有进程, 即 broadcast。

## • 2. sleep中为啥要一个lock参数?

首要目标：使得判断和阻塞进程成为一个原子操作。

### - 2.0 为何要sleep?

看一看讲义给出的例子：

```
1 struct q
2 {
3     void * ptr;
4 };
5 //-----ver 1-----
6 void*send(struct q *q, void *p)
7 {
8     while(q->ptr != 0); // receiver在接受数据
9     q->ptr = p;         // sender存入新数据
10 }
11
12 void*recv(struct q *q)
13 {
14     void *p;
15     while((p = q->ptr) == 0)    // 等待接收数据
16         ;
17     q->ptr = 0;
18     return p;
19 }
```

如果不使用sleep和wakeup，且sender和reciever运行在不同核，功能上没毛病，但在第15行rec不断check是否有发来的新信息，浪费资源。故我们希望引入sleep和wakeup。

### - 2.1 如果sleep不带锁——lost wake up

```
1 // p 是sender要送给receiver的信息
2 void*send(struct q *q, void *p)
3 {
4     while(q->ptr != 0); // q里还有数据
5     q->ptr = p;         // sender存入新数据
6     wakeup(q); /* wake recv */ // 让receiver接受数据
7 }
8
9 void*recv(struct q *q)
10 {
11     void *p;
```

```

12         while((p = q->ptr) == 0)
13             sleep(q);    // recv让出cpu, 在channel q上睡下了
14         q->ptr = 0;      // q的数据槽清空
15         return p;       // 返回接收到的数据p
16     }

```

若recv在11判断  $(p = q->ptr) == 0$  为真，正准备执行sleep(q)时切换到send，send执行了  $q->ptr = p$  后并调用了wakeup（此时会发现），此时recv会无脑进入sleep状态，然后一睡不醒 :)（sender执行完wakeup后已经拍拍屁股走人了）。

## - 2.2 如果只是把sleep放进临界区——死锁

```

1  //-----ver 3-----
2  // p 是sender要送给receiver的信息
3  void*send(struct q *q, void *p)
4  {
5      acquire(&q->lock);
6      while(q->ptr != 0); // q里还有数据
7      q->ptr = p;        // sender存入新数据
8      wakeup(q); /* wake recv */ // 让receiver接受数据
9      release(&q->lock)
10 }
11
12 void*recv(struct q *q)
13 {
14     void *p;
15     acquire(&q->lock);
16     while((p = q->ptr) == 0)
17         sleep(q);    // recv让出cpu, 在channel q上睡下了
18     q->ptr = 0;      // q的数据槽清空
19     release(&q->lock)
20     return p;       // 返回接收到的数据p
21 }

```

假设recv在16处获得了锁，并判断出  $(p = q->ptr) == 0$  后再次进入sleep状态，sender将永远无法获得锁，更无法执行wakeup，recv还是一睡不醒.....

## - 2.3 需求：sleep时不要占着茅坑不拉shi🤔

我们希望recv进入睡眠状态时自觉把锁释放，返回后再自动重新拿回锁。故而引入参数lock：

```

1  //-----ver 4-----
2  // p 是sender要送给receiver的信息

```

```

3     void*send(struct q *q, void *p)
4     {
5         acquire(&q->lock);
6         while(q->ptr != 0); // q里还有数据
7         q->ptr = p;         // sender存入新数据
8         wakeup(q); /* wake recv */ // 让receiver接受数据
9         release(&q->lock);
10    }
11
12    void*recv(struct q *q)
13    {
14        void *p;
15        acquire(&q->lock);
16        while((p = q->ptr) == 0)
17            sleep(q, &q->lock); // recv让出cpu, 在channel q上睡下后
18        // 同时释放锁, 返回后再拿回锁
19        q->ptr = 0; // q的数据槽清空
20        release(&q->lock);
21        return p; // 返回接收到的数据p
22    }

```

## • 3. sleep中ptable.lock的作用

### - 3.1 三千弱水，只取ptable的锁

在操作系统中，所有的进程信息 `struct proc` 都存储在 `ptable` 中，`ptable` 的定义如下

```

1 struct {
2     struct spinlock lock;
3     struct proc proc[NPROC];
4 } ptable;

```

但为什么明明进程手里已经有一把锁，却一定要ptable的锁？要回答这个问题，要先看看sleep后续要做什么：

```

1 if(lk != &ptable.lock){ //DOC: sleeplock0
2     acquire(&ptable.lock); //DOC: sleeplock1
3     release(lk);
4 }
5 ....
6     sched();
7 ....

```



即其要调用调度器 `sched` 做上下文切换：

```
1 // Enter scheduler. Must hold only ptable.lock
2 // and have changed proc->state. Saves and restores
3 // intena because intena is a property of this
4 // kernel thread, not this CPU. It should
5 // be proc->intena and proc->ncli, but that would
6 // break in the few places where a lock is held but
7 // there's no process.
8 void
9 sched(void)
10 {
11     ...
12     if(!holding(&ptable.lock))
13         panic("sched ptable.lock");
14     ...
15     swtch(&p->context, mycpu()->scheduler);
16     ...
17 }
18
```

如果当前没有获得ptable的锁，`sched`会抛出一个panic信息，然后异常终止程序。其根本原因在于后续要做上下文的切换 `swtch`，在这个过程中要进行当前context的保存以及下一个进程的context恢复，`ptable.lock` 会保证进程的 `state` 和 `context` 在运行 `swtch` 时保持不变。

如果在 `swtch` 中无需持有 `ptable.lock`，可能在多核的情况下引发这样的问题：在另一个核上运行的程序对当前进程调用了一下 `wakeup`，将当前进程的状态又置为了 `RUNNABLE`——这显然会导致混乱。而如果我们确保 `wakeup` 和 `sleep` 等修改进程 `state` 等信息的操作都在临界区进行，就可以避免上述问题。

• e.g, 其他 `swtch` 需持有锁的例证：yield 将某个进程状态设置为 `RUNNABLE` 之后，但又在 `swtch` 让它停止在其内核栈上运行之前，有另一个 CPU 的scheduler检测到该进程位 `RUNNABLE`，然后调用 `swtch`，结果将是**两个 CPU 都运行在同一个栈上**，这显然是不该发生的。

## - 3.2 两把锁通力协作：避免lost wake up

- lk保证要么recv不sleep，要么必须在wakeup前sleep：
  - 假设sender先抢到锁：

```
1         acquire(&q->lock);
2         while(q->ptr != 0); // q里还有数据
3         q->ptr = p;         // sender存入新数据
4         wakeup(q); /* wake recv */ // 让receiver接受数据
5         release(&q->lock);
```

- 若q里有数据，死等；
- 若q里无数据，填充数据p，唤醒recv：

```

1      acquire(&q->lock);
2      while((p = q->ptr) == 0)
3          sleep(q, &q->lock);    // recv让出cpu，在
channel q上睡下后同时释放锁，返回后再拿回锁
4      q->ptr = 0;    // q的数据槽清空
5      release(&q->lock);

```

此时q里有数据，故recv不会执行sleep。

- 假设recv先抢到锁：
  - 若q里没有数据，睡！→切换到sender→wakeup
  - 若q里有数据，干活
- 中途获取了ptable.lock后才释放了lk：没人会修改进程的状态（详见4.2）

## • 4.sleep函数抢锁时的细节

### - 4.1 if的作用？

```

1      if(lk != &ptable.lock){    //DOC: sleeplock0
2          acquire(&ptable.lock);    //DOC: sleeplock1
3          release(lk);
4      }

```

- lk==&ptable.lock：啥也不干
- lk!=&ptable.lock：获取ptable.lock，释放lk；

最终达到的效果都是：**当前进程只抢了一把ptable的锁。**

### - 4.2 acquire和release的顺序？

必须要先获得ptable的锁再释放lk。

以下直接照搬去年ppt，待修改：

- acquire ptable.lock可以保护临界区资源p->state，release lk是为了防止receiver在睡眠时持有锁，从而造成死锁。拿到ptable.lock锁之后，即使有另外的进程调用wakeup函数，wakeup函数也会等待直到它拿到了ptable.lock，wakeup函数一定是在sleep执行完毕之后。所以acquire ptable.lock之后release lk是安全的。
- 如果交换顺序，在两个指令之间存在未被锁保护的区域。如果在此处有新的进程进行wakeup函数，将会漏掉正在进行sleep操作的进程，产生错误

## • 5. sleep的状态转换

sleep中将进程挂起，把进程p放进channel chan内，状态置为sleep

```
1 // Go to sleep.
2 p->chan = chan;
3 p->state = SLEEPING;
```

下一步调用sched：

```
1 sched();
```

```
1 void
2 sched(void)
3 {
4     int intena;
5     struct proc *p = myproc();    // 获取当前进程
6     ...
7     swtch(&p->context, mycpu()->scheduler);    // 调用调度器scheduler
8     mycpu()->intena = intena;
9 }
10
```

调度器检查当前进程表，挑选出一个 **RUNNABLE** 的并把之状态置为 **RUNNING**：

```
1 void
2 scheduler(void)
3 {
4     struct proc *p;
5     struct cpu *c = mycpu();
6     c->proc = 0;
7
8     for(;;){
9         // Enable interrupts on this processor.
10        sti();
11
12        // Loop over process table looking for process to run.
13        acquire(&ptable.lock);
14        for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
15            if(p->state != RUNNABLE)
16                continue;
17
18            // Switch to chosen process. It is the process's job
19            // to release ptable.lock and then reacquire it
20            // before jumping back to us.
21            c->proc = p;
```

```

22     switchvm(p);
23     p->state = RUNNING;
24
25     switch(&(c->scheduler), p->context);
26     switchkvm();
27
28     // Process is done running for now.
29     // It should have changed its p->state before coming back.
30     c->proc = 0;
31 }
32 release(&ptable.lock);
33
34 }
35 }

```

值得一提的是其调用了 `switch(&(c->scheduler), p->context)`，其作用是保存上下文（这将确保后续能重新回到sleep处继续执行）：

```

1  # Context switch
2  #
3  # void swtch(struct context **old, struct context *new);
4  #
5  # Save the current registers on the stack, creating
6  # a struct context, and save its address in *old.
7  # Switch stacks to new and pop previously-saved registers.
8
9  .globl swtch
10 swtch:
11     movl 4(%esp), %eax
12     movl 8(%esp), %edx
13
14     # Save old callee-saved registers
15     pushl %ebp
16     pushl %ebx
17     pushl %esi
18     pushl %edi
19
20     # Switch stacks
21     movl %esp, (%eax)
22     movl %edx, %esp
23
24     # Load new callee-saved registers
25     popl %edi
26     popl %esi
27     popl %ebx
28     popl %ebp
29     ret

```

后续何时醒？取决于何时调用wakeup，其重新将chan的进程的状态置为

RUNNABLE：

```
1 for(p = ptable.proc; p < &ptable.proc[NPROC]; p++)
2     if(p->state == SLEEPING && p->chan == chan)
3         p->state = RUNNABLE;
```

## • 6. 分成wakeup 和 wakeup1?

可能有时需要连续执行若干次wakeup1，此时若调用 wakeup，将反复获取释放锁，开销👉。如proc.c中的exit：

```
1 void
2 exit(void)
3 {
4     ...
5     acquire(&ptable.lock);
6
7     // Parent might be sleeping in wait().
8     wakeup1(curproc->parent);
9
10    // Pass abandoned children to init.
11    for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
12        if(p->parent == curproc){
13            p->parent = initproc;
14            if(p->state == ZOMBIE)
15                wakeup1(initproc);
16        }
17    }
18
19    // Jump into the scheduler, never to return.
20    curproc->state = ZOMBIE;
21    sched();
22    panic("zombie exit");
23 }
```

第8行进行了wakeup1后，还需进入循环再进行若干次wakeup1，同时最后还需调用sche，故无需release ptable.lock。

## • 7. channel里多个/无进程的情况？

wakeup1逐个将channel中的进程置为 `RUNNABLE`：

```
1   for(p = ptable.proc; p < &ptable.proc[NPROC]; p++)
2       if(p->state == SLEEPING && p->chan == chan)
3           p->state = RUNNABLE;
```

### - 7.1 多个被唤醒的进程的执行顺序？

取决于调度器算法，xv6使用了最朴素的Round Robin，故其先后顺序只取决于其在进程表中的顺序。

疑问：

为什么在wakeup1中只修改了其 `state` 而没有把chan域重新置零？即使其对判断进程状态没有影响，但基于其官方给出的注释，是不是置零比较好？

```
1   struct proc {
2       ...
3       void *chan;           // If non-zero, sleeping on chan
4       ...
5   };
```

### - 7.2 是否符合Mesa semantics

何为Mesa Semantics：

```
1   Mesa proposes that when a signal call is made, a thread will be
    taken out from the waiting queue, and will be put on the ready
    queue. However, this doesn't mean that the thread will run
    immediately. The scheduler can choose when to run this thread, which
    means that another thread could have run and changed the state of
    the object!
```

符合，wakeup1逐个将chan中进程置为 `RUNNABLE` 后就调用了 `sched` 去做调度。

拓展（只是为了更好地帮助理解这题在为什么）：

- Brinch Hansen：要求signal发生在线程A离开管程的时候，也就是说，signal之后，线程B就离开管程了，线程A就自然进入管程。（无需经过调度器）
- Hoare：除却上述两种queue，还有signal queue。等的线程A在wait queue，signal发生时，线程B被从管程中移到signal queue中，而线程A则从wait queue移到管程中，等线程A离开管程后线程A再回来。

## - 7.3 无进程是否会阻塞wakeup?

唤醒时的循环如下：

```
1   for(p = ptable.proc; p < &ptable.proc[NPROC]; p++)
2       if(p->state == SLEEPING && p->chan == chan)
3           p->state = RUNNABLE;
```

其循环初始条件为 `p = ptable.proc`，退出循环的条件是 `p < &ptable.proc[NPROC]`，循环的每轮更新操作是 `p++`（等价于 `for(int i=0; i<NPROC; i++)`，其显然是不会卡死的）。

## 二.Linux 信号量

请阅读Linux kernel中关于信号量的代码 [include\linux\semaphore.h, kernel\locking\semaphore.c] 以及查找相关资料，回答如下问题：

1. Linux中信号量的数据结构及对应的PV操作；
2. 说明Linux扩展的各类 down 操作的用途；
3. 简要分析 down 和 up 操作的实现。

### • 1.Linux中信号量的数据结构&PV操作

#### - 1.1 Linux中信号量的数据结构

Linux Kernel 除了提供了自旋锁，还提供了睡眠锁，信号量就是一种睡眠锁。信号量的特点是，如果一个任务试图获取一个已经被占用的信号量，他会被推入等待队列，让其进入睡眠。

```
1  /* Please don't access any members of this structure directly */
2  struct semaphore {
3      raw_spinlock_t    lock;
4      unsigned int      count;
5      struct list_head  wait_list;
6  };
7
```

补充资料——信号量和互斥锁（另一种睡眠锁）的区别：

	semaphore	mutex
几把锁	任意，可设置	1
谁能解锁	别的程序、中断等都可以	谁加锁，就得由谁解锁
多次解锁	可以	不可以，因为只有 1 把锁
循环加锁	可以	不可以，因为只有 1 把锁
任务在持有锁的期间可否退出	可以	不建议，容易导致死锁
硬件中断、软件中断上下文中使用	可以	不可以

## 1.2 信号量的PV操作

- p操作：对信号量减1，若结果大于等于0，则进程继续，否则执行p操作的进程被阻塞等待释放。
- v操作：对信号量加1，若结果小于等于0，则唤醒队列中一个因为p操作而阻塞的进程。

函数名	作用
void down(struct semaphore *sem)	获得信号量，如果暂时无法获得就会休眠返回之后就表示肯定获得了信号量在休眠过程中无法被唤醒，即使有信号发给这个进程也不处理
int down_trylock(struct semaphore *sem)	尝试获得信号量，不会休眠，返回值：0：获得了信号量 1：没能获得信号量
void up(struct semaphore *sem)	释放信号量，唤醒其他等待信号量的进程

## 2. 扩展的各类 down 操作的用途

## 3. down 和 up 操作的实现

一篇很棒的blog（注意版本不一致）：[\(3条消息\) Linux Kernel Semaphore代码解析Murick的博客-CSDN博客kernel semaphore](#)

down操作统一调用了\_\_down\_common：

```

1 static ninline void __sched __down(struct semaphore *sem)
2 {
3     __down_common(sem, TASK_UNINTERRUPTIBLE, MAX_SCHEDULE_TIMEOUT);
4 }
5
6 static ninline int __sched __down_interruptible(struct semaphore
7 *sem)
8 {
9     __down_common(sem, TASK_INTERRUPTIBLE, MAX_SCHEDULE_TIMEOUT);
10 }

```



```

8     return __down_common(sem, TASK_INTERRUPTIBLE,
MAX_SCHEDULE_TIMEOUT);
9 }
10
11 static ninline int __sched __down_killable(struct semaphore *sem)
12 {
13     return __down_common(sem, TASK_KILLABLE, MAX_SCHEDULE_TIMEOUT);
14 }
15
16 static ninline int __sched __down_timeout(struct semaphore *sem,
long timeout)
17 {
18     return __down_common(sem, TASK_UNINTERRUPTIBLE, timeout);
19 }

```

查看\_\_down\_common:

```

1 static inline int __sched __down_common(struct semaphore *sem, long
state,
2                                     long timeout)
3 {
4     struct semaphore_waiter waiter;
5
6     list_add_tail(&waiter.list, &sem->wait_list);
7     waiter.task = current;
8     waiter.up = false;
9
10    for (;;) {
11        if (signal_pending_state(state, current))
12            goto interrupted;
13        if (unlikely(timeout <= 0))
14            goto timed_out;
15        __set_current_state(state);
16        raw_spin_unlock_irq(&sem->lock);
17        timeout = schedule_timeout(timeout);
18        raw_spin_lock_irq(&sem->lock);
19        if (waiter.up)
20            return 0;
21    }
22
23    timed_out:
24        list_del(&waiter.list);
25        return -ETIME;
26
27    interrupted:
28        list_del(&waiter.list);
29        return -EINTR;
30 }
31

```

\_\_down\_common使用的API:

list\_add\_tail: 见<include/linux/list.h>, 其作用是将entry加入head指示的循环列表的表尾:

```
1 static inline void
2 list_add_tail(struct list_head *entry, struct list_head *head)
3 {
4     __list_add(entry, head->prev, head);
5 }
6 static inline void __list_add(struct list_head *new,
7                               struct list_head *prev,
8                               struct list_head *next)
9 {
10     if (!__list_add_valid(new, prev, next))
11         return;
12
13     next->prev = new;
14     new->next = next;
15     new->prev = prev;
16     WRITE_ONCE(prev->next, new);
17 }
```

signal\_pending\_state: 判断进程是否需要立即回到RUNNING状态。

如果进程是INTERRUPTIBLE状态 或 进程是KILLABLE且收到kill信号则进程需要立即返回RUNNING 状态。其他情况则不需要。源码位于<include/linux/sched/signal.h>:

```
1 static inline int signal_pending_state(long state, struct task_struct
2 *p)
3 {
4     if (!(state & (TASK_INTERRUPTIBLE | TASK_WAKEKILL)))
5         return 0;
6     // signal_pending(p)检查当前进程是否有信号处理, 返回不为0表示有信号需要
7     处理
8     if (!signal_pending(p))
9         return 0;
10     return (state & TASK_INTERRUPTIBLE) || __fatal_signal_pending(p);
11 }
```

unlikely: 位于<include/linux/compiler.h>, 理解可见[linux中likely\(\)和unlikely\(\) - yuxi o - 博客园 \(cnblogs.com\)](linux中likely()和unlikely() - yuxi o - 博客园 (cnblogs.com))

```
1 # define unlikely(x)  (__branch_check__(x, 0,
2 __builtin_constant_p(x)))
```

```
1 #define __branch_check__(x, expect, is_constant) ({          \
2     int ____r;                                                \
3     static struct ftrace_likely_data                          \
4     __attribute__((__aligned__(4)))                          \
```

```

5         __attribute__((section("_ftrace_annotated_branch"))) \
6         _____f = { \
7             .data.func = __func__, \
8             .data.file = __FILE__, \
9             .data.line = __LINE__, \
10        }; \
11        _____r = __builtin_expect(!(x), expect); \
12        ftrace_likely_update(&_____f, _____r, \
13                             expect, is_constant); \
14        _____r; \
15    })

```

raw\_spin\_unlock\_irq:

```

1  #define raw_spin_unlock_irq(lock)    _raw_spin_unlock_irq(lock)
2  ...
3  #define _raw_spin_unlock_irq(lock) __raw_spin_unlock_irq(lock)
4  static inline void __raw_spin_unlock_irq(raw_spinlock_t *lock)
5  {
6      spin_release(&lock->dep_map, 1, _RET_IP_);
7      do_raw_spin_unlock(lock);
8      local_irq_enable(); // 允许本地中断
9      preempt_enable();   // 允许内核抢占
10 }

```

### 三. QEMU跟踪sleep和wakeup的过程

参考资料:

[Ubuntu+QEMU+Xv6环境搭建 - AlexAlex - 博客园 \(cnblogs.com\)](#)

[MIT 6.S081 xv6 调试指北 - 操作系统 \(2021秋季\) | 哈工大 \(深圳\) \(gitee.io\)](#)

查看Makefile的结构:

```

1  UPROGS=\
2      _cat\
3      _echo\
4      _forktest\
5      _grep\
6      _init\
7      _kill\
8      _ln\
9      _ls\
10     _mkdir\

```

```

11     _rm\
12     _sh\
13     _stressfs\
14     _usertests\
15     _wc\
16     _zombie\

```

在此处指定了给用户的API，添加上我们写的 `_mytest`。

由于xv6似乎没提供创建线程的API，改为使用fork创建进程：

```

1

```

但是头铁直接调用内核API是不行的：

```

wc          2 16 14164
zombie      2 17 12408
mytest      2 18 14896
console     3 19 0
$ mytest
pid 4 mytest: trap 6 err 0 on cpu 1 eip 0x20 addr 0x0--kill proc
pid 5 mytest: trap 6 err 0 on cpu 0 eip 0x20 addr 0x0--kill proc
$ zombie!
QEMU: Terminated
ubuntu@ubuntu:~/Documents/xv6-public$ S

```

起初尝试内联汇编手动进入内核态：

```

objdump -t _zombie | sed -i,1,7s/0x/0/ > zombie.sym
gcc -fno-pic -static -fno-builtin -fno-strict-aliasing -O2 -Wall -MD -ggdb -m32
-Werror -fno-omit-frame-pointer -fno-stack-protector -fno-pie -no-pie -c -o my
test.o mytest.c
mytest.c: Assembler messages:
mytest.c:21: 错误: no such instruction: `ecall'
<内置>: recipe for target 'mytest.o' failed
make: *** [mytest.o] Error 1

```

但不知道为啥不能识别 `ecall` 指令。于是尝试在内核初始化时尚未进入用户态时进行测试，但是初始化卡死：

```

xv6...
cpu1: starting 1
xv6...
cpu1: starting 1
xv6...
cpu1: starting 1
xv6...
cpu1: starting 1
xv6...
cpu1: starting 1
xv6...
cpu1: starting 1
xv6...

```

gdb跟踪发现变量被optimised out了，查找资料：

- [\(3条消息\) C++调试时出现“optimized out”的原因、解决办法y\\_m\\_h的博客-CSDN博客qt](#)
- [optimized out](#)

在Makefile中注明使用-O0