# Linux Character Device Driver

*Sunbeam Infotech*

# Register character device

- Each char device is represented by struct cdev.

- *cdev_map* is global array (hash table) to keep track of all char devices.

- cdev_map is object of struct kobj_map.
  - key = device major number
  - hash function = major % 255
  - value = struct probe
    - void *data = struct cdev

- (2) Add device into char device database *cdev_map*.
  - cdev_init(&cdev, &fops);
  - cdev_add(&cdev, devno, dev_count);

- Added device can be removed.
  - cdev_del(&cdev);

*file_operations*

| |
|---|
| Owner |
| open * |
| release * |
| read * |
| write * |
| ... |

→ Pchar_open()
    {...}

→ Pchar_close()
    {...}

→ Pchar_read()
    {...}

→ Pchar_write()
    {...}

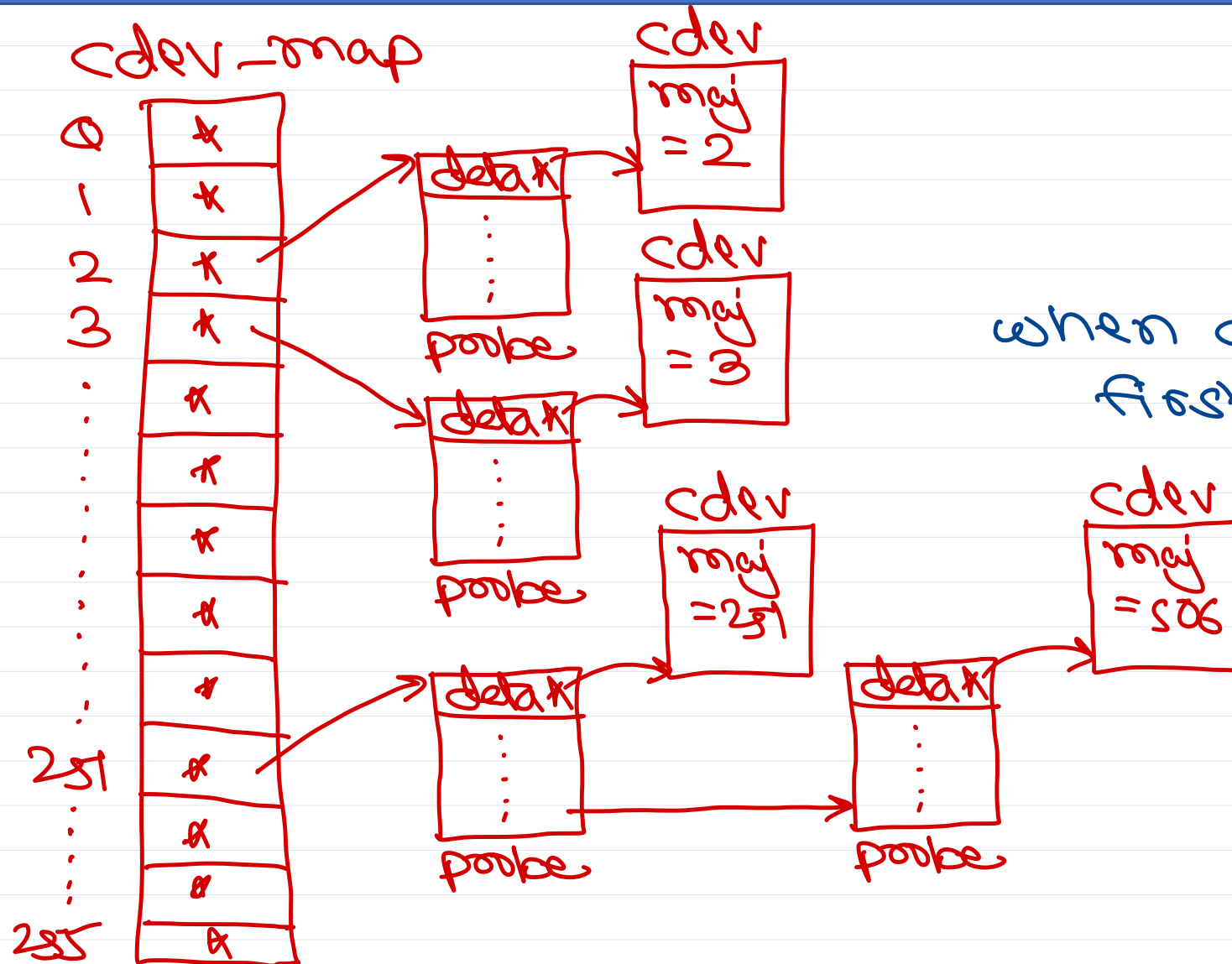dev no

1

```
struct kobj_map {
    struct probe {
        dev_t dev;
        unsigned long range;
        ...
        void *data;
    } *probes[255];
    ...
};
```

*struct kobj_map cdev_map;*

```
struct cdev {
    struct kobject kobj;
    struct module *owner;
    struct file_operations *ops;
    struct list_head list;
    dev_t dev;
    unsigned count;
};
```
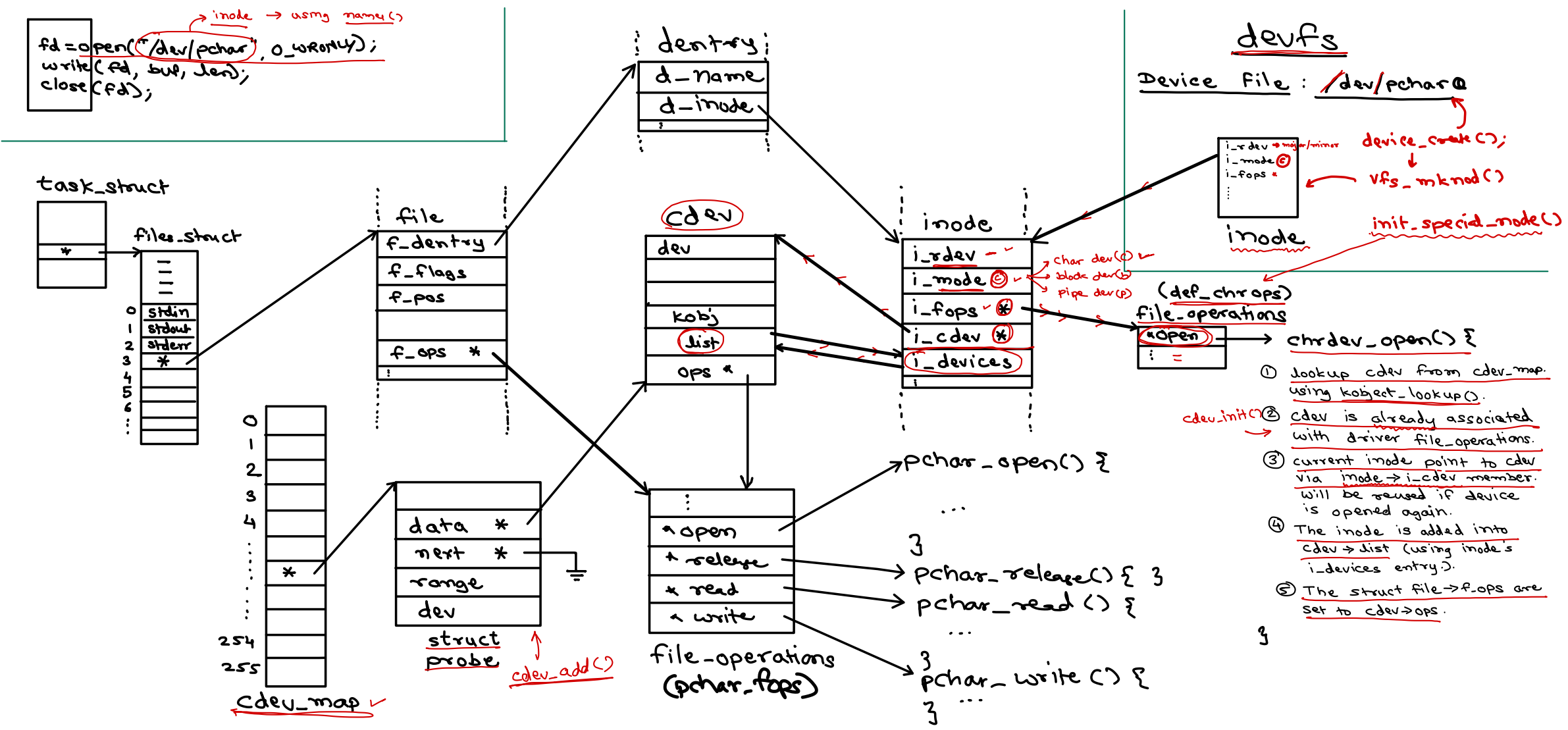
# cdev_add()



cdev_map

0
1
2
3
...
251
...
255

cdev
maj = 2

cdev
maj = 3

cdev
maj = 251

cdev
maj = 506

poobe
poobe
poobe
poobe

key = major e.g. 506
hash fn = key % 255
e.g. 506 % 255 = 251

when device file is opened for first time, its cdev will be searched from this cdev_map using dev no (i_rdev).

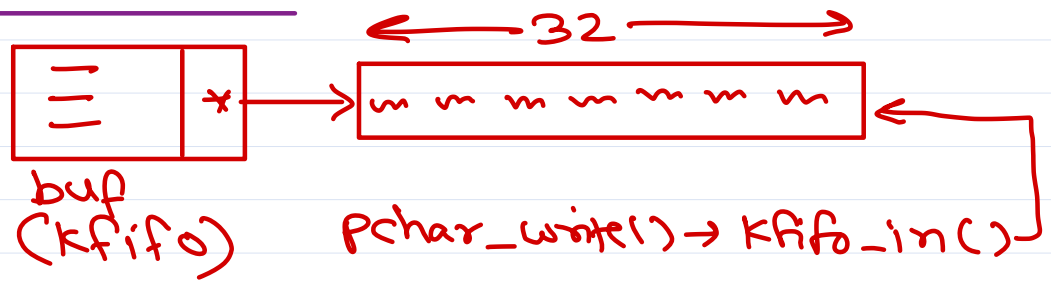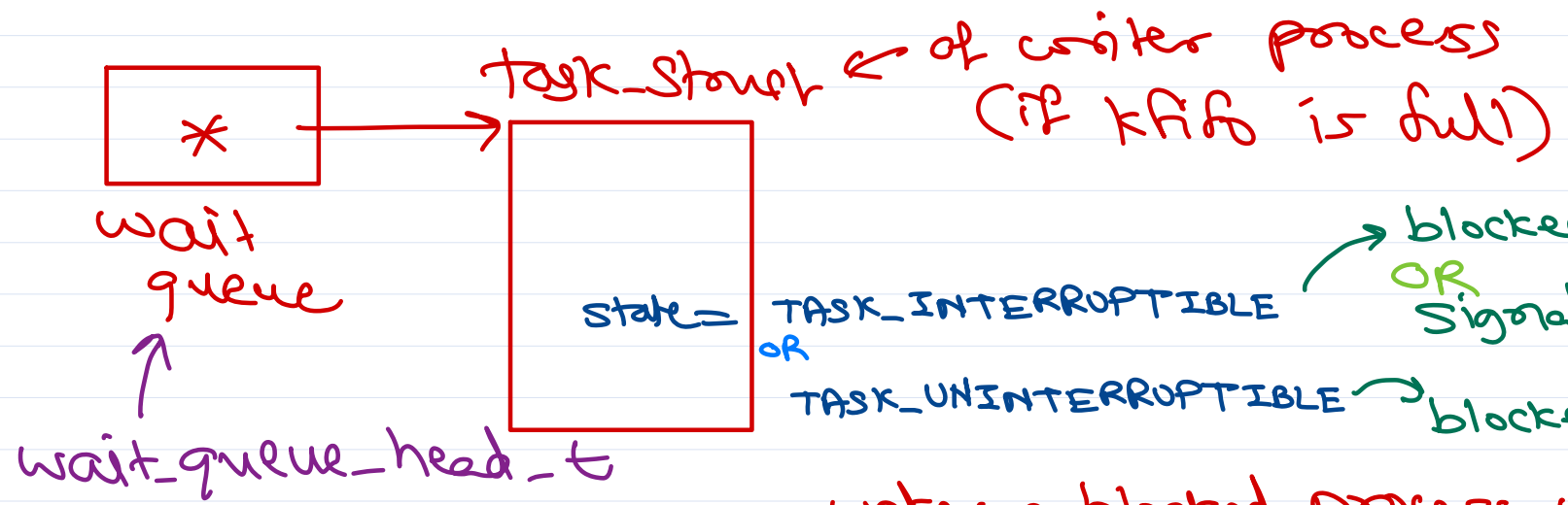# Execution Flow of Pseudo Char Device Driver



Sunbeam Infotech
www.sunbeaminfo.com

# Waiting Queue

pchar device



buf
(kfifo)

pchar_write() → kfifo_in()

←—— 32 ——→

if kfifo is full, then block the writing process.
↳ wait_event_interruptible()
OR
wait_event()

→ kfifo_is_full()

task_struct ← of writer process (if kfifo is full)

wait queue

↑
wait_queue_head_t

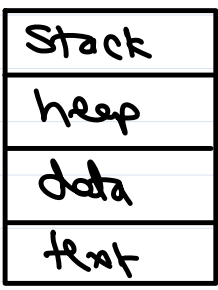state = TASK_INTERRUPTIBLE
OR
TASK_UNINTERRUPTIBLE

→ blocked until data is popped from fifo.
OR
Signal is received.

→ blocked until data is popped from fifo.

wake up blocked process, when some data is removed from fifo → kfifo_out() – read op.
OR
wake_up()
wake_up_interruptible()

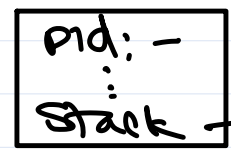# Kernel stack - to create FARs of kernel functions
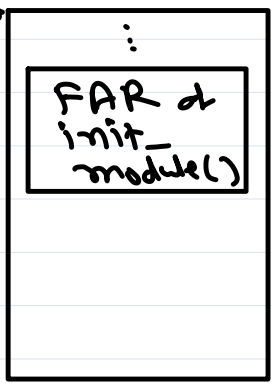
cmd> sudo insmod pchar.ko

insmod 2
Process

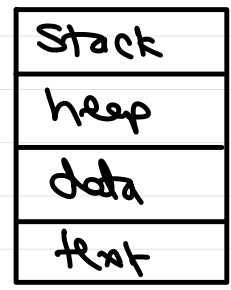| Stack |
|-------|
| heep |
| data |
| text |

→ load_module ()

Sys_load_module ()

① check pinv
② create struct module
   state = COMING.
③ load module in
   kernel space.
④ call its init_module()
⑤ module state = LIVE;

}

task_struct
(insmod)

| pid: – |
| ⋮ |
| Stack |

→ kernel Stack

⋮

| FAR of init_module() |

↓

---

cmd> sudo cat /dev/pchar

cat 2
Process

| Stack |
|-------|
| heep |
| data |
| text |

→ read ()

Sys_read () ?

① get OFT entry
   ie. struct file
② call
   file→f_ops→read()

}

pchar_read() ?
   ≡

}

task_struct
(Cat)

| pid: – |
| ⋮ |
| Stack |

→ kernel Stack

⋮

| FAR of pchar_read() |

# Get current process

\* Each process (task) have a kernel stack.
  - on x86-32 arch, kernel stack size = 8KB.
  - the kernel stack is 8KB aligned i.e.
     its addr is multiple of 8KB (in kernel space).
  - task_struct has pointer to kernel stack.

\* thread_info obj placed at the bottom of kernel
   stack (on x86).
  - it contains info about current thread &
     pointer to current task_struct.

\* current_thread_info() returns addr of current
    task's thread_info object.
         → sp & -8192

\* current macro returns current task_struct.
   #define current get_current() ⟶
   returns current process's        current_thread_info()→task
     task_struct address. ⟵

# task_struct and thread_info and kernel stack

# Synchronization in Kernel space

* multiple user processes may access kernel code (in drivers/syscalls/...) simultaneously and may cause race condition.

* Kernel provided a few sync objects:
 ① Semaphore.
 ② mutex (since 3.18+)
 ③ Spinlock

# Semaphore in kernel space | mutex in kernel space

**Semaphore** `<linux/semaphore.h>`

struct semaphore S;

| |
|---|
| lock |
| count |
| wait_list |

↳ init sem with given cnt
Sema_init (&S, init);
                        (cnt)

down (&S); → decr cnt & if
              -ve block process
      in uninterruptible Sleep

down_interruptible(&S);
↳ decr cnt & if -ve block process
   in interruptible Sleep

up (&S); → incr count & wake up blocked
              process.

---

**mutex** 3.18 + `<linux/mutex.h>`

struct mutex m;
                    =;

| |
|---|
| owner |
| wait_lock |
| wait_list |

mutex_init (&m);

mutex_lock (&m);

mutex_lock_interruptible (&m);

mutex_unlock(&m);

mutex_destroy (&m);

---

Sunbeam Infotech                    www.sunbeaminfo.com

# Semaphore Internals

## Solution 1: disable intrs

### Process 1

```
sem_p ( ) {
    disable intrs;
    count --;
    =
    enable intrs;
```

3 →

→ intr handler ( )
   {
     =
     scheduler ( );
   } dispatcher ( );
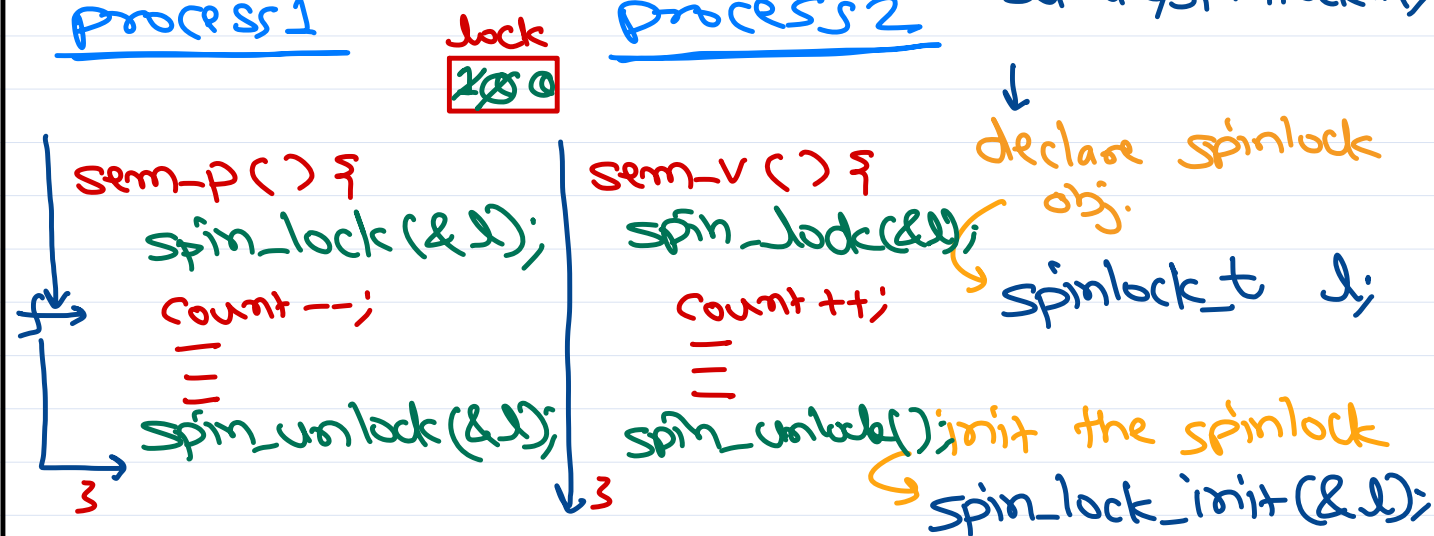
### Process 2

```
sem_v ( ) {
    disable intrs;
    count ++;
    =
    enable intrs;
```

3

Limitation: In multiprocessor env disabling intrs will disable them on current CPU only. Other process on other CPU Can still cause the race condition.

## Solution 2: spin locks

### Process 1

```
sem_p ( ) {
    spin_lock (&l);
    count --;
    =
    spin_unlock (&l);
```

3

lock
2Ø Ø

### Process 2

```
sem_v ( ) {
    spin_lock(&l);
    count ++;
    =
    spin_unlock();
```

3

#include <
linux/spinlock.h>

↓
declare spinlock
obj.
↳ spinlock_t  l;

init the spinlock
↳ spin_lock_init (&l);

lock the spinlock:
↳ spin_lock (&l);

unlock the spinlock:
↳ spin_unlock (&l);

# Spinlock

- hw based sync mechanism.
- it uses special instructions of exclusive access/bus holding.
  1. ARM7 → SWP instruction.
     - Reference: Slass book.
  2. ARM Cortex A/M →
     - LDREX, STREX, CLREX
     - Reference: Yiu book.
- these special instruction - test_and_set() kind.
  - testing (read) var value and setting var value is done in exclusive access to bus.
  - only one CPU can access bus at a time.

Spinlock pseudo code:

⬜ lock

1. spinlock_init():
   lock = 0;

2. spin_lock():
   while(lock == 1)
      ;
   lock = 1;

3. spin_unlock():
   lock = 0;

ARM instructions:

LDREX: Load var value in CPU regr and mark that addr for exclusive access.

STREX: Store CPU regr value in given var. STREX will be successful only if it is done after LDREX for that addr.

CLREX: clear exclusive access for given addr (marked by LDREX).

spinlock(lock) {     addr of lock var in RAM
   do {                        ⬜ 1
      while(__LDREX(lock) != 0)
         ;
      status = __STREX(1, lock);
   } while (status != 0);
   __DMB();
}

3

# Spinlock → in Linux Kernel.

```
#include <linux/spinlock.h>

spinlock_t   lock;

init:
        spin_lock_init (&lock);
                ↳ lock=0;
lock:
        spin_lock(&lock);
                ↳ while (lock != 0)
                   ;
                   lock = 1;
unlock:
        spin_unlock (&lock);
                ↳ lock = 0;
```

# Kernel debugging techniques

① debugging by printing

printk() → log levels

@ .... 7

EMERG...DEBUG

kernel log daemon → syslogd or klogd.

printk() → [ log buffer ] → kernel log file (kern.log or messages)

4k to 1MB

② debugging by querying

① ioctl() operation          ② procfs entry.

③ debugging by watching

strace command → shows which sys calls are called

④ System faults/hangs

① kernel OOPs message → register values + stack trace.

② kernel panic - crash - scheduler stopped

⑤ kernel debuggers

① gdb          ② kdb

# *Thank you!*

Nilesh Ghule <nilesh@sunbeaminfo.com>