Linux Kernel Programming
**Time Management**

Pierre Olivier

Systems Software Research Group @ Virginia Tech

March 19, 2017

Virginia Tech

## Outline

Virginia
Tech

Pierre Olivier (SSRG@VT)                    LKP - Time Management                    March 19, 2017      2 / 38

# Outline

1. **Kernel notion of time**

2. Tick rate and Jiffies

3. hardware clocks and timers

4. Timers

5. Delaying execution

6. Time of day

Virginia Tech

Pierre Olivier (SSRG@VT)  |  LKP - Time Management  |  March 19, 2017  3 / 38

## Kernel notion of time

- ▶ Having the notion of time passing in the kernel is essential in multiple cases:
    - ▶ Perform periodic tasks (ex: CFS time accounting)
    - ▶ Delay some processing at a relative time in the future
    - ▶ Give the time of the day
- ▶ **Absolute** vs **relative** time
- ▶ Central role of the **system timer**
    - ▶ Periodic interrupt, *system timer interrupt*
        - ▶ Update system uptime, time of day, balance runqueues, record statistics, etc.
    - ▶ Pre-programmed frequency, timer tick rate
    - ▶ tick = 1/(tick rate) seconds
- ▶ **Dynamic timers** to schedule event a relative time from now in the future

Virginia Tech

# Outline

Virginia
Tech

# Tick rate and Jiffies
Tick rate: HZ

- ▶ The **tick rate** (system timer frequency) is defined in the HZ variable
- ▶ Set to CONFIG_HZ in include/asm-generic/param.h
  - ▶ Kernel compile-time configuration option
- ▶ *Default* value is per-architecture:

| Architecture | Frequency (in Hertz) | Period (ms) |
|--------------|----------------------|-------------|
| x86          | 100                  | 10          |
| arm          | 100                  | 10          |
| Alpha        | 1024                 | 1           |
| ...          |                      |             |

Virginia Tech

# Tick rate and Jiffies
Tick rate: the ideal `HZ` value

- **High vs low system timer frequency**
- **High timer frequency pros:**
  - High *precision* for:
    - Kernel timers (finer resolution)
    - System call with timeout value (ex: `poll`)
      - Significant performance improvement for some applications
    - Timing measurements
  - *Process preemption occurs more accurately*
    - Low frequency allows processes to potentially get (way) more CPU time after the expiration of their timeslices
- **Cons:**
  - More interrupts, more overhead
    - Not very significant on modern hardware

Virginia Tech

# Tick rate and Jiffies
Tickless OS

- ▶ Option to compile the kernel as a **tickless system**
  - ▶ NO_HZ family of compilation options
- ▶ The kernel dynamically reprogram the system timer according to the current timer status
  - ▶ Situation in which there are no events for hundreds of milliseconds
- ▶ Overhead reduction
- ▶ **Energy savings**
  - ▶ CPUs spend more time in low power idle states

# Tick rate and Jiffies
Jiffies

- ▶ **jiffies** is a global variable containing the number of timer ticks since the system booted
- ▶ **unsigned long**
- ▶ include/linux/jiffies.h:

```
1 extern unsigned long volatile __jiffy_data jiffies;
```
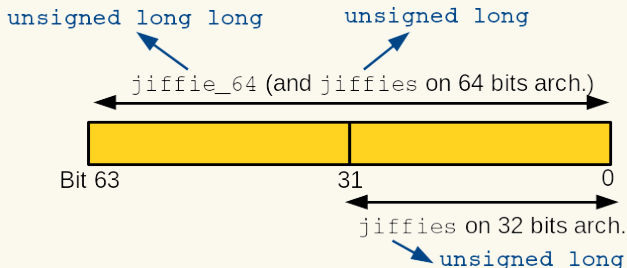
- ▶ Conversions:
  - ▶ Seconds → jiffies: (seconds * HZ)
  - ▶ jiffies → seconds: (jiffies / HZ)

```
1 unsigned long time_stamp = jiffies;      /* Now */
2 unsigned long next_tick = jiffies + 1;   /* One tick from now */
3 unsigned long later = jiffies + 5*HZ;    /* 5 seconds from now */
4 unsigned long fraction = jiffies + HZ/10; /* 100 ms from now */
```

# Tick rate and Jiffies
Jiffies: internal representation

- `unsigned long` size is 32 bits on 32 bits architectures, and 64 bits for 64 bits architectures
  - On a 32 bits variable with `HZ == 100`, overflows in 497 days
  - Still on 32 bits with `HZ == 1000`, overflows in 50 days
  - But on a 64 bits variable, no overflow for a very long time
    - Want access to a 64 bits variable while still maintaining an `unsigned long` on both architectures → linker magic

unsigned long long          unsigned long

jiffie_64 (and jiffies on 64 bits arch.)

Bit 63          31          0

jiffies on 32 bits arch.

unsigned long

Virginia Tech

# Tick rate and Jiffies
Jiffies: wraparound

- An unsigned integer going over its maximum value wraps around to zero
  - On 32 bits, $0xFFFFFFFF + 0x1 == 0x0$

```
1  unsigned long timeout = jiffies + HZ/2; /* timeout in .5 seconds */
2
3  /* do some work ... */
4
5  /* then check if we timed out */
6  if (jiffies < timeout) {
7    /* we did not time out */
8  } else {
9    /* timeout, error */
10 }
```

- If `jiffies` wraps around, chances are it will be inferior to `timeout` even in the case of an actual timeout

Virginia Tech

# Tick rate and Jiffies
Jiffies: wraparound (2)

- ▶ Macros are available in `include/linux/jiffies.h` to handle `jiffies` wraparound:

```
1 #define time_after(a,b)
2 #define time_before(a,b)
3 #define time_after_eq(a,b)
4 #define time_before_eq(a,b)
```

```
1 unsigned long timeout = jiffies + HZ/2; /* timeout in .5 seconds */
2 /* ... */
3 if (time_before(jiffies, timeout)) {
4   /* we did not time out */
5 } else {
6   /* timeout, error */
7 }
```

Virginia
Tech

# Tick rate and Jiffies

Userspace and HZ

- ▶ Values in ticks can be sent to userspace
  - ▶ Some applications grew to rely on a hard-coded value of HZ to convert in seconds
    - ▶ The fact that HZ can change caused some malfunction
- ▶ The kernel defines a constant value for the tick rate viewed from userspace: USER_HZ
  - ▶ For example it is 100 for x86
- ▶ In order to export a value in ticks (kernel space) to userspace, conversion is needed:

```
1  clock_t jiffies_to_clock(unsigned long x);
2  clock_t jiffies_64_to_clock_t(u64 x);
```

Virginia Tech

# Outline

Virginia
Tech

Pierre Olivier  (SSRG@VT)                LKP - Time Management                March 19, 2017      14 / 38

# hardware clocks and timers
RTC and the system timer

- **System timer**
  - Programmable hardware timer sending an interrupt at regular intervals
    - Programmed at boot time by the kernel to send an interrupt at `HZ` frequency
  - Other time sources on x86:
    - CPU timestamp counter (TSC) incremented every CPU clock cycle (read through `RDTSC`)
    - Local APIC (intrerrupt controller) timer
- **Real-Time Clock** (RTC):
  - Stores the **wall-clock time** (still incremented when the computer is powered off)
    - Backed-up by a small battery on the motherboard
    - Linux stores the wall-clock time in a data structure at boot time

# hardware clocks and timers
Timer interrupt processing

- ▶ Constituted of two parts: (1) architecture-dependent and (2) architecture-independent
- ▶ **Architecture-dependent** part is registered as the *handler* (top-half) for the timer interrupt
  - ▶ Generally performs those steps:
    1. Acknowledge the system timer interrupt (reset if needed)
    2. Save the wall clock time to the RTC
    3. Call the architecture independent function (still executed as part of the top-half)
- ▶ **Architecture independent part**: tick_handle_periodic()
  - ▶ Call tick_periodic()
    - ▶ Increment jiffies64
    - ▶ Update statistics for the currently running process and the entire system (load average)
    - ▶ Run dynamic timers
    - ▶ Run scheduler_tick()

Virginia Tech

# hardware clocks and timers

Timer interrupt processing: `tick_periodic()`, `do_timer`

- ▶ `kernel/`
  `time/tick-common.c`:

```
1  static void tick_periodic(int cpu)
2  {
3    if (tick_do_timer_cpu == cpu) {
4      write_seqlock(&jiffies_lock);
5
6      /* Keep track of the next tick event */
7      tick_next_period =
8        ktime_add(tick_next_period, tick_period
           );
9
10     do_timer(1); /* ! */
11     write_sequnlock(&jiffies_lock);
12     update_wall_time(); /* ! */
13   }
14
15   update_process_times(
16     user_mode(get_irq_regs())); /* ! */
17   profile_tick(CPU_PROFILING);
18 }
```

- ▶ `kernel/`
  `/time/timekeeping.c`:

```
1  void do_timer(unsigned long ticks)
2  {
3    jiffies_64 += ticks;
4    calc_global_load(ticks);
5  }
```

Virginia Tech

# hardware clocks and timers
Timer interrupt processing: `update_process_times()`

- ▶ **update_process_times()** in `kernel/timer/timer.c`

1. Call `account_process_tick()` to add one tick to the time passed:
   - ▶ In a process in user space
   - ▶ In a process in kernel space
   - ▶ In the idle task
2. Call `run_local_timers()` and run expired timers
   - ▶ Raise a softirq
3. Call `scheduler_tick()`
   - ▶ Call the `task_tick()` function of the currently running process's scheduler class
     - ▶ Update timeslices information
     - ▶ Set `need_resched` if needed
   - ▶ Perform CPU runqueues load balancing (raise the `SCHED_SOFTIRQ` softirq)

Virginia Tech

# Outline

Virginia
Tech

# Timers
Presentation

- ▶ **Timers == dynamic timers == kernel timers**
  - ▶ Used to **delay the execution of some piece of code *for a given amount of time***
    - ▶ Contrary to bottom-halves that are deferring work in a "just not now" fashion

- ▶ **struct timer_list** in includes/linux/timer.h
- ▶ entry: linked list of timers
- ▶ expires: timer expiration date in jiffies
- ▶ function: handler

```
1  struct timer_list {
2    struct hlist_node entry;
3    unsigned long expires;
4    void (*function)(unsigned long);
5    unsigned long data;
6    u32 flags;
7    /* ... */
8  }
```

- ▶ data: handler parameters

- ▶ flags: TIMER_IRQSAFE (executed with interrupts disabled), TIMER_DEFERRABLE (does not wake up an idle CPU [1])

# Timers
Using timers

- ► **Declaring**, **initializing** and **activating** a timer:

```
1  void handler_name(unsigned long data)
2  {
3    /* executed when the timer expires */
4    /* ... */
5  }
6
7  void another function(void)
8  {
9    struct timer_list my_timer;
10
11   init_time(&my_timer);                 /* initialize internal fields */
12   my_timer.expires = jiffies + 2*HZ;    /* expires in 2 secs */
13   my_timer.data = 42;                   /* 42 passed as parameter to the handler */
14   my_timer.function = handler_name;
15
16   /* activate the timer: */
17   add_timer(&my_timer);
18 }
```

- ► **Modify the expiration date** of an already running timer:

```
1  mod_timer(&my_timer, jiffies + another_delay);
```

Virginia
Tech

# Timers
## Using timers (2)

- ▶ **Deactivate a timer** prior to its expiration:

```
1 del_timer(&my_timer);
```

  - ▶ Returns 0 if the timer is already inactive, and 1 if the timer was active
  - ▶ Potential *race condition on SMP* when the handler is currently running on another core
    - ▶ Solution: del_timer_sync()

    ```
    1 del_timer_sync(&my_timer);
    ```

    - ▶ Waits for a potential currently running handler to finishes before removing the timer
    - ▶ Can be called from interrupt context only if the timer is **irqsafe** (declared with TIMER_IRQSAFE)
      - Interrupt handler interrupting the timer handler and calling del_timer_sync() → deadlock

Virginia Tech

# Timers
Race conditions

- ▶ Timers run **asynchronously** with the currently running code
    - ▶ They run in softirq context
    - ▶ Several potential race conditions exist
- ▶ Do not directly modify the `expire` field as a substitution for `mod_timer()`:

```
1  /* unsafe on SMP: */
2  del_timer(&my_timer);
3  my_timer->expires = jiffies + new_delay;
4  add_timer(&my_timer);
```

- ▶ Use `del_timer_sync()` rather than `del_timer()`
- ▶ Protect data shared by the handler and other entities

Virginia
Tech

# Timers
Implementation

- ▶ In the system timer interrupt handler, update_process_times()
  is called
    - ▶ Calls run_local_timers()
        - ▶ Raises a softirq (TIMER_SOFTIRQ)
- ▶ Softirq handler is run_timer_softirq()
    - ▶ Calls __run_timers()
        - ▶ Grab expired timers through collect_expired_timers()
        - ▶ Executes function handlers with data parameters for expired
          timers with expire_timers()
- ▶ **Timer handlers are executed in interrupt (softirq) context**

# Timers
## Example

```
1  #include <linux/module.h>
2  #include <linux/kernel.h>
3  #include <linux/init.h>
4  #include <linux/timer.h>
5
6  #define PRINT_PREF  "[TIMER_TEST] "
7
8  struct timer_list my_timer;
9
10 static void my_handler(unsigned long data)
11 {
12   printk(PRINT_PREF "handler executed!\n")
       ;
13 }
14
15 static int __init my_mod_init(void)
16 {
17   printk(PRINT_PREF "Entering module.\n");
18
19   /* initialize the timer data structure
        internal values: */
20   init_timer(&my_timer);
```

```
21   /* fill out the interesting fields: */
22   my_timer.data = 0;
23   my_timer.function = my_handler;
24   my_timer.expires = jiffies + 2*HZ; /*
        timeout == 2secs */
25
26   /* start the timer */
27   add_timer(&my_timer);
28   printk(PRINT_PREF "Timer started\n");
29
30   return 0;
31 }
32
33 static void __exit my_mod_exit(void)
34 {
35   del_timer(&my_timer);
36   printk(PRINT_PREF "Exiting module.\n");
37 }
38
39 module_init(my_mod_init);
40 module_exit(my_mod_exit);
```

Virginia Tech

# Outline

Virginia
Tech

## Delaying execution

- Sometimes the kernel needs to wait for some time without using timers (bottom-halves)
    - For example drivers communicating with the hardware
    - Needed delay can be quite small, sometimes inferior to the timer tick period
    - Several solutions:
        1. **Busy looping**
        2. **Small delays**
        3. **schedule_timeout()**

# Delaying execution
Busy looping

- ▶ **Busy looping**: spin on a loop until a given amount of ticks has elapsed

```
1  unsigned long timeout = jiffies + 10; /* timeout in 10 ticks */
2
3  while(time_before(jiffies, timeout)); /* spin until now > timeout */
```

- ▶ Can use HZ to specify a delay in seconds:

```
1  unsigned long delay = jiffies + 2*HZ; /* 2 seconds */
2
3  while(time_before(jiffies, timeout));
```

- ▶ Amount of time to wait must be a multiple of the timer period
- ▶ This technique is generally sub-optimal as the waiting process monopolizes the CPU

Virginia Tech

# Delaying execution
Busy looping (2)

▶ A better solution is to leave the CPU while waiting:

```
1  unsigned long delay = jiffies + 2*HZ;
2
3  while(time_before(jiffies, delay))
4    cond_resched();
```

▶ cond_resched() invokes the scheduler only if the need_resched flag is set
▶ Cannot be used from interrupt context (not a schedulable entity)
  ▶ Pure busy looping is probably also not a good idea from interrupt handlers as they should be fast
▶ **Busy looping can severely impact performance while a lock is help or while interrupts are disabled**

Virginia Tech

# Delaying execution
Small delays and BogoMIPS

- ▶ What if one wants to sleep for a **time inferior to the system timer period**?
    - ▶ HZ is 100 → period is 10ms
    - ▶ HZ is 1000 → period is 1ms

    - ▶ include/linux/delay.h:

    ```
    1  void mdelay(unsigned long msecs);
    2  void udelay(unsigned long usecs);
    3  void ndelay(unsigned long nsecs);
    ```

- ▶ Implemented as a busy loop
    - ▶ Kernel knows **how many loop iterations the kernel can be done in a given amount of time**: *BogoMIPS*
        - ▶ Unit: iterations / jiffy
        - ▶ Calibrated at boot time
        - ▶ Can be seen in /proc/cpuinfo

- ▶ udelay/ndelay should only be called for delays <1ms
    - ▶ Risk of overflow

# Delaying execution
schedule_timeout()

- **schedule_timeout()** put the calling task to sleep for *at least n ticks*
  - Usage:

```
1  set_current_state(TASK_INTERRUPTIBLE); /* can also use TASK_UNINTERRUPTIBLE */
2
3  schedule_timeout(2 * HZ); /* go to sleep for at least 2 seconds */
```

- Calling task must be in TASK_INTERRUPTIBLE or TASK_UNINTERRUPTIBLE otherwise calling schedule_timeout() has no effect
- schedule_timeout() should be called:
  1. From process context
  2. Without any lock held

Virginia Tech

# Delaying execution

schedule_timeout(): implementation

```
1  signed long __sched schedule_timeout(
       signed long timeout)
2  {
3    struct timer_list timer;
4    unsigned long expire;
5
6    switch (timeout)
7    {
8    case MAX_SCHEDULE_TIMEOUT:
9      schedule();
10     goto out;
11   default:
12     if (timeout < 0) {
13       printk(KERN_ERR "schedule_timeout:
           wrong timeout "
14         "value %lx\n", timeout);
15       dump_stack();
16       current->state = TASK_RUNNING;
17       goto out;
18     }
19   }
20
21   expire = timeout + jiffies;
```

```
22   setup_timer_on_stack(&timer,
       process_timeout, (unsigned long)
       current);
23   __mod_timer(&timer, expire, false);
24   schedule();
25   del_singleshot_timer_sync(&timer);
26
27   /* Remove the timer from the object
       tracker */
28   destroy_timer_on_stack(&timer);
29
30   timeout = expire - jiffies;
31
32 out:
33   return timeout < 0 ? 0 : timeout;
34 }
```

► When the timer expires,
  process_timeout()
  calls
  wake_up_process()

Virginia Tech

# Delaying execution
Sleeping on a waitqueue with a timeout

- ▶ Tasks can be placed on wait queues to wait for a specific event
- ▶ To wait for such an event *with* a timeout:
  - ▶ Call `schedule_timeout()` instead of `schedule()`

Virginia Tech

# Outline

Virginia Tech

Pierre Olivier (SSRG@VT)                LKP - Time Management                March 19, 2017        34 / 38

# Time of day
struct timespec **and** ktime_t

- ▶ Linux provides plenty of function to get / set the time of the day
- ▶ Several data structures to represent a given point in time
  - ▶ Two important ones are struct timespec **and** ktime_t

- ▶ uapi/linux/time.h:

```
1  struct timespec {
2    __kernel_time tv_sec;  /* seconds */
3    long         tv_nsec; /* nanoseconds */
4    /* __kernel_time_t is long on x86_64 */
5  }
```

- ▶ include/linux/
  ktime.h:

```
1  union ktime {
2    s64 tv64;  /* nanoseconds */
3  };
4
5  typedef union ktime ktime_t;
```

- ▶ include/
  linux/time64.h:

```
1  #define timespec64 timespec
```

Virginia Tech

# Time of day
## API usage examples

```
1  #include <linux/module.h>
2  #include <linux/kernel.h>
3  #include <linux/init.h>
4  #include <linux/timekeeping.h>
5  #include <linux/ktime.h>
6  #include <asm-generic/delay.h>
7
8  #define PRINT_PREF  "[TIMEOFDAY] "
9
10 extern void getboottime64(struct
       timespec64 *ts);
11
12 static int __init my_mod_init(void)
13 {
14   unsigned long seconds;
15   struct timespec64 ts, start, stop;
16   ktime_t kt, start_kt, stop_kt;
17
18   printk(PRINT_PREF "Entering module.\n"
       );
19
20   /* Number of seconds since the epoch
       (01/01/1970) */
21   seconds = get_seconds();
22   printk("get_seconds() returns %lu\n",
       seconds);
```

```
23   /* Same thing with seconds + nanoseconds
         using struct timespec */
24   ts = current_kernel_time64();
25   printk(PRINT_PREF "current_kernel_time64()
         returns: %lu (sec),"
26     "i %lu (nsec)\n", ts.tv_sec, ts.tv_nsec);
27
28   /* Get the boot time offset */
29   getboottime64(&ts);
30   printk(PRINT_PREF "getboottime64() returns:
         %lu (sec),"
31     "i %lu (nsec)\n", ts.tv_sec, ts.tv_nsec);
32
33   /* The correct way to print a struct
         timespec as a single value: */
34   printk(PRINT_PREF "Boot time offset: %lu.%09
         lu secs\n", ts.tv_sec, ts.tv_nsec);
35   /* Otherwise, just using %lu.%lu transforms
         this:
36    * ts.tv_sec  == 10
37    * ts.tv_nsec == 42
38    * into: 10.42 rather than 10.000000042 */
```

Virginia Tech

# Time of day
## API usage examples (2)

```
39  /* another interface using ktime_t (
          number of nsec since boot) */
40  kt = ktime_get();
41  printk(PRINT_PREF "ktime_get() returns
          %llu\n", kt.tv64);
42
43  /* Subtract two struct timespec */
44  getboottime64(&start);
45  stop = current_kernel_time64();
46  ts = timespec64_sub(stop, start);
47  printk(PRINT_PREF "Uptime: %lu.%09lu
          secs\n", ts.tv_sec, ts.tv_nsec);
48
49  /* measure the execution time of a
          piece of code */
50  start_kt = ktime_get();
51  udelay(100);
52  stop_kt = ktime_get();
53
54  kt = ktime_sub(stop_kt, start_kt);
55  printk(PRINT_PREF "Measured execution
          time: %llu usecs\n", (kt.tv64)
          /1000);
56
57  return 0;
58  }
```

```
59  static void __exit my_mod_exit(void)
60  {
61    printk(PRINT_PREF "Exiting module.\n");
62  }
63
64  module_init(my_mod_init);
65  module_exit(my_mod_exit);
66
67  MODULE_LICENSE("GPL");
```

```
1  obj-m += timeofday.o
2
3  all:
4    make -C /lib/modules/$(shell uname -r)/
          build M=$(PWD) modules
5
6  test: all
7    sudo rmmod timeofday.ko &> /dev/null ||
          true
8    sudo insmod timeofday.ko
9    sudo rmmod timeofday.ko
10
11  clean:
12    make -C /lib/modules/$(shell uname -r)/
          build M=$(PWD) clean
```

# Bibliography I

[1] CORBET, J.
Deferrable timers.
https://lwn.net/Articles/228143/.
Accessed: 2017-03-14.

Virginia Tech