

# Systems Programming

## Kernel Development

H. Turgut Uyar Şima Uyar

2001-2014

1 / 44

## Topics

### Kernel

Architecture  
Kernel Development  
Kernel Modules

### Process Management

Data Structures  
Synchronization  
Scheduling

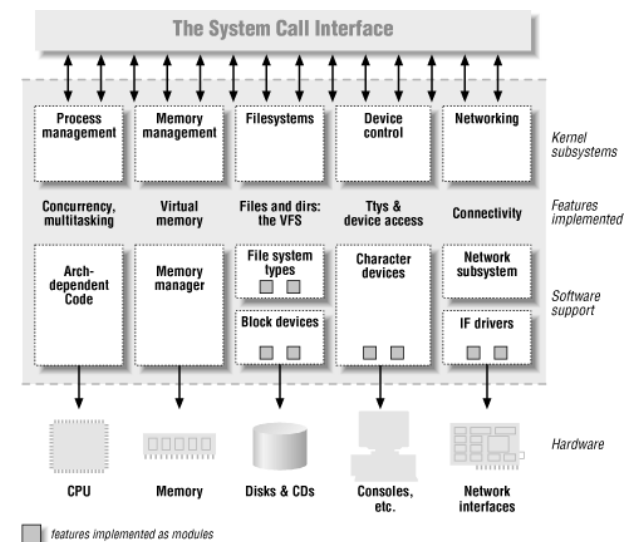
2 / 44

## Kernel

- ▶ provides programs with a consistent view of the hardware
- ▶ protects against unauthorized access to resources
- ▶ kernel runs in supervisor mode (kernel space), applications run in user mode (user space)
- ▶ switching to kernel space:
  - ▶ system calls: synchronous, in the process context
  - ▶ interrupts: asynchronous

3 / 44

## Kernel Subsystems



4 / 44

## Kernel Subsystems

- ▶ process management
  - ▶ creating and destroying processes
  - ▶ communication between processes
  - ▶ scheduling
- ▶ memory management
  - ▶ virtual address space for each process
- ▶ filesystems
  - ▶ structured filesystem on top of unstructured hardware
- ▶ device control
- ▶ networking
  - ▶ delivering data packets across program and network interfaces
  - ▶ routing and address resolution

5 / 44

## Kernel Architecture

- ▶ **monolithic**: all functionality in one big chunk of code
- ▶ **microkernel**: organized as layers
  - ▶ most functionality in user space
  - ▶ too much communication overhead

6 / 44

## Kernel Development

- ▶ recompile the kernel
- ▶ reboot the computer
- ▶ test the new kernel
- ▶ reboot to the original kernel
- ▶ very slow development cycle!
- ▶ no external libraries

7 / 44

## Example: Adding a System Call

- ▶ add an entry to the system call table:  
system call number, name, function to invoke, ...
- ▶ add prototype to the system calls header file
- ▶ implement system call

8 / 44

## Example: Adding a System Call

- ▶ new system call: add two integers
- ▶ add an entry to the system call table

arch/x86/kernel/syscall\_table\_32.S

```
.long sys_mycall
```

- ▶ append entry for system call

arch/x86/include/asm/unistd\_32.h

```
#define __NR_mycall 333
```

9 / 44

## Example: Adding a System Call

- ▶ add prototype to the system calls header file

include/linux/syscalls.h

```
asmlinkage int sys_mycall(int i, int j);
```

- ▶ implement system call

mycall.c

```
asmlinkage int sys_mycall(int i, int j)
{
    return i + j;
}
```

10 / 44

## Example: Test Program

```
#define __NR_mycall 333

int main(int argc, char **argv)
{
    int x1 = 10, x2 = 20, y;

    y = syscall(__NR_mycall, x1, x2);
    printf("%d\n", y);
    return 0;
}
```

11 / 44

## Data Transfer

- ▶ special functions for transferring data between kernel space and user space
- ▶ kernel → user:  
copy\_to\_user(user\_buf, kernel\_buf, length)
- ▶ user → kernel:  
copy\_from\_user(kernel\_buf, user\_buf, length)

12 / 44

## Example: Data Transfer

- ▶ new system call: get the time passed since 1970
- ▶ kernel structure for representing time

```
struct timeval {  
    long tv_sec;    /* seconds */  
    long tv_usec;   /* microseconds */  
};
```

- ▶ global variable that keeps the current time

```
struct timeval xtime;
```

13 / 44

## Example: Data Transfer

```
asm linkage int sys_ptime(struct timeval *tm)  
{  
    copy_to_user(tm, &xtime, sizeof(struct timeval));  
    return 0;  
}
```

14 / 44

## Example: Test Program

```
#define __NR_ptime 334  
  
int main(int argc, char **argv)  
{  
    struct timeval utime;  
    int res;  
  
    res = syscall(__NR_ptime, &utime);  
    printf("%d\n", (int) utime.tv_sec);  
    sleep(2);  
    res = syscall(__NR_ptime, &utime);  
    printf("%d\n", (int) utime.tv_sec);  
    return 0;  
}
```

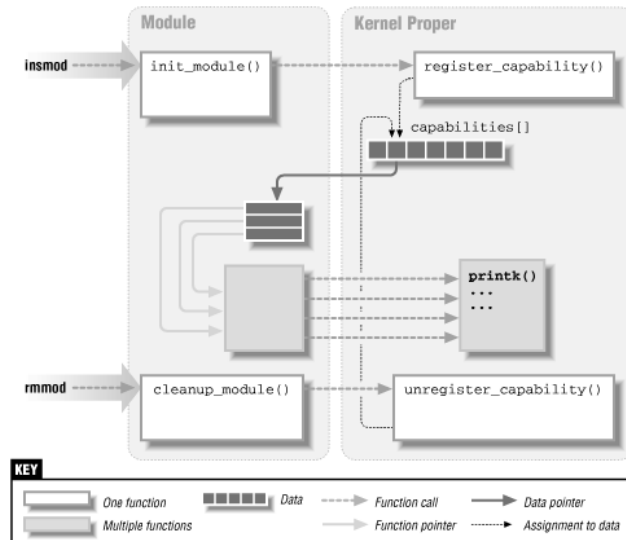
15 / 44

## Modular Kernel

- ▶ monolithic architecture
- ▶ modules added or removed at runtime
- ▶ no need to reboot: faster development cycle

16 / 44

## Module Registry



17 / 44

## Example: Hello, world!

```
#include <linux/init.h>
#include <linux/module.h>

MODULE_LICENSE("Dual BSD/GPL");

static int hello_init(void) { ... }

static void hello_exit() { ... }

module_init(hello_init);
module_exit(hello_exit);
```

18 / 44

## Example: Hello, world!

```
static int hello_init(void)
{
    printk(KERN_ALERT "Hello, world!\n");
    return 0;
}

static void hello_exit()
{
    printk(KERN_ALERT "Goodbye, cruel world!\n");
}
```

19 / 44

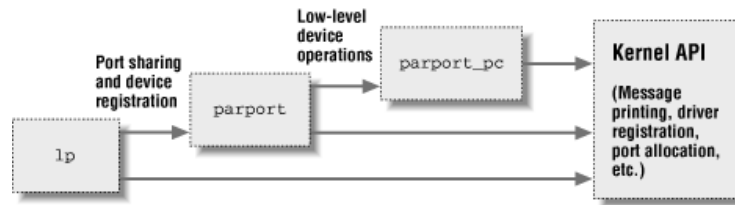
## Kernel Symbol Table

- ▶ kernel symbol table contains addresses of global symbols
- ▶ when loading a module:
- ▶ unresolved symbols are linked to the kernel symbol table
- ▶ exported symbols become part of the kernel symbol table

20 / 44

## Module Stacking

- ▶ modules can use symbols exported by other modules



21 / 44

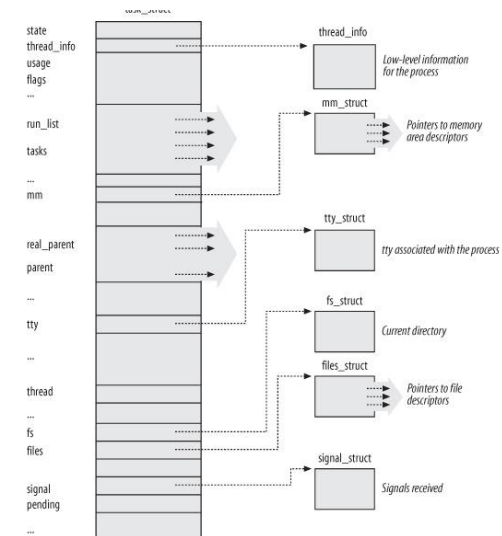
## Reading Material

- ▶ Corbet-Rubini-Hartman, 3/e
  - ▶ Chapter 2: Building and Running Modules

22 / 44

## Process Descriptor

- ▶ a process descriptor for each process:  
`struct task_struct`
- ▶ process state
- ▶ process identification (pid, uid, euid, ...)



23 / 44

## Process Descriptor

24 / 44

## Process List

- ▶ doubly linked list of all process descriptors
- ▶ `tasks` field of `task_struct`
- ▶ `current` macro gives the process descriptor of the running process: e.g. `current->pid`

25 / 44

## Process 0

- ▶ a.k.a. the *idle process* or the *swapper*
- ▶ the first entry in the process list
- ▶ created during the initialization stage of the kernel
- ▶ the only process created without using the *fork* system call
- ▶ the ancestor of all processes

26 / 44

## Process 0

- ▶ uses a statically allocated data structure
  - ▶ process descriptor stored in the `init_task` variable
  - ▶ initialized by the `INIT_TASK` macro
- ▶ executes the `start_kernel()` function
  - ▶ initializes all data structures needed by kernel
  - ▶ enables interrupts
  - ▶ creates process 1 (commonly known as the *init process*)
- ▶ executes the `cpu_idle()` function

27 / 44

## Creating Processes

- ▶ *fork* is implemented as the `clone` system call
- ▶ `do_fork()` function handles the `clone` system call:
- ▶ allocates a `pid` for the child process
- ▶ uses `copy_process()` to set up the process descriptor and other kernel data structures for new process
  - ▶ uses `dup_task_struct()` to allocate a new process descriptor and to copy parent process' process descriptor info
- ▶ adjusts some parameters of parent and child processes
- ▶ returns `pid` of child process

28 / 44

## Destroying Processes

- ▶ through the `_exit()` system call
- ▶ uses the `do_exit()` function

29 / 44

## Synchronization

- ▶ critical sections and race conditions also exist for kernel code
- ▶ synchronization is needed
- ▶ several kernel level synchronization primitives
- ▶ primitive must be chosen based on requirements of operation

30 / 44

## Synchronization Primitives

- ▶ atomic read-modify-write operations
- ▶ memory barriers (to avoid instruction reordering)
- ▶ spin locks (locks with busy waiting)
- ▶ kernel semaphores (lock with blocking wait)
- ▶ interrupt disabling (local CPU)

31 / 44

## Atomic Operations

- ▶ instructions that execute atomically
- ▶ no interrupts
- ▶ to implement counters
- ▶ to atomically perform an operation and test results:  
e.g. `atomic_dec_and_test`

```
typedef struct {  
    volatile int counter;  
} atomic_t;
```

32 / 44



## Memory Barriers

- ▶ kernel may reorder assembly instructions for optimization
- ▶ reordering must be avoided when synchronization is needed
- ▶ barrier ensures that the instructions before the primitive are completed before the instructions after the primitive
- ▶ read memory barrier: `rmb()`
- ▶ write memory barrier: `wmb()`
- ▶ memory barrier: `barrier()` – same as `wmb()`

33 / 44

## Spin Locks

- ▶ for locking access to shared data (critical sections)
- ▶ for multiprocessor environments
- ▶ uses busy waiting
  - ▶ kernel resources usually locked for very short periods
  - ▶ more time consuming to release and reacquire cpu
- ▶ represented by a `spinlock_t` structure
- ▶ macros used for working with spin locks
- ▶ read and write spin locks to increase concurrency: `rwlock_t` structure

34 / 44

## Semaphores

- ▶ sleeping locks
- ▶ suited for locks that are held for a long time
- ▶ not optimal for locks that are held for short periods
- ▶ kernel preemption not disabled, i.e. no adverse effects on scheduling latency
- ▶ allows arbitrary number of simultaneous lock holders: counting semaphores
- ▶ two atomic operations: `P()` - `V()`  
`down()` - `up()`

35 / 44

## Scheduling

- ▶ divide the finite resource of processor time between the runnable processes on the system
- ▶ conflicting goals:
  - ▶ fast process response time (low latency)
  - ▶ maximal system utilization (high throughput)
- ▶ processor bound processes - I/O bound processes
  - ▶ Linux favors I/O bound processes, i.e. optimizes for low latency

36 / 44

## O(1) Scheduler

- ▶ constant-time algorithm for timeslice calculation and per processor runqueues
- ▶ scalable
- ▶ ideal for large server workloads
- ▶ problems for interactive processes

37 / 44

## CFS Scheduler

- ▶ Completely Fair Scheduler
- ▶ aims at improving scheduling for interactive processes

38 / 44

## Linux Scheduler

- ▶ different algorithms to schedule different types of processes
- ▶ scheduler classes with priorities
- ▶ iterate over each scheduler class in order of priority
- ▶ CFS for normal processes
- ▶ two policies for real time processes:
  - ▶ SCHED\_FIFO
  - ▶ SCHED\_RR

39 / 44

## CFS

- ▶ assign processes a *proportion* of processor
- ▶ nice value (priority) acts as weight to determine proportion of processor time
- ▶ preemptive (based on proportions of processor time consumed)

40 / 44

## CFS

- ▶ *timeslice* proportional to process' weight over sum of weights of all runnable processes
- ▶ targeted latency
- ▶ minimum granularity

41 / 44

## CFS Implementation

- ▶ for process accounting:  
`struct sched_entity`
- ▶ member of `struct task_struct`
- ▶ *virtual runtime* (vruntime):  
actual runtime (in ns) of a process  
normalized by the number of runnable processes
- ▶ in a perfectly multitasking system all processes should have the same virtual runtime
- ▶ updated periodically by the system timer  
and also whenever a process becomes runnable or is blocked

42 / 44

## CFS Implementation

- ▶ the runnable process with the smallest vruntime is selected to run
- ▶ red-black tree to manage list of runnable processes:  
search in  $O(\log n)$ 
  - ▶ leftmost node has lowest vruntime
  - ▶ leftmost node is cached

43 / 44

## Reading Material

- ▶ Linux Kernel Development, 3rd Edition
  - ▶ Author: Robert Love
  - ▶ Publisher: Addison-Wesley Professional
  - ▶ Year: 2010
  - ▶ Chapters: 3, 4, 5, 9 and 10
  - ▶ accessible on Safari e-books through the ITU Library

44 / 44