

Systems Programming

Devices

H. Turgut Uyar Şima Uyar

2001-2014

1 / 76

Topics

I/O Subsystem

Introduction
Device Types
I/O Software
Accessing Devices

Device Drivers

Interface
Implementation
Device Access

2 / 76

I/O Devices

- ▶ O/S controls all I/O devices
- ▶ issues commands to devices
- ▶ catches interrupts
- ▶ handles errors
- ▶ provides interface

3 / 76

Device Controllers

- ▶ devices consist of:
 - ▶ mechanical components
 - ▶ electronic components: *device controller*
- ▶ O/S deals with controller
 - ▶ connected through a standard interface
 - ▶ SCSI, USB, Firewire, ...

4 / 76

Controller Registers

- ▶ CPU communicates with the controller through registers
- ▶ data register: for sending/receiving data
- ▶ control register: for sending commands to device
- ▶ status register: for getting/setting the state of device

5 / 76

I/O Architecture

- ▶ **ports**: special address space for I/O
 - ▶ separate lines for I/O ports
 - ▶ special instructions for I/O
- ▶ **memory-mapped**: registers part of regular address space
 - ▶ directly-mapped: part of address space reserved for I/O
 - ▶ software-mapped: I/O space part of virtual memory

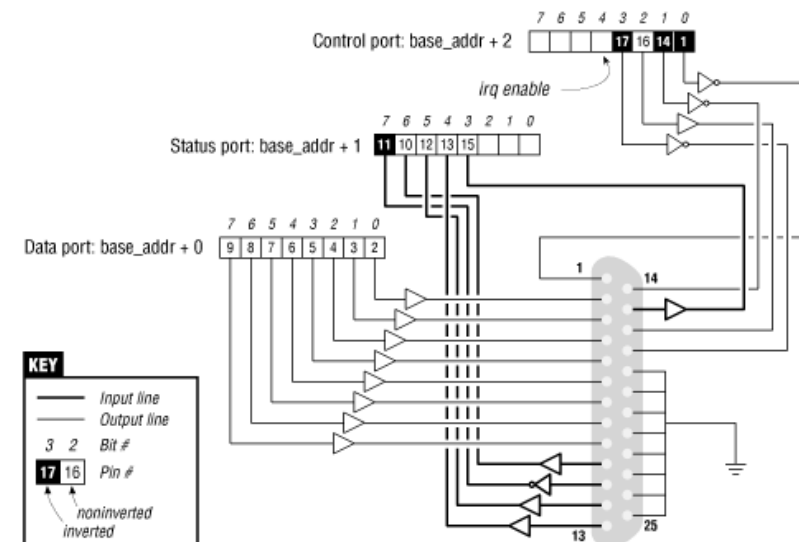
6 / 76

PC Parallel Interface

- ▶ parallel interface base addresses on a PC: 0x378, 0x278
- ▶ ports:
 - ▶ +0: bidirectional data register
 - ▶ +1: status register (read-only)
online, out-of-paper, busy
 - ▶ +2: control register (write-only)
enable/disable interrupts

7 / 76

Parallel Interface



8 / 76

Device Types

- ▶ character devices
- ▶ block devices
- ▶ network interfaces
- ▶ clocks and timers

9 / 76

Character Devices

- ▶ a character device acts like a stream of characters
- ▶ arbitrary-sized data transfer
- ▶ not addressable: no seek operation

examples

- ▶ console, mouse
- ▶ sound card
- ▶ serial port, parallel port

10 / 76

Block Devices

- ▶ a block device can host a filesystem
- ▶ data transfer in fixed-size blocks
- ▶ each block has its own address
- ▶ read/write each block independently

example

- ▶ disks

11 / 76

Device Type

- ▶ the device type is more the characteristic of the driver rather than the device itself

example: disk

- ▶ usually a block device
- ▶ it can also be used as a character device: tar

12 / 76

I/O Software

- ▶ blocking vs interrupt-driven
 - ▶ better for CPU to work interrupt-driven fashion
 - ▶ better for user-space programs to work in blocking fashion
 - ▶ easier to develop programs that work in blocking fashion
 - ▶ O/S makes interrupt-driven operations look blocking
- ▶ standardized interface
- ▶ uniform naming

13 / 76

Unix Device Naming

- ▶ in Unix, every device has a **device node**
- ▶ under the /dev folder
- ▶ /dev/sda: first SCSI disk
- ▶ /dev/sdb: second SCSI disk
- ▶ /dev/sdb1: first partition of the second SCSI disk
- ▶ /dev/sdb2: second partition of the second SCSI disk
- ▶ /dev/parport0: first parallel port

14 / 76

Unix Device Naming

- ▶ device nodes have major and minor numbers
- ▶ major number identifies the driver
- ▶ minor number identifies the physical device
- ▶ all /dev/sd* devices have the same major number
- ▶ they all have different minor numbers
- ▶ (recently) major number alone doesn't identify driver
- ▶ major number + region of minor numbers

15 / 76

I/O Services

- ▶ copy semantics: transfer the snapshot of data at the time of the I/O request
- ▶ scheduling: issue order may not be the best execution order
- ▶ buffering: adapt between different data transfer sizes
- ▶ caching
- ▶ spooling: deal with dedicated devices (e.g. printers)
 - ▶ a daemon for controlling the device
 - ▶ a spooling directory
- ▶ error handling

16 / 76

Software Layers

- ▶ top-down:
- ▶ user-space applications
- ▶ device-independent software
- ▶ device drivers
- ▶ interrupt handlers

17 / 76

Device-Independent Software

- ▶ functions common to all devices
- ▶ uniform interface to user-level software
- ▶ device naming
- ▶ device protection
- ▶ provide device-independent block sizes
- ▶ buffering
- ▶ allocating and releasing dedicated devices
- ▶ error reporting

18 / 76

Device Drivers

- ▶ device-dependent code
- ▶ a driver for each device type
- ▶ accept request from device-independent software
- ▶ decide on sequence of controller operations

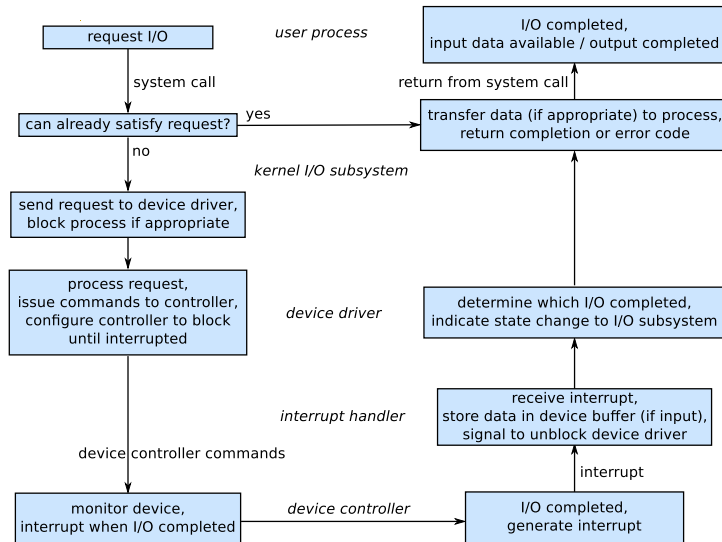
19 / 76

Interrupt Handlers

- ▶ interrupts hidden from rest of system
 - ▶ requesting process is blocked until I/O is completed
- ▶ when I/O is completed, interrupt occurs
 - ▶ process is made to unblock

20 / 76

I/O Life Cycle



21 / 76

Accessing Devices

- ▶ directly: using ports or memory
- ▶ through device drivers: using the device driver interface

22 / 76

Direct Access

- ▶ input: `inb`, `inw`, `inl`
- ▶ output: `outb`, `outw`, `outl`
- ▶ get permission from O/S: `ioperm` system call

```
int ioperm(unsigned long from,
           unsigned long num,
           int turn_on);
```

23 / 76

Direct Access Example

output to parallel interface

```
ioperm(0x378, 1, 255);
outb(0xff, 0x378);
```

24 / 76

Reading Material

- ▶ Silberschatz, 8/e
 - ▶ Chapter 13: **I/O Systems**

25 / 76

Unix Device Driver Interface

- ▶ in Unix, the device driver interface is similar to the file interface
- ▶ open, close
- ▶ read, write

26 / 76

Device Specific Operations

- ▶ some operations are neither read nor write
- ▶ ioctl: issue command specific to device

device-specific operation examples

- ▶ eject CDROM
- ▶ make the speaker beep
- ▶ set communication parameters for modem

27 / 76

System Calls

- ▶ open:

```
int open(const char *pathname,  
         int flags,  
         mode_t mode);
```

- ▶ flags:
 - ▶ O_RDONLY O_WRONLY O_RDWR
 - ▶ O_CREAT O_APPEND
- ▶ mode: permissions
- ▶ returns: file descriptor

28 / 76

System Calls

► close:

```
int close(int fd);
```

- returns: success / failure status

29 / 76

System Calls

► read:

```
ssize_t read(int fd,  
             void *buf,  
             size_t count);
```

► returns: number of bytes read (x)

- $x = \text{count}$: successful completion
- $x = 0$: end-of-file
- $x < 0$: error
- $0 < x < \text{count}$: partial transfer, retry remaining part

30 / 76

System Calls

► write:

```
ssize_t write(int fd,  
             const void *buf,  
             size_t count);
```

► returns: number of bytes written (x)

- $x = \text{count}$: successful completion
- $x = 0$: end-of-file
- $x < 0$: error
- $0 < x < \text{count}$: partial transfer, retry remaining part

31 / 76

System Calls

► ioctl:

```
int ioctl(int fd,  
         int request,  
         ...);
```

- parameter and return values depend on request

32 / 76

Device Access Example

output to parallel port

```
fd = open("/dev/parport0", O_WRONLY);
if (fd == -1)
{
    perror("cannot access device");
    exit(EXIT_FAILURE);
}
write(fd, buffer, len);
close(fd);
```

33 / 76

Device Specific Command Example

make the speaker beep

```
fd = open("/dev/console", O_RDWR);
status = ioctl(fd, KDMKTONE, 0x100011AA);
if (status == -1)
{
    perror("cannot generate beep");
    exit(EXIT_FAILURE);
}
close(fd);
```

34 / 76

Implementing Device Drivers

- ▶ implement system calls for device
- ▶ convert system calls to device specific I/O instructions

35 / 76

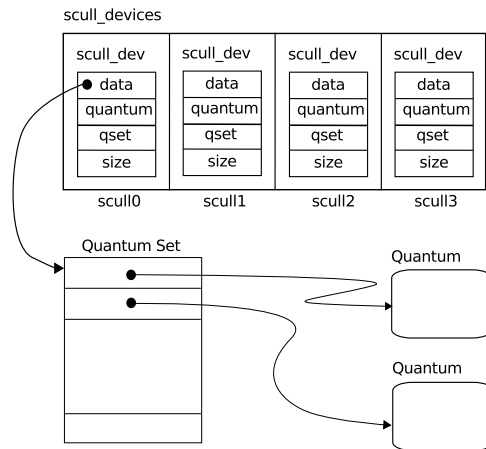
Device Driver Example

simplified scull

- ▶ use memory as device
 - ▶ /dev/scull0
 - ▶ /dev/scull1
- ▶ each device can hold data up to a limit
 - ▶ data persists during module's lifetime

36 / 76

Memory Layout



- ▶ each device has a quantum set
- ▶ each quantum contains the actual data
- ▶ memory is allocated as data is written

37 / 76

Global Definitions

`scull.h`

```
#define SCULL_MAJOR 0
#define SCULL_NR_DEVS 4
#define SCULL_QUANTUM 4000
#define SCULL_QSET 1000
```

38 / 76

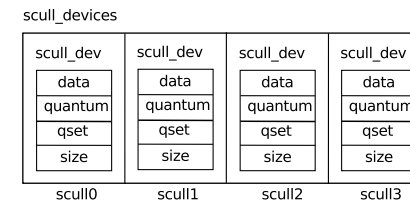
Module Parameters

```
int scull_major = SCULL_MAJOR;
int scull_minor = 0;
int scull_nr_devs = SCULL_NR_DEVS;
int scull_quantum = SCULL_QUANTUM;
int scull_qset = SCULL_QSET;

module_param(scull_major, int, S_IRUGO);
module_param(scull_minor, int, S_IRUGO);
module_param(scull_nr_devs, int, S_IRUGO);
module_param(scull_quantum, int, S_IRUGO);
module_param(scull_qset, int, S_IRUGO);
```

39 / 76

Data Structures



```
struct scull_dev {
    char **data;
    int quantum;
    int qset;
    unsigned long size;
    struct semaphore sem;
    struct cdev cdev;
};

struct scull_dev *scull_devices;
```

40 / 76

Module Initialization

- ▶ allocate I/O region
 - ▶ base address
 - ▶ number of ports
- ▶ register driver with the kernel
 - ▶ major and minor numbers
 - ▶ capabilities: file operations

41 / 76

Module Initialization

driver registration: major and minor numbers

```
if (scull_major)
{
    devno = MKDEV(scull_major, scull_minor);
    result = register_chrdev_region(devno,
                                    scull_nr_devs, "scull");
}
else /* dynamic */
{
    result = alloc_chrdev_region(&devno,
                                scull_minor, scull_nr_devs, "scull");
    scull_major = MAJOR(devno);
}
```

42 / 76

Module Initialization

data structure allocation

```
scull_devices = kmalloc(scull_nr_devs *
                        sizeof(struct scull_dev), GFP_KERNEL);
if (!scull_devices)
{
    result = -ENOMEM;
    goto fail;
}
memset(scull_devices, 0,
       scull_nr_devs * sizeof(struct scull_dev));
```

43 / 76

File Operations

- ▶ map system calls to functions:
`struct file_operations`

```
struct file_operations scull_fops = {
    .open    = scull_open,
    .release = scull_release,
    .read    = scull_read,
    .write   = scull_write,
    .llseek  = scull_llseek,
    .ioctl   = scull_ioctl,
};
```

44 / 76

Module Initialization

driver activation

```
for (i = 0; i < scull_nr_devs; i++)
{
    dev = &scull_devices[i];
    dev->quantum = scull_quantum;
    dev->qset = scull_qset;
    init_Mutex(&dev->sem);

    devno = MKDEV(scull_major, scull_minor + i);
    cdev_init(&dev->cdev, &scull_fops);
    dev->cdev.owner = THIS_MODULE;
    dev->cdev.ops = &scull_fops;
    cdev_add(&dev->cdev, devno, 1);
}
```

45 / 76

Module Cleanup

```
dev_t devno = MKDEV(scull_major, scull_minor);

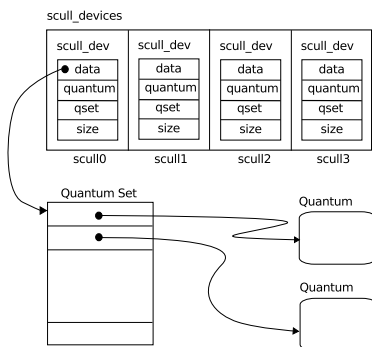
if (scull_devices)
{
    for (i = 0; i < scull_nr_devs; i++)
    {
        scull_trim(scull_devices + i);
        cdev_del(&scull_devices[i].cdev);
    }
    kfree(scull_devices);
}

unregister_chrdev_region(devno, scull_nr_devs);
```

46 / 76

Module Cleanup

data structure deallocation



```
if (dev->data)
{
    for (i = 0; i < dev->qset; i++)
    {
        if (dev->data[i])
            kfree(dev->data[i]);
    }
    kfree(dev->data);
}
dev->data = NULL;
dev->quantum = scull_quantum;
dev->qset = scull_qset;
dev->size = 0;
```

47 / 76

Kernel Data Structures

- ▶ a structure for each device node:
`struct inode`
- ▶ a structure for each open file:
`struct file`
 - ▶ `f_mode`: readable, writable, both
 - ▶ `f_pos`: current reading/writing position
 - ▶ `f_flags`
 - ▶ `f_op`: operations associated with the file
 - ▶ `private_data`: pointer to allocated data

48 / 76

Open

- ▶ identify actual device
- ▶ check for device-specific errors
- ▶ initialize device
- ▶ allocate and initialize data structures

49 / 76

Kernel System Call Interface

- ▶ open system call:

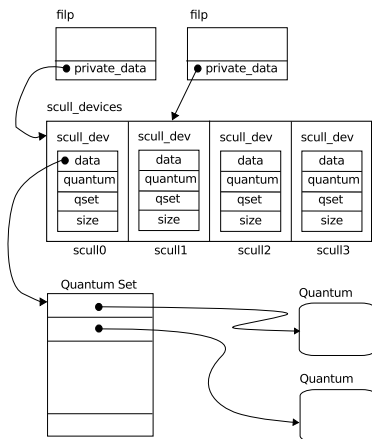
```
int open(const char *pathname,  
         int flags,  
         mode_t mode);
```

- ▶ kernel function to implement:

```
int foo_open(struct inode *inode,  
             struct file *filp);
```

50 / 76

Open



```
struct scull_dev *dev;  
  
dev = container_of(  
    inode->i_cdev,  
    struct scull_dev,  
    cdev  
);  
filp->private_data = dev;
```

51 / 76

Kernel System Call Interface

- ▶ close system call:

```
int close(int fd);
```

- ▶ kernel function to implement:

```
int foo_release(struct inode *inode,  
               struct file *filp);
```

52 / 76

Kernel System Call Interface

- ▶ write system call:

```
ssize_t write(int fd,
              const void *buf,
              size_t count);
```

- ▶ kernel function to implement:

```
ssize_t foo_write(struct file *filp,
                  const char __user *buf,
                  size_t count,
                  loff_t *f_pos);
```

53 / 76

Write

```
struct scull_dev *dev = filp->private_data;
ssize_t retval = -ENOMEM;
```

```
if (down_interruptible(&dev->sem))
    return -ERESTARTSYS;
```

```
/* determine position */
/* allocate quantum if necessary */
/* copy from user space */
/* update size */
```

out:

```
up(&dev->sem);
return retval;
```

54 / 76

Write

determine position

```
int quantum = dev->quantum, qset = dev->qset;
int s_pos, q_pos;
```

```
if (*f_pos >= quantum * qset)
{
    retval = 0;
    goto out;
}
```

```
s_pos = (long) *f_pos / quantum;
q_pos = (long) *f_pos % quantum;
```

55 / 76

Write

allocate quantum if necessary

```
if (!dev->data)
{
    dev->data = kmalloc(qset * sizeof(char *),
                       GFP_KERNEL);

    if (!dev->data)
        goto out;
    memset(dev->data, 0, qset * sizeof(char *));
}
if (!dev->data[s_pos])
{
    dev->data[s_pos] = kmalloc(quantum, GFP_KERNEL);
    if (!dev->data[s_pos])
        goto out;
}
```

56 / 76

Write

copy from user space

```
/* adjust write amount */
if (count > quantum - q_pos)
    count = quantum - q_pos;

if (copy_from_user(dev->data[s_pos] + q_pos,
                    buf, count))
{
    retval = -EFAULT;
    goto out;
}
```

57 / 76

Write

update size

```
*f_pos += count;
retval = count;

if (dev->size < *f_pos)
    dev->size = *f_pos;
```

58 / 76

Kernel System Call Interface

► read system call:

```
ssize_t read(int fd,
             void *buf,
             size_t count);
```

► kernel function to implement:

```
ssize_t foo_read(struct file *filp,
                 char __user *buf,
                 size_t count,
                 loff_t *f_pos);
```

59 / 76

Read

```
struct scull_dev *dev = filp->private_data;
ssize_t retval = 0;
```

```
if (down_interruptible(&dev->sem))
    return -ERESTARTSYS;
```

```
/* determine position */
/* copy to user space */
```

```
out:
    up(&dev->sem);
    return retval;
```

60 / 76

Read

determine position

```
int quantum = dev->quantum;
int s_pos, q_pos;

if (*f_pos >= dev->size)
    goto out;
if (*f_pos + count > dev->size)
    count = dev->size - *f_pos;

s_pos = (long) *f_pos / quantum;
q_pos = (long) *f_pos % quantum;

if (dev->data == NULL || ! dev->data[s_pos])
    goto out;
```

61 / 76

Reading from the Device

copy to user space

```
/* adjust read amount */
if (count > quantum - q_pos)
    count = quantum - q_pos;

if (copy_to_user(buf, dev->data[s_pos] + q_pos, count))
{
    retval = -EFAULT;
    goto out;
}
*f_pos += count;
retval = count;
```

62 / 76

Kernel System Call Interface

► lseek system call:

```
off_t lseek(int fd,
            off_t offset,
            int whence);
```

► kernel function to implement:

```
loff_t foo_llseek(struct file *filp,
                  loff_t off,
                  int whence);
```

63 / 76

Seek

calculate new position

```
switch(whence)
{
    case 0: /* SEEK_SET */
        newpos = off;
        break;
    case 1: /* SEEK_CUR */
        newpos = filp->f_pos + off;
        break;
    case 2: /* SEEK_END */
        newpos = dev->size + off;
        break;
    default: /* can't happen */
        return -EINVAL;
}
```

64 / 76

Seek

set new position

```
if (newpos < 0)
    return -EINVAL;
filp->f_pos = newPos;
return newPos;
```

65 / 76

Kernel System Call Interface

- ▶ ioctl system call:

```
int ioctl(int fd,
          int request,
          ...);
```

- ▶ kernel function to implement:

```
int foo_ioctl(struct inode *inode,
              struct file *filp,
              unsigned int cmd,
              unsigned long arg)
```

66 / 76

Device-Specific Commands

- ▶ SCULL_IOCTLRESET:
assign default values to quantum set size and quantum size
- ▶ SCULL_IOCSEQUANTUM: set quantum size from pointer
- ▶ SCULL_IOCTLQUANTUM: (tell) set quantum size from value
- ▶ SCULL_IOCQQUANTUM: get quantum size to pointer
- ▶ SCULL_IOCQQUANTUM: (query) return quantum size
- ▶ SCULL_IOCXQUANTUM: (exchange) set + get
- ▶ SCULL_IOCHQUANTUM: (shift) tell + query
- ▶ similar operations for quantum set size

67 / 76

Device Operations

```
switch(cmd)
{
    case SCULL_IOCTLRESET:
        scull_quantum = SCULL_QUANTUM;
        scull_qset = SCULL_QSET;
        break;

    /* other cases */
}
```

68 / 76

Device Operations

setting quantum size

```
case SCULL_IOCSEQUANTUM:
    if (!capable (CAP_SYS_ADMIN))
        return -EPERM;
    retval = __get_user(scull_quantum,
                       (int __user *) arg);

    break;

case SCULL_I OCTQUANTUM:
    if (!capable (CAP_SYS_ADMIN))
        return -EPERM;
    scull_quantum = arg;
    break;
```

69 / 76

Device Operations

getting quantum size

```
case SCULL_IOC GQUANTUM:
    retval = __put_user(scull_quantum,
                       (int __user *) arg);

    break;

case SCULL_IOC QQUANTUM:
    return scull_quantum;
```

70 / 76

Device Driver Example

short: read/write I/O ports

- ▶ each device node accesses a different port:
 - ▶ /dev/short0: port at base
 - ▶ /dev/short1: port at base+1
- ▶ module parameters:
 - ▶ major number (default dynamic)
 - ▶ base address (default 0x378)

71 / 76

Region Allocation

module initialization

```
if (!request_region(short_base, SHORT_NR_PORTS,
                   "short"))
{
    return -ENODEV;
}
```

module cleanup

```
release_region(short_base, SHORT_NR_PORTS);
```

72 / 76

Read

```
int retval = count;
int minor = iminor(filp->f_dentry->d_inode);
unsigned long port = short_base + (minor & 0x0f);
unsigned char *kbuf, *ptr;

kbuf = kmalloc(count, GFP_KERNEL);
if (!kbuf)
    return -ENOMEM;

/* do the I/O */

kfree(kbuf);
return retval;
```

73 / 76

Read

do the I/O

```
ptr = kbuf;
while (count-- > 0)
{
    *(ptr++) = inb(port);
    rmb();
}
if ((retval > 0) && copy_to_user(buf, kbuf, retval))
    retval = -EFAULT;
```

74 / 76

Write

```
if (copy_from_user(kbuf, buf, count))
    return -EFAULT;
ptr = kbuf;
while (count-- > 0)
{
    outb(*(ptr++), port);
    wmb();
}
```

75 / 76

Reading Material

- ▶ Corbet-Rubini-Hartman, 3/e
 - ▶ Chapter 3: Char Drivers
 - ▶ Chapter 9: Communicating with Hardware

76 / 76