

Chapter 9: Filesystems

Introduction

Till this point, we have seen how a kernel of an operating system works as a resource manager. Given that both the processor and the main memory are resources in the system, 539kernel manages these resources ¹ and provides them to the different processes in the system. Another role of a kernel is to provide a way to communicate with external devices, such as the keyboard and hard disk. *Device drivers* are the way of realizing this role of the kernel ². The details of the external devices and how to communicate with them are low-level and may be changed at any time. The goal of a device driver is to communicate with a given device by using the device's own language and the other goal of a device driver is to provide an interface for any other component of the system that wish to use the given device, most probably the low-level details of this given device will be hidden behind the interface that the device driver provides, that means the user of the device drivers doesn't need to know anything about how the device really work.

The matter of hiding the low-level details with something higher-level is too important and can be found in, basically, everywhere in computing and the kernels are not an exception of that. Of course, there is virtually no limit of providing higher-level concepts based a previous lower-level concept, also upon something that we consider as a high-level concept we can build something even higher-level. Beside the previous example of device drivers, one of obvious example where the kernels fulfill the role of hiding the low-level details and providing something higher-level, in other words, providing an *abstraction*, is a filesystem which provides the well-known abstraction, a file.

In this chapter we are going to cover these two topics, device drivers and filesystem by using 539kernel. As you may recall, it turned out that accessing to the hard disk is an important aspect for virtual memory, so, to be able to implement virtual memory, the kernel itself needs to access the hard disk which makes it an important component in the kernel, so, we are going to implement a device driver that communicate with the hard disk in this chapter. After getting the ability of reading from the hard disk or writing to it, we can explore the idea of providing abstractions by the kernel through writing a filesystem that uses the hard disk device driver and provides a higher-level view of the hard disk, that we all familiar with, instead of the physical view of the hard disk which has been described previously in chapter . The final result of this chapter is version NE of 539kernel which has as we mentioned a hard disk device driver and a filesystem.

¹Incompletely of course, to keep 539kernel as simple as possible, only the basic parts of resources management were presented.

²At least a monolithic kernel.

ATA Device Driver

No need to say the hard disks are too common devices that are used as secondary storage devices. Also, there are a lot of manufacturers that manufacture hard disks and sell them. Imagine for a moment that each hard disk from a different manufacturer use its own way for the communication, that is, the method `X` should be used to be able to communicate with hard disks from manufacturer `A` while the method `Y` should be used with hard disks from manufacturer `B`. Given that there are too many manufacturers, this will be a nightmare. Each hard disk will need its own device driver which talks a different language from the other hard disk device drivers. Fortunately, this is not the case, at least for the hard disks, in this type of situations, standards are here to the rescue. A manufacturer may design the hard disk hardware in anyway, but when it comes to the part of the communication between the hard disk and the outside world, a standard can be used, so, any device driver that works with this given standard, will be able to communicate with this new hard disk. There are many well-known standards that are related to the hard disks, *small computer system interface* (SCSI) is one of them, another one is *advanced technology attachment* (ATA)³. While the latter one is more common in personal computers we are going to focus on it here and write a device driver for it, the former one is more common in servers.

As in PIC which is been discussed in chapter , ATA hard disks can be communicated with by using port-mapped I/O communication through the instructions `in` and `out` that we have covered previously. But before discussing the ATA commands that let us to issue a read or write request to the hard disk, let's write two routines in assembly that can be used in C code and perform the same functionality for the instructions `in` and `out`. If you didn't recall, the instruction `out` is used to write some bytes to a given port number, so, if we know the port number that a device receives that commands from, we can use the instruction `out` to write a valid command to that port, on the other hand, the instruction `in` reads data from a given port, for example, sometimes after we send a command to a device, it responds through writing something on a specific port, the instruction `in` can be used to read this value. The assembly code of the both routines that we are going to define next should reside in `starter.asm` anywhere between `bits 32` and the beginning of `start_kernel` routine. As we have said previously, the goal of these two routines is to make it possible for C code to use the instructions `in` and `out` through calling these routines. The following is the code of `dev_write` which can be used by C kernel code to write to a given port. In C, we can see that it has this prototype: `dev_write(int port, int cmd)`.

```
dev_write:
    ; Part 1
```

³Another well-known name for ATA is *Integrated Drive Electronics* (IDE). The older ATA standard is now known as Parallel ATA (PATA) while the newer version of ATA is known is Serial ATA (SATA).

```

push edx
push eax

; Part 2
xor edx, edx
xor eax, eax

; Part 3
mov dx, [esp + 12]
mov al, [esp + 16]

; Part 4
out dx, al

; Part 5
pop eax
pop edx

ret

```

The core part of this routine is part four which contains the instruction `out` that send the value which is stored in `al` to the port number which is stored in `dx`. Because we are using these two registers ⁴, we push their previous values into the stack as we did in the first part of the routine, pushing the previous values of these registers lets us restore them easily after the routine finishes its work, this restoration is performed in the fifth part of the routine right before returning from the routine, this is an important step to make sure that when the routine returns, the environment of the caller is same as the one before calling the routine. After storing the previous values of `eax` and `edx` we can use them freely, so, the first step after that is to clear their previous values by setting the value 0 to the both of them, as you can see, we have used `xor` and the both operands of it are the same register that we wish to clear, this is a well-known way in assembly programming to clear the value of a register ⁵. After that, we can move the values that have been passed to the routine as parameters to the correct registers to be used with `out` instruction, this is performed in the third part of the routine ⁶. The following is the code of the routine `dev_read` which uses the instruction `in` to read the data from a given port and return them to the caller, its prototype can be imagined as `char dev_read(int port)`.

```

dev_read:
    push edx

```

⁴Which are as you know parts of the registers `eax` and `edx` respectively.

⁵To my best knowledge its more performant than the normal way of using `mov`.

⁶The readers who have previous knowledge in x86 assembly programming may notice that I've omitted the epilogue of routines which creates a new stack frame, this decision has been made to make the matters simpler and you are absolutely free to use the calling convention.

```

xor edx, edx
xor eax, eax

mov dx, [esp + 8]

in ax, dx

pop edx

ret

```

For the same reason of restoring the previous environment when returning to the caller, the routine pushes the value of `edx` into the stack, then both of `edx` and `eax` are cleared since they will be used by the instruction `in`. After that, the value of the passed parameter which represents the port number that caller wishes to read from, is stored in `dx`. Finally, `in` is called, the result is stored in `ax`⁷, the previous value of `edx` is restored and the routine returns. You may ask, why did we only stored and restored the previous value of `edx` while the register `eax` is also used also, why didn't we store and restore the previous value of `eax`? The reason is that `dev_read` is a function that returns a value, and according to `cdecl` convention the returned values should be stored in the register in `eax`, so, the value of `eax` is intended to be changed when return to the caller, therefore, it will not be correct, logically, to restore the the previous value of `eax` when `dev_read` returns. The ultimate goal of defining both `dev_write` and `dev_read` is to make them available to be used in C code, so, the lines `global dev_write` and `global dev_read` should be written in the beginning of `starter.asm`.

The Driver

One ATA bus in the computer's motherboard makes it possible to attach two hard disks into the system, one of them is called master drive which is the main one that the computer boots from, the other disk is known as slave drive. Usually, a computer comes with two ATA buses instead of just one, which means up two four hard disks can be attached into the computer. The first one of those buses is known as the primary bus while the second one is known as the secondary bus. The port numbers that are used to communicate with the devices that are attached into the primary bus start from `0x1F0` and ends in `0x1F7` each one of them has its own functionality while the ports number from `0x170` to `0x177` are used to communicate with devices that are attached into the secondary bus, so, there are eight ports for each ATA bus. The terms that combine a bus name and a device name are used to specify exactly which device is being discussed, for example, primary master means the master hard disk that is connected to the

⁷Since the first operand of `in` is `ax` and not `al` then a **word** will be read from the port and not a single byte. The decision on using `ax` instead of `al` was made here because of our needs as you will see later, if you need to read just one byte for some reason you can define another routine for that.

primary bus while secondary slave means the slave hard disk which is connected to the secondary bus. For the sake of simplicity, our device driver is going to assume that there is only a primary master and all read and write requests should be oriented to this primary master, therefore, our device driver uses the port number 0x1F0 as the base port to send the commands via PIC. You may ask, why are we calling this port number a base port? As you know that all the following port numbers are valid to communicate with the primary ATA bus: 0x1F0, 0x1F1, 0x1F2, 0x1F3, 0x1F4, 0x1F5, 0x1F6, 0x1F7, we can add any numbers from 0 through 7 to the base port number of the primary bus 0x1F0 to get a correct port number to communicate with the primary bus, the same is correct with the secondary ATA bus which its base port number is 0x170. So, we can define the base port as a macro⁸ as we will see in our device driver, and then we can use this macro by adding a specific value to it from 0 through 7 to get a specific port, the advantage of doing so is the easiness of changing the value of the base port to another port without the need of changing the code itself. Before starting in the implementation of the driver, let's create new two files: `ata.h` and `ata.c` which will contain the code of the device driver which communicates to ATA devices and provides an interface for the rest of the kernel to write and read data from the disk. The following is the content of `ata.h` and the details of the functions will be discussed in the next subsections.

```
#define BASE_PORT 0x1F0
#define SECTOR_SIZE 512

void *read_disk( int );
void write_disk( int, void * );

void *read_disk_chs( int );
void write_disk_chs( int, void * );
```

Addressing Mode

As in the main memory, the hard disks use addresses to read the data that are stored in a specific area of the disk, the same is applicable in write operation, the same address can be used to write on the same specific area. There are two schemes of hard disk addresses, the older one is known as *cylinder-head-sector* addressing (CHS) while the newer one which more dominant now is known as *logical block addressing* (LBA). In chapter we have covered the physical structure of hard disks and we know from that discussion that the data are stored in small blocks known as sectors, also, there are tracks which each one of them consists of a number of sectors, and finally, there are heads that should be positioned on a specific sector to read from it or to write to it. The scheme CHS uses the same concepts of physical structure of hard disk, the address of given sector on the hard disk can be composed by combining three numbers together, the cylinder

⁸Or even variable.

(track) that this sector reside on, the sector that we would like to access and the head that is able to access this sector, however, this scheme is obsolete now and LBA is used instead of it. In LBA, a logical view of a hard disk is used instead of the physical view. This logical view states that the hard disk is composed of a number of logical blocks with a fixed size, say, n bytes. These blocks are contagious in a similar way of the main memory, to reach any block you can use its own address and the addresses start from 0, the block right after the first one has the address 1 and so on. As you can see, addressing in LBA is more like the addressing of the main memory, the main difference here is that in current computers each address of the main memory points to a byte in memory while an address in LBA points to a block which can be a sector (512 bytes) or even bigger.

Reading from Disk

In this subsection we are going to implements both `read_disk_chs` and `read_disk` which send commands to an ATA hard disk via the available ports in order to read a sector/block from the disk. The first one of those functions uses CHS scheme while the second one uses LBA scheme. In the next discussions, I'm going to use the symbol `base_port` to indicate the base port of one of ATA ports, in our case, as we mentioned earlier, the base port is `0x1F0` since we are going to use the primary bus in our device driver, but what we are discussing is applicable to any ATA bus with any base port number.

To issue a read command the command `0x20` should be sent to `base_port + 7`, but before doing that, a number of values should be set in the other ports in order to specify that address that we would like to read from, therefore, these values should be set before issuing the read command. In `base_port + 2` the number of sectors/blocks that we would like to read in this operation should be set. The value that should be written to `base_port + 6` specifies more than one thing, bit 6 of this value specifies whether we are using CHS or LBA in the current read request, when the bit's value is 0 then CHS is being used while the value 1 means that LBA is being used. The bits 5 and 7 of this value should always be 1. The bit 4 is used to specify the drive that we would like to read from, the value 0 for the master drive while 1 for the slave drive. In the case that we are using CHS, then the first four bits (0 to 3) of this value is used to specify the head and in LBA the bits 24 to 27 of the address should be stored in these bits. When the current addressing mode is CHS, the sector number that we would like our read operation to start from should be sent to `base_port + 3`, the low part of the cylinder number should be sent to `base_port + 4` and the high part of the cylinder number should be sent to `base_port + 5`. Once the read command is issued with the right parameters passed to the correct ports, we can read the value of `base_port + 7` to check if the disk finished the reading operating or not by reading the eighth bit (bit 7) of the value, when the value of this bit is 1 that means the drive is busy, once it becomes 0 that means that the operation completed. When the reading operation is completed

successfully, the data are brought to `base_port` which means we need to read from it and put the required data in the main memory. The following is the code of `read_disk_chs` that should reside in `ata.c`⁹.

```
void *read_disk_chs( int sector )
{
    // Part 1

    dev_write( BASE_PORT + 6, 0x0a0 );
    dev_write( BASE_PORT + 2, 1 );
    dev_write( BASE_PORT + 3, sector );
    dev_write( BASE_PORT + 4, 0 );
    dev_write( BASE_PORT + 5, 0 );
    dev_write( BASE_PORT + 7, 0x20 );

    // ... //

    // Part 2

    int status = 0;

    do
    {
        status = dev_read( BASE_PORT + 7 );
    } while ( ( status ^ 0x80 ) == 128 );

    // ... //

    // Part 3

    short *buffer = kalloc( SECTOR_SIZE );

    for ( int currByte = 0; currByte < ( SECTOR_SIZE / 2 ); currByte++ )
        buffer[ currByte ] = dev_read( BASE_PORT );

    return buffer;
}
```

In the first part of `read_disk_chs` we send the required values to the appropriate ports as we have described above. In the port `base_port + 6` we set that the drive is 0 and the head is 0, also, we set that the addressing mode that we are using is CHS. In port `base + 2` we set that we would like to read one sector, that is, the function `read_disk_chs` reads 512 bytes from disk with each call. In port `base_port + 3` we set the number of the sector that we would like to read, as you can see, this number is passed through the parameter `sector`, so,

⁹Don't forget to include `ata.h` when you create `ata.c`

the caller can specify the sector that it would like to read. In both `base_port + 4` and `base_port + 5` we specify the cylinder that we would like to read from, which is cylinder 0. Finally, we issue a read request to the ATA bus by writing `0x20` to `base_port + 7`. As you can see, for the sake of simplicity, we have used fixed values for a number of parameters here, in real situations, those fixed values should be more flexible. However, using LBA will provide us more flexibility with the simplicity that we are interested on. After issuing a read request in `read_disk_chs`, the second part of the function starts to work, it reads the value of the port `base_port + 7` which is going to contain the status of the drive, as we mentioned earlier, the eighth bit (bit 7) of this value indicates whether the drive is busy or not, with each iteration of the loop in part 2 the function reads this value and checks whether the value of the eighth bit is 1 by using bitwise operations, if this is the case that means that the drive is still busy in completing the reading operation, so, we need to wait for it until it finishes and we spend this time of waiting in reading the value of the port and check the eighth bits over and over again until the drive finishes, this technique of checking the device's status over and over again until it finishes an operation is known as *busy waiting*. When the read operation finishes, we can read the target content word by word^{10 11} from the base port and that's the job of the third part of `read_disk_chs` which reads a word each time and stores it in a buffer that we dynamically allocate from the kernel's heap. We have used the type `short` for the variable `buffer` because the size of this type in 32-bit architecture is 2 bytes, that is, a word. In the condition of the `for` loop we used `currByte < (SECTOR_SIZE / 2)` instead of `currByte < SECTOR_SIZE` due to the fact that we are reading two bytes in each iteration instead of one byte. Finally, the memory address which `buffer` points to, is going to contain the sector that we have just read from the disk, this memory address will be returned to the caller.

Now, we can turn to the other read function `read_disk` which uses LBA scheme instead of CHS. In fact, this new function will be almost similar to `read_disk_chs`, the main differences will be with some values that we are passing to the ports of ATA bus. In LBA scheme, `base_port + 6` should be changed to indicate that the scheme that we are using is LBA and we can do that as we mentioned before by setting the value 1 to bit 6. The other difference in this port value is that the bits 0 to 3 should contain the bits 24 to 27 of the logical block address that we would like to read from, the other parts of the addresses are divided to the following ports: `base_port + 3` contains bits 0 to 7 of that address, `base_port + 4` contains bits 8 to 15, `base_port + 5` contains bits 16 to 23. Both ports `base_port + 2` and `base_port + 7` stay the same. The function `read_disk` receives a parameter names `address` instead of `sector`, that is, the logical block address that the caller would like to read the data from should be passed to the function, by using bitwise operations the value of this parameter can be divided into the described parts to be filled in the appropriate

¹⁰Recall that a "word" in x86 architecture is 2 bytes.

¹¹The fact that we need to read a word here instead of a byte is the reason of using the register `ax` instead of `al` as the first operand for `in` instruction in the function `dev_read`.

ports. The rest of `read_disk` is exactly same as `read_disk_chs`. To not get a lot of space, the following is the beginning of the new function and the only part that has differences with `read_disk_chs`.

```
void *read_disk( int address )
{
    dev_write( BASE_PORT + 6, ( 0x0e0 | ( ( address & 0xF000000 ) >> 24 ) ) );
    dev_write( BASE_PORT + 2, 1 );
    dev_write( BASE_PORT + 3, address & 0x000000FF );
    dev_write( BASE_PORT + 4, ( address & 0x0000FF00 ) >> 8 );
    dev_write( BASE_PORT + 5, ( address & 0x00FF0000 ) >> 16 );
    dev_write( BASE_PORT + 7, 0x20 );
}
```