

OS Book

Neil Klingensmith

September 2, 2022

Contents

1	Introduction	5
1.1	What is an Operating System?	5
1.1.1	Managing Hardware with Abstraction	5
1.1.2	Operating System Abstractions	6
1.2	A History of x86 CPUs	8
1.3	The i386 Programmer's Model	9
1.3.1	16-Bit Real Mode	9
1.3.2	32-Bit Protected Mode	21
1.4	Exercises	21
2	x86 Boot Process	23
2.1	BIOS	23
2.1.1	BIOS Driver Interface	24
2.1.2	Printing Characters to the Screen	25
2.1.3	DOS	26
2.2	MBR	27
2.3	GRUB	27
2.3.1	Partitions	27
2.3.2	Binary File Formats	27
2.3.3	GRUB Hello World	27
2.3.4	Linker Scripts	30
2.4	Exercises	30
	Index	33

Chapter 1

Introduction

1.1 What is an Operating System?

As we will see throughout this book, computer hardware is complicated and messy. Most programmers would not like to interact directly with the hardware because of all its idiosyncracies. And since every CPU is different, porting software from one platform to another would be nightmarish without some kind of intermediate layer between the application and the hardware.

In fact, in all but the simplest computers, there are many layers of abstraction between an application and the hardware platform that it runs on. Details like how much memory the computer has, which applications are allowed to use which portions of memory, what type of display is available, how the computer is connected to the internet, and many others are managed by these intermediate layers of software.

1.1.1 Managing Hardware with Abstraction

The Android mobile ecosystem is a great example of abstraction between application and hardware. Many mobile devices run the Android operating system, with a huge amount of diversity in hardware capability among them. On the top end of the spectrum live high-performance tablets, mostly with price points in the range of \$1,000-\$2,000. These devices have big screens, fast WiFi and cellular networking interfaces, large memories, and fast processors. Many of them have detachable keyboards, styluses, and other Bluetooth accessories. Most high-end tablets also come with facial recognition cameras to authenticate their users without a password. On the low end of the spectrum are small, low-cost and lightweight phones, many of which cost less than \$100. They have smaller screens, slower CPUs and network interfaces, less memory, and fewer connectivity options.

Across all of these hardware platforms, even a seemingly simple task of displaying a menu on the screen is implemented completely differently. Once the layout of the menu screen is set by the application software, objects like text, images, and buttons are rendered into a rasterized image in memory. In the rasterization process, each pixel's color is set to an RGB value which is stored in an array in memory. The dimensions of the rasterized screen array depend on the size of the screen, which is different on each device. Once this rasterized array is constructed in memory, it needs to be transmitted from the main CPU to the display unit, which usually contains a separate coprocessor and memory to hold the contents of the rasterized screen buffer. The interface between the main CPU and the display coprocessor is also different on each device.

In the Android operating system, there are many layers of software involved in rendering

the display buffer and transmitting the rendered image to the screen. Some of those software modules are hardware-independent and common to all Android systems. These common modules live inside the operating system, and their purpose is mostly to make the application programmer’s job easier by automating some common, repetitive tasks like rasterizing text. Those common tasks need to be done by nearly every app that runs on the device, so it makes sense to have one common piece of code within the operating system to serve every app. Doing so reduces bugs (app developers can’t introduce bugs in software they don’t write) and creates a common interface so application software is easy to write, even by unskilled developers.

Some of the software modules involved in screen rendering are hardware-dependent, such as the communication between the mobile device’s main CPU and the display coprocessor. As we will discuss later in the book, we don’t trust application developers with any kind of direct access to the system’s hardware. Access to network interfaces, display controllers, power management, disks, and every other hardware resource is always mediated by the operating system. With direct access to the hardware, malicious or buggy code could do a lot of damage: erasing files from the disk, starving other applications of CPU time, attacking other computers on the network, shutting the machine down without warning, and any number of other kinds of problematic behavior.

1.1.2 Operating System Abstractions

So the operating system’s role is essentially to serve as a trusted mediator between the application and the hardware. In doing so, it provides a lot of higher-level abstractions that free up the application from having to worry about details of the hardware. Android apps don’t have to worry about how their screens are rendered, rasterized, and transmitted to the display—the OS does it for them. Apps don’t have to worry about how data is organized into files and stored onto the device’s SD card—the OS implements a file system that does it for them. And apps don’t need to worry about sharing CPU time. The OS does that for them as well. In general, we have a few desired functions for an operating system:

- Allow multiple programs to run simultaneously on the CPU.
- Allocate memory to programs.
- Provide support for files and filesystems.
- Mediate access to displays, input devices, network interfaces, and other I/O devices.
- Provide a user interface with a consistent look and feel.

The common thread among all functions of an OS is that it provides an abstract view of the computer hardware to the applications. People often say that an operating system’s job is to present a runtime environment to the computer’s applications that appears as if each application was the only one using the computer. Most operating systems provide some version of this vision of solitary execution to their applications, but none of the popular OSes in widespread use fully embrace it.

Why would an application wish to run in an abstracted environment like this? In the early days of computing, programs ran in *batch mode*, in which the computer would run one program at a time from start to finish. The program would have complete access to all of the computer’s resources, including its entire memory, its persistent storage¹, and all of its I/O devices. Programmers became accustomed to writing programs that had free rein of the computer, but there was a problem: it was not possible for two or more programs to coexist

¹Computers in those days used tapes instead of disk drives for long-term storage.

```

top - 19:57:42 up 3 days, 6:39, 2 users, load average: 0.10, 0.15, 0.14
Tasks: 323 total, 1 running, 320 sleeping, 2 stopped, 0 zombie
%Cpu(s): 0.7 us, 0.3 sy, 0.0 ni, 99.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0 st
MiB Mem : 3891.1 total, 335.3 free, 2181.4 used, 1374.4 buff/cache
MiB Swap: 0.0 total, 0.0 free, 0.0 used. 1357.2 avail Mem

  PID USER  PR  NI   VIRT   RES   SHR S  %CPU  %MEM    TIME+  COMMAND
46603 root   20   0     0     0     0 I   1.0   0.0    0:00.35 kworker/0:2+
46541 root   20   0     0     0     0 I   0.7   0.0    0:05.81 kworker/1:0+
46696 neil   20   0  10620   4008  3144 R   0.7   0.1    0:00.10 top
   553 root   rt    0 289324 27108  9076 S   0.3   0.7    0:58.28 multipathd
   640 root   20   0 241204  7024  5340 S   0.3   0.2   10:52.40 vmtoolsd
  1593 neil   20   0 312596  6356  5308 S   0.3   0.2    0:24.96 gvfs-afc-vo+
  1684 neil   20   0 4097088 191444 71664 S   0.3   4.8   75:33.33 gnome-shell
  2152 neil   20   0 4905840 551744 157796 S   0.3  13.8  115:13.77 firefox
27732 neil   20   0 2521696 178500  87048 S   0.3   4.5    3:51.30 Isolated We+
46455 root   20   0     0     0     0 I   0.3   0.0    0:18.64 kworker/0:0+
    1 root   20   0 167968  11588  6400 S   0.0   0.3    0:17.33 systemd
    2 root   20   0     0     0     0 S   0.0   0.0    0:00.18 kthreadd
    3 root    0 -20     0     0     0 I   0.0   0.0    0:00.00 rcu_gp
    4 root    0 -20     0     0     0 I   0.0   0.0    0:00.00 rcu_par_gp
    6 root    0 -20     0     0     0 I   0.0   0.0    0:00.00 kworker/0:0+
    9 root    0 -20     0     0     0 I   0.0   0.0    0:00.00 mm_percpu_wq
   10 root   20   0     0     0     0 S   0.0   0.0    0:00.00 rcu_tasks_r+
   11 root   20   0     0     0     0 S   0.0   0.0    0:00.00 rcu_tasks_t+
   12 root   20   0     0     0     0 S   0.0   0.0    0:30.62 ksoftirqd/0
   13 root   20   0     0     0     0 I   0.0   0.0    0:27.14 rcu_sched

```

Figure 1.1: The output of `top` showing many process running on the system.

on a computer at once. Operating systems were invented in large part to allow multiple programs to “time share” computers.

But by the time operating systems came into wider use, programmers were already accustomed to writing solitary programs that did not need to coexist with one another. So operating system programmers aimed to present a runtime environment that resembled the batch programming their users were comfortable with. Although this illusion of a completely isolated program has deteriorated over the years as programmers and operating systems have evolved, the legacy of the batch-mode programming model lives on in some respects. Isolated memory spaces—the concept that two programs sharing the CPU have no access to one another’s code or variables—is a clear vestige of batch-mode programming.

Figure 1.1.2 shows the output of the Linux `top` program, which displays the computer’s CPU and memory usage along with a sorted list of active processes. When this `top` snapshot was captured, there were 321 processes sharing the CPU, all of which were running under this illusion of isolation provided by the Linux operating system. Each process had its own memory space containing variables and code that only it could access. Each process also had the illusion that the CPU was under its complete control—there is no need for any of the processes to actively relenquish control of the CPU so other programs could run. Although not all 321 process could be actively executing simultaneously, the operating system provides the illusion of concurrency by allowing one to run for a while, then stopping that process and giving access to the CPU to a different process for a while. The Linux operating system is responsible for switching between actively running processes. From the user’s perspective,

all the processes appear to run simultaneously, even though only one or two of them can be active at a time².

In this book, we build a simple operating system that provides basic hardware virtualization: memory isolation, time sharing and file systems. We build this operating system on the Intel 386 PC platform. The PC platform is convenient because it is well-documented and widely available.

1.2 A History of x86 CPUs

Starting in the 1940s and continuing through about 1975, all computers were one-off designs. When a new computer was to be built, a team of 20 or so hardware engineers would spend three to five years laying out its specifications and designing its hardware. In these early days, logic components like AND gates and flip flops were built from discrete transistors wired together on printed circuit boards the size of a sheet of paper. A separate team of software engineers would write a custom operating system, compiler, and other tools specifically for the new design.

All this changed in the mid 1970s with the invention of the integrated circuit, which could pack multiple transistors on a single piece of Silicon. It was possible to fit an entire computer on a single chip and mass produce them at a fraction of the cost of mainframes and supercomputers. Integration as it was called was the key technological advancement that made personal computers economical.

In 1968, a small group of engineers left Fairchild Semiconductor to found Intel. At first they produced memory chips, but in 1971 they released the first fully integrated microprocessor, the 4004. While fairly useless by today's standards, the 4004 enabled miniturized electronics like calculators. It was soon succeeded by the 8008 and 8080 the first 8-bit microprocessors.

In 1976, Intel released the 8086, the product that would lead them to market dominance in PCs and servers. The 8086 was a fairly versatile 16-bit computer that made sense as the central component for PCs and minicomputers. However, as we will see later, it has a nonintuitive programmer's model and memory structure that is still being used in modern PCs today.

One important problem that designers of the 8086 had to overcome was memory addressability in a 16-bit computer. 16-bit addresses only would allow the programmer to access $2^{16} = 65536$ bytes of RAM, which was considered too little to be useful at the time. Programmers expected to have at least one megabyte of addressable memory, but that would require 20 address bits.

The solution they devised, which even at the time was widely seen as a kludge, was to use segmented addressing in which complete 20-bit addresses would be formed as the sum of a segment register plus a 16-bit pointer. Although crude, programmers learned to work with segmented addressing in x86 processors.

A marketing campaign by Intel targeted at manufacturers of PCs, minicomputers, and electronics manufacturers resulted in wide adoption of the 8086 and its successors. Despite its shortcomings, the 8086 was a success in the marketplace. Intel's price and first-mover advantages, not the technical merits of their product, made them successful.

A few years later, in 1979, Motorola released a competitive CPU called the 68000, so named because of the number of transistors on its die. Their design improved upon the 8086 in many ways. First, all registers in the 68000 were 32 bits long, making it possible to

²Each CPU core can actively execute one process. A computer with four CPU cores can actively execute four processes simultaneously, even though many more processes may be running on the system. The remaining processes are waiting for the operating system to give them access to one of the cores.

address significantly more memory without using segmentation. Second, the instruction set and programmer's model of the 68000 was much more intuitive than that of the 8086.

What happened next—part cost saving measure and part logistics solution—ultimately shaped the PC and server markets for decades. IBM, planning the new release of a home computer, chose the 8086 to be the central component of its new product. This decision was made because of the lower cost and wider availability of the 8086 and its peripheral components relative to the 68000. Apple, which was in process of designing the first Macintosh, chose to work with the 68000 for performance reasons.

When IBM shipped its first PCs, it also published a full set of schematics inside the user manual. Apple, a much more secretive organization, published very little information about the hardware and software inside the Macintosh. This had two effects. First, the availability of documentation about hardware and software on the PC made it much easier for third-party developers to write software for it. Even though the Macintosh had a much sleeker user interface, the lack of third party software support through the 80s and early 90s nearly lead to the company's demise. Second, low-cost clones of the IBM PC—made by Compaq and others—were widely available. As a result, the cost of ownership of the IBM PC was much lower and its utility was much higher. Riding on the success of the IBM PC, Intel's 8086 and successors became the most popular CPUs in desktop computers. Motorola eventually sold off and shut down its semiconductor division, including the 68000.

This book discusses operating system concepts and implementation techniques, with a focus on implementation on 8086-family CPUs. In particular, we work directly with hardware and software details of the i386 variant—the first to offer full support for most modern operating system features. Later variants, including the Pentium family and the AMD64 architecture offer architectural enhancements that give us performance gains—at the expense of quite a bit of additional complexity for the programmer. But the core concepts can all be implemented on the i386 CPU.

1.3 The i386 Programmer's Model

Before we begin discussing core operating systems concepts, it is important to have some base understanding of the low-level programmer's model of our CPU. Many of the operating system features we will introduce later in the book—virtual memory, time sharing, etc—will require us to manipulate the CPU at the assembly level. Although most of our operating system will be programmed in C, we will need to make use of a few assembly language instructions, and even many of the features of the OS implemented in C are directly manipulating some low-level features of the CPU.

We will start by discussing the i386's 16-bit real mode, which is a compatibility mode that allows the processor to run programs and operating systems written for the older 8086 CPU. Most of the concepts and features of the 8086 such as interrupts, registers, and stack frames should be familiar. The 8086 family CPUs do present a bit of a curve ball in the way they manage memory, which we will also discuss in the following sections.

1.3.1 16-Bit Real Mode

Figure 1.2 shows the programmer's view of the 8086 registers, all of which are 16 bits long. The 8086 has two distinct types of registers: integer registers **AX**, **BX**, **CX**, **DX**, and pointer registers **SI** and **DI**. The integer registers can be used for local storage of variables, and the pointer registers store addresses in memory. The 8086 instruction pointer **IP** stores the address of the instruction that is currently being executed.

All 8086 assembly instructions operate on two operands: one source and one destination. **The destination operand is always on the left** in 8086 assembly. For example, to set

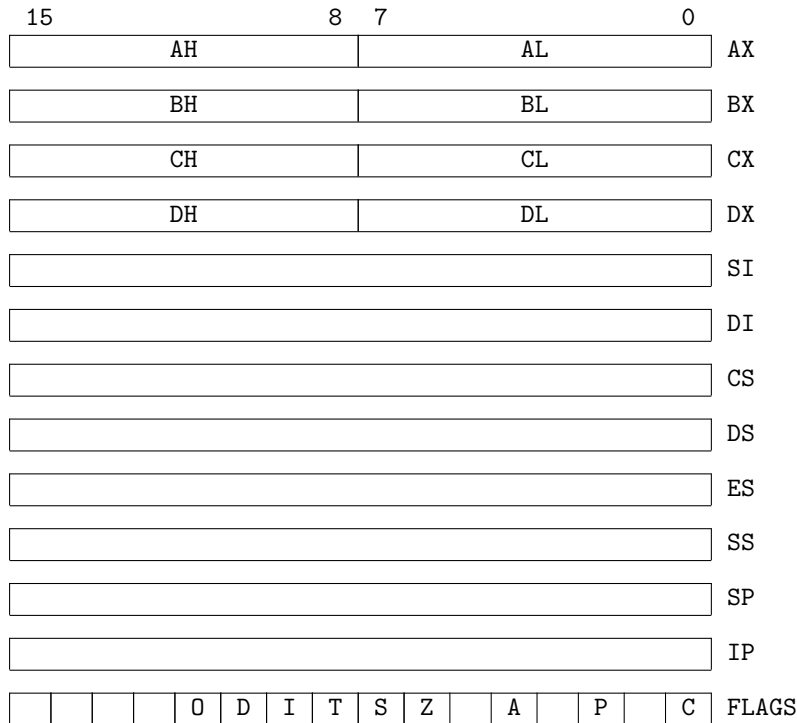


Figure 1.2: Registers of the 8086.

the value of register AX to 5:

```
mov ax,5
```

In this `mov` instruction, register AX is the destination, and 5 is the source. Putting the destination on the left is similar in notation to the assignment operator (=) in higher level languages, which would represent the same operation as follows:

```
ax = 5;
```

Example: Adding Two Numbers Together in 8086 Assembly The following example illustrates how to add two numbers together in 8086 assembly.

```

1  int i = 2; // in reg ax
2  int j = 3; // in reg bx
3  i += j;

1 [BITS 16]
2
3 boot: ; This file starts executing here.
4
5     mov ax,2 ; ax := 2
6     mov bx,3 ; bx := 3
7     add ax,bx ; ax := ax + bx
8     hlt
9
10    times 510-($-$$) db 0 ; Skip to end of boot sector
11    db 0x55 ; Magic Numbers
12    db 0xaa ; To make disk bootable

```

In the above example, we use integer registers `AX` and `BX` to hold the values of `i` and `j` respectively. The `ADD` instruction adds the two together, storing the result in destination register `AX`. Source register `BX` is not modified by the `add`.

Assembling and Running Real-Mode Programs This book assumes you are using a Linux system to build and test your code. Like all the examples in this chapter, you can enter the assembly listing into a text editor and compile it with `nasm` to produce a runnable 16-bit 8086 program:

```
user@system ~ $ nasm -f elf32 -F dwarf mbr.asm -o mbr.o
user@system ~ $ ld -Ttext=0x7c00 -melf_i386 mbr.o -o mbr.elf
user@system ~ $ objcopy -O binary mbr.elf mbr.img
```

The first command assembles the source file `mbr.asm` to object file `mbr.o`. The second command uses the `ld` linker to locate the object code to address `0x7C00`, creating output file `mbr.elf`. The third command converts the ELF file to a flat binary file `mbr.img`, which can be booted directly by `qemu`. Your program can be emulated in 16-bit real mode using `qemu`:

```
user@system ~ $ qemu-system-i386 -hda mbr.img
```

The `qemu` command emulates your `mbr.img` file on a 386 computer. Your `mbr.img` file is treated as the main boot hard disk on the emulated machine. Of course it is a very small hard disk, with only one 512-byte sector.

Loops in 8086 Assembly

As with most assembly languages, building a loop in 8086 consists of three steps: initialization, loop body, and conditional branch instructions. Let's start with a simple example that clears the contents of a 10-byte buffer with zeros.

Example: Filling an Array with Zeros

```
char buf[10];
for(int k = 0; k < 10; k++) {
    buf[k] = 0;
}
```

```
1 [BITS 16]
2   lea di,buf // di points to buf
3   mov bx,0 // bx is index reg
4   loop:
5   mov byte [di,bx],0 // buf[bx] := 0
6   add bx,1 // increment bx
7   cmp bx,10 // bx < 10 ?
8   jl loop // loop if bx < 10
9   hlt
10  buf:
11  db 10 dup 0xff // 10 bytes initialized to 0xff
12
13  times 510-($-$$) db 0 ; Skip to end of boot sector
14  db 0x55 ; Magic Numbers
15  db 0xaa ; To make disk bootable
```

Instruction	Condition
<code>ja</code>	Jump if Above (Unsigned Greater Than)
<code>jae</code>	Jump if Above or Equal (Unsigned Greater or Equal)
<code>jb</code>	Jump if Below (Unsigned Less Than)
<code>jbe</code>	Jump if Below or Equal (Unsigned Less or Equal)
<code>jl</code>	Jump if Less Than (signed)
<code>jle</code>	Jump if Less or Equal (signed)
<code>kg</code>	Jump if Greater Than (signed)
<code>jge</code>	Jump if Greater or Equal (signed)
<code>je</code>	Jump if Equal
<code>jne</code>	Jump if Not Equal

Table 1.1: Commonly used 8086 jump instructions.

In the above example, we begin with a 10-byte buffer initialized with `0xff`. Lines 2 and 3 are the initialization in which we point the `DI` register to the beginning of the buffer and initialize the index register `BX` to zero. Line 4 has a label that represents the beginning of the body of our loop. On each round of the loop, we will jump to the `loop` label and execute the body of the loop. On lines 7 and 8 we have a compare and conditional jump. The compare instruction on line 7 checks to see if we are finished with all iterations of the loop by comparing the value in register `BX` to 10. The conditional jump on line 8 jumps back to the beginning of the loop if `BX` is less than 10. The most commonly used 8086 jump instructions are shown in Table 1.1.

Example: if Blocks

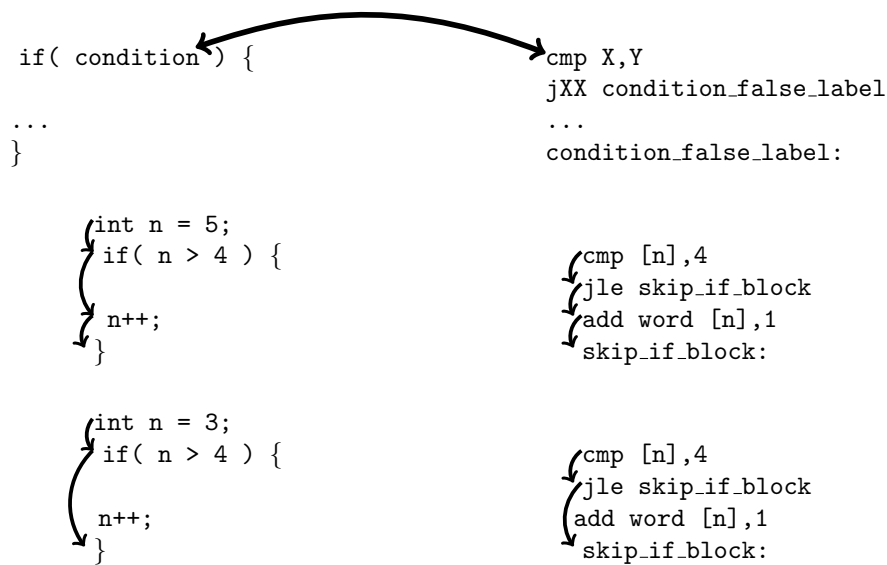
```
int n = 5;
if(n > 4) {
    n++;
}
```

```
1 [BITS 16]
2 boot:
3   cmp word [n],4 ; n > 4 ?
4   jle skip_if_block ; Skip following instruction if n <= 4
5   add word [n],1 ; n++
6 skip_if_block:
7   hlt
8 n:
9   dw 5
10 times 510-($-$$) db 0 ; Skip to end of boot sector
11 db 0x55 ; Magic Numbers
12 db 0xaa ; To make disk bootable
```

In this example, we are using a conditional jump instruction to implement an `if` block. The canonical form for implementing an `if` block in assembly is shown in the code snippet below. First, the `if` condition is evaluated with a `cmp` instruction. After the `cmp`, we use a conditional jump to skip past the body of the `if` block if the condition evaluates to false. After the `jmp`, we write the body of the `if` block that will be executed if the condition evaluates to true.

If the condition evaluates to true, the conditional jump will be taken, and the code inside the braces will be skipped. If the condition evaluates to false, the conditional jump will not be taken. Conditional jumps that are not taken have no effect on the program: instruction

continue to execute in order after the untaken conditional jump.



The 8086 Stack and Calling Conventions

8086 systems use the stack to pass parameters to functions and store local variables. The SP register points to the top of the stack, which grows downward in memory (toward lower memory addresses) as new items are pushed. The `push` and `pop` instructions can be used to push and pop data to the stack. The `call` and `ret` instructions are used for function calls and returns: `call` pushes the return address to the stack before jumping to a function, and `ret` pops the return address into the instruction pointer (IP).

Example: Calling Functions In 8086, we use the `call` instruction to call functions. `call` pushes the return address onto the stack and jumps to the specified function. When a function is finished executing, it uses the `ret` instruction to return. `ret` pops the return address off the stack into the IP. The example below assumes that `char *s` is passed by the caller in the SI register, and the result is returned in AX.

```

int strlen(char *s) {
    int k = 0;

    while(s[k] != '\0') {
        k++;
    }
    return k;
}

```

Passing the parameter to `strlen` works well in this example, but what about functions that take more than one or two parameters? The 8086 only has four integer registers and two pointer registers, leaving a total of at most six registers for us to use at one time. And some of those will be taken up by local variables and scratch storage. Clearly we need to find a different place for local temporary storage. Main memory is the obvious choice. Next, we discuss a data structure called a *stack frame* that can be used for organizing local storage in memory.

```

1  [BITS 16]
2  boot:
3    lea si,string ; Get addr of string in SI
4    call strlen ; Call strlen function
5    hlt ; Halt CPU
6  strlen:
7    mov bx,-1 ; Initialize bx
8  loop:
9    add bx,1 ; increment bx
10   cmp [si,bx],0 ; check for NULL term
11   jne loop ; loop
12   mov ax,bx ; put ret val in AX
13   ret
14  str:
15   db 'this is a string',0 ; 10 bytes initialized to 0xff
16
17  times 510-($-$$) db 0 ; Skip to end of boot sector
18  db 0x55 ; Magic Numbers
19  db 0xaa ; To make disk bootable

```

Listing 1: A simple implementation of `strlen` in 8086 assembly.

Stack Frames In addition to tracking a function’s return address, functions also allocate their local variables on the stack. Like most processors, the 8086 provides some primitive instructions for managing data on its stack: the `push` instruction adds a new element to the top of the stack and `pop` removes an element from the top of the stack. The CPU’s stack pointer (SP) holds the address in main memory of the last element that was added to the stack.

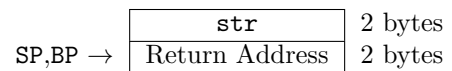
Let’s rewrite the `strlen` function, this time using the stack to pass parameters from `boot` to `strlen`. In `boot`, just before we call `strlen`, we will push the address of `string` onto the stack instead of placing it in SI. Then, in `strlen`, we will retrieve the address of the string from the stack, placing it into the SI register before we begin computing the string’s length. Other than how we communicate `strlen`’s parameter from the caller to the callee, everything else is the same. We still use the SI register to hold the address of the string inside `strlen`, and we still use BX to index into the string.

The function call in this example is simple: one line 3, we push the parameter we want to pass to `strlen` onto the stack. On line 4 we call `strlen`. Then, on line 5 after `strlen` returns, we remove the parameter that we pushed from the stack. The pointer argument that we pass to `strlen` is 2 bytes long, so we add 2 to the SP to remove it from the stack. If we don’t remove parameters from the stack after function calls, the stack will just keep growing, taking up more and more memory.

Inside `strlen`, we need to get the function’s parameter from the stack into SI before we begin computing the string’s length. But the 8086 CPU doesn’t allow us to use the SP register directly to read from the stack. We need to use a different register instead. On line 9, we copy the address from the SP into the BP register, which is used to read values from the stack.

Then, on line 10, we read the value of the parameter from the stack into SI.

Figure 1.3 shows `strlen`’s stack frame. Inside of `strlen`, the SP and BP both point to the same place on the stack because on Line 9 we copy the value of SP to the BP. That slot on the stack contains the return address to the caller, which was pushed by the `call` instruction. Just above the return address on the stack is the parameter we passed to `strlen`. That parameter lives at SP+2 because the return address is 2 bytes long.

Figure 1.3: Stack frame for `strlen` in Listing 2.

```

1 [BITS 16]
2 boot:
3   push str ; Push the address of string onto the stack
4   call strlen ; Call strlen function
5   add sp,2 ; Remove the address of string from the stack
6   hlt ; Halt CPU
7 strlen:
8   mov bx,-1 ; Initialize bx
9   mov bp,sp ; Copy SP to BP
10  mov si,[bp+2] ; Read the parameter from the stack
11 loop:
12  add bx,1 ; increment bx
13  cmp byte [si,bx],0 ; check for NULL term
14  jne loop ; loop
15  mov ax,bx ; put ret val in AX
16  ret
17 str:
18  db 'this is a string',0 ; 10 bytes initialized to 0xff
19  times 510-($-$$) db 0 ; Skip to end of boot sector
20  db 0x55 ; Magic Numbers
21  db 0xaa ; To make disk bootable

```

Listing 2: XXX

The code in Listing 2 has one problem that needs to be fixed: the `strlen` function overwrites the BP, SI and BX registers on lines 8, 9 and 10 without saving their values. If there were important values in BP, SI or BX, they will be overwritten. We need to save the value of the both registers before we clobber them. That is simple enough to do by just pushing the BP, SI and BX onto the stack before we clobber their values. Then we can use all three registers in `strlen` to refer to values on the stack and index into our string.

Immediately before `strlen` returns, we will pop BP, SI and BX off the stack, restoring their original values. On a side note, we also clobber whatever value was in AX, but since that register is used to store `strlen`'s return value, there is no need to save its old value on the stack before we clobber it. The final implementation of the `strlen` function is shown in Listing 7, with an updated stack frame shown in Figure 1.4.

The process of saving copies of the registers before we overwrite them in a function is called the function's *prologue*. `strlen`'s prologue is on lines 8-1 of Listing 7. Restoring the registers is called the function's *epilogue*, shown on lines 19-21 of Listing 7.

To make our lives easier, the 8086 implements the `pusha` and `popa` instructions that respectively push all registers and pop all registers. Using these instructions saves both typing and thinking. It saves typing by reducing the number of `push` and `pop` instructions we need to write in the prologue and epilogue. It saves thinking by saving all registers on the stack, freeing us from accounting for which registers our function is using.

Example: Passing Parameters to Functions In the 8086, parameters are passed to a function on the stack. To call the `strlen` function in the previous example, we need to pass a pointer to a string on the stack.

Local Variables

The stack allows us to stash and restore the contents of all the registers, freeing the registers up for local storage within a function without needing to worry about overwriting important

	str	2 bytes
	Return Address	2 bytes
BP →	Caller's BP	2 bytes
	Caller's SI	2 bytes
SP →	Caller's BX	2 bytes

Figure 1.4: Stack frame for `strlen` in Listing 7 with registers saved on the stack.

```

1  [BITS 16]
2  boot:
3    push str ; Push the address of string onto the stack
4    call strlen ; Call strlen function
5    add sp,2 ; Remove the address of string from the stack
6    hlt ; Halt CPU
7  strlen:
8    push bp ; Stash caller's BP on the stack
9    mov bp,sp ; Copy SP to BP
10   push si ; Stash caller's SI on the stack
11   push bx ; Stash caller's BX on the stack
12   mov bx,-1 ; Initialize bx
13   mov si,[bp+4] ; Read the parameter from the stack
14  loop:
15   add bx,1 ; increment bx
16   cmp byte [si,bx],0 ; check for NULL term
17   jne loop ; loop
18   mov ax,bx ; put ret val in AX
19   pop bx ; Restore BX, SI, and BP before return
20   pop si
21   pop bp
22   ret
23  str:
24   db 'this is a string',0 ; 10 bytes initialized to 0xff
25  times 510-($-$$) db 0 ; Skip to end of boot sector
26  db 0x55 ; Magic Numbers
27  db 0xaa ; To make disk bootable

```

Listing 3: XXX

values. But what about functions with lots of local variables? The 8086 only provides six registers—four integer registers and two pointer registers—enough for trivial functions like `strlen`, but lots of functions need more than six local variables.

There is also a way to use the stack for local variable storage, which allows us to have more local variables than registers. Local variables are customarily allocated on a function's stack frame. Consider the `max` function below, which finds the maximum value in an integer array:

```

1  int max(int *buf, unsigned int len) {
2    int maxval = -32767;
3    int k;
4    for(k = 0; k < len; k++) {
5      if(buf[k] > maxval) {
6        maxval = buf[k];
7      }
8    }
9    return maxval;
10 }

```

Listing 4: XXX

This function has two local variables: `maxval` and `k` along with two parameters passed by the caller (`buf` and `len`). Although this function is not at risk of exhausting the 8086's registers, we can use it to demonstrate how to allocate local variables on the stack. We will allocate the `maxval` variable on the stack, and we will keep the array index `k` in register `BX` as before. The stack frame for our implementation of `max()` is drawn in Figure 1.5. We can access the value of the `maxval` variable on the stack at location `BP-2`. The prologue we will

use to create this stack frame is as follows:

1. `push bp`
2. `mov bp,sp`
3. `sub sp,2`
4. `push bx`

Instruction 1 saves the value of the CPU BP register on the stack so we can overwrite it. Instruction 2 copies the SP register to the BP. The BP register is a pointer from which all accesses to the stack frame are made—if we want to read or write the value of a local variable or a parameter passed to the function, we use the BP register. After instruction 2, the SP and BP registers both point to the same location on the stack. Instruction 3 subtracts 2 bytes from the SP, allocating space for the 2-byte integer `maxval`. Instruction 4 pushes the value of the BX register, which allows us to overwrite it in the function without losing its contents.

Instructions 3 and 4 cause the SP to change, but the BP continues to point to the same slot on the stack for the remainder of the function. This is a convenient feature of the BP: it never moves during the course of a function's execution. The SP, on the other hand, does move during a function's execution. For example, if we call a function, we will push the function's arguments onto the stack, causing the SP register to change. If we tried to refer to variables relative to the SP, their locations relative to the SP would change every time we push or pop something to the stack. Since we reference all of our variables relative to the BP, the variable's locations relative to the BP never change.

In summary, the procedure for function calls is as follows:

1. The caller pushes the function's arguments onto the stack in reverse order.
2. The caller uses a `call` instruction to jump to the function. The `call` instruction pushes the return address onto the stack before jumping.
3. The function being called (the *callee*) pushes the BP register and copies the SP to the BP. BP now points to the middle of the stack frame, just below the return address.
4. The callee creates space for its own local variables by subtracting from the SP the number of bytes needed to store its locals. For example, if callee has two local `int` variables, it would subtract $2 \times 2 = 4$ bytes from the SP.

Interrupts

There are many hardware events that occur in a system that need to be dealt with immediately. Urgent events are usually triggered by I/O devices like the keyboard, mouse, network interface, etc., and if they're not dealt with immediately, the computer will appear laggy and unresponsive to the user. Most CPUs offer a mechanism called *interrupts* to temporarily transfer control from the application that is currently running to a function called an *interrupt service routine* that can process data from the I/O device that caused the interrupt. An interrupt service routine is a regular function in every sense except that it is called by hardware, not by software.

Most functions are called by software during the normal course of execution of a program. That is, a programmer plans to call a function. Interrupt service routines are by definition unplanned. Since we don't know when an interrupt event will occur, we do not know when during a program's execution the CPU will invoke an interrupt service routine.

```

1  [BITS 16]
2  boot:
3      push 10 ; Push the length of buf onto the stack
4      push buf ; Push the address of buf onto the stack
5      call max ; Call max function
6      add sp,4 ; Remove the address of buf and length of buf from the stack
7      hlt ; Halt CPU
8  max:
9      push bp ; Prologue
10     mov bp,sp
11     sub sp,2
12     pusha
13     mov word [bp-2],-32768 ; Initialize maxval
14     mov bx,0 ; Initialize k = 0
15     mov si,[bp+4] ; Get pointer to buf in SI
16     shl word [bp+6],1 ; Multiply length of buf by 2 since each element of buf is 2 bytes long
17  max_loop:
18     mov ax,[si+bx]
19     cmp ax,[bp-2] ; Check if maxval < buf[k]
20     jl not_greater ; If maxval > buf[k], do not update maxval
21     mov [bp-2],ax ; maxval := buf[k]
22  not_greater:
23     add bx,2 ; Add 2 to index register because each element of buf is 2 bytes long
24     cmp bx,[bp+6] ; Check if we have reached end of buf
25     jl max_loop ; If not, go back to beginning of loop
26     popa ; Epilogue
27     mov sp,bp
28     pop bp
29     ret
30  buf: dw 167,99,10000,-2598,31000,32000,-31000,-6000,2000,0 ; 10 an array with 10 integers
31  times 510-($-$$) db 0 ; Skip to end of boot sector
32  db 0x55 0xaa ; Magic Numbers to make disk bootable

```

Listing 5: XXX

What if the application is executing some critical piece of code when the CPU takes an interrupt? We do not want the application's state—the register contents, the stack contents, etc—to be disrupted by the interrupt service routine. So when an interrupt occurs, the interrupt service routine must not modify the CPU state.

To avoid modifying the CPU state, the interrupt service routine pushes the contents of all registers onto the stack before responding to the hardware event. Preserving the CPU state is actually done in two phases: first, the CPU automatically pushes the return address and the `FLAGS` register onto the stack before calling the ISR. Then, in software, the ISR pushes the contents of all other registers onto the stack.

Since the 8086 has this weird segmented addressing where each 16-bit register can only represent 64k of the 1M address space, we need to save both the `CS` and `IP` registers to record the return address. Together, they take up a total of four bytes on the interrupt stack frame. Figure 1.6 shows a diagram of the 8086 interrupt stack frame.

	unsigned int len	BP+6
	int *buf	BP+4
	Return Address	BP+2
BP →	Caller's BP	BP+0
	int maxval	BP-2
SP →	Caller's BX	BP-4

Figure 1.5: Stack frame for the `max` function

Memory Segmentation

One problem the 8086 had was its ability to support large memories. Its 16-bit registers could only address 64kbytes of memory ($2^{16} = 65536$), which isn't really enough to hold large programs. The 8086 architects wanted to support 1 mbyte of memory, but for that they would need 20 address bits—4 more than they had.

To solve this problem, they created segment registers—extra registers that were appended to the 16-bit address registers to hold the high-order 4 bits of a 20-bit address. Table 1.2 illustrates how segment registers are combined with address registers to form a physical address. The value in the segment register is shifted left by 4 bits and added to the address register. The sum becomes the 20-bit physical address that is presented to the memory system.

Segment registers are named `CS`, `DS`, `ES`, `SS` (see Table ??). The code segment register `CS` is appended to the instruction pointer `IP` when fetching instructions from memory. The data segment register `DS` is appended to `SI` and `DI` by default when fetching data with source index and destination index (`SI` and `DI`) registers. The stack segment register `SS` is appended to the stack pointer `SP` when reading and writing from the stack. The extra segment register `ES` is not used by default, but it can be used as an additional segment register.

	FLAGS
	CS
SP →	IP

Figure 1.6: Interrupt stack frame created by the 8086 CPU.

The main problem with memory segmentation is address aliasing—many combinations of `SEGMENT:OFFSET` can refer to the same physical address. For example, segment `0x07C0` offset `0x0800` refers to the same address as `0x07E0` offset `0x0000` (both generate physical address `0x7E00`). Programmers of 16-bit 8086 systems had to be very careful when generating pointers!

Memory segmentation makes programming for the 8086 processor uniquely weird, and it is widely regarded by programmers as a blunder. But from a practical perspective, memory segmentation was a success because it allowed the 8086 to support large memories with a relatively simple 16-bit processor. The result was a cheap and versatile CPU (although difficult to program).

A	0	0	0	0	DS					
+					0	1	0	0	0	SI
A	0	1	0	0	0	Effective Address				

Table 1.2: Memory segmentation in the 8086.

Example: 8086 Hardcoded Terminal Output One of the cool things about the x86 PC platform is it allows us to write directly to the display with very little setup. If you're lazy, you can use the BIOS's built-in terminal driver (which we talk about later in the book). Or you can write characters directly to the screen at any location—no `printf` call needed. Writing to the screen is done using *video memory*.

Video memory is basically just a big array in the PC's memory map that is dedicated to video output. Any characters you write to that array is displayed on the screen. By default, the screen is divided into 80 columns \times 24 rows of characters. Each character is represented in video memory as 2 bytes: one byte of ASCII character and one byte of color.

	Column 1	Column 2	Column 3	...	Column 80		
	↓	↓	↓		↓		
B8000	ascii	color	ascii	color	...	ascii	color
B8050	ascii	color	ascii	color	...	ascii	color
B80A0	ascii	color	ascii	color	...	ascii	color

The character at position (1,1) is located in the upper left corner of the screen, and row numbers increase toward the bottom of the screen. The character at position (1,1) lives at address `0xB8000`. The next character on the first row—at position (1,2) lives at address `0xB8002` and so on.

Let's write a program to draw the character `a` in the upper-left corner of the screen. To do that, we'll be writing two bytes starting at address `0xB8000`. The first byte to write will be the ASCII character (`a`), and the second byte will be the color to display. The color code for gray text on black background is `0x07`.

The only problem is that the address of video memory (`0xB8000`) is 20 bits long: bigger than the 16-bit registers in the 8086. We'll need to use segment registers to get the right address. In this code, we are going to use the Extra Segment (`ES`) register to store the beginning of video memory. We will put the value `0xB800` into register `ES` and use `SI` as the offset into video memory (depicted in Table 1.3). When we write to `ES:[SI]`, the store will be done to address `0xB8000+SI`. The assembly code listing below shows how to do this in 8086 assembly.

Even a simple operation like writing a single character to the screen is extremely clunky in 8086 assembly. Memory segmentation makes programming very painful. If only the CPU's

B	8	0	0	0	ES					
+					0	0	0	0	0	SI
B	8	0	0	0	0	Effective Address				

Table 1.3: Using the `ES` register to address video memory.

```

1 [BITS 16]
2 boot:
3     mov ax,0xb800 ; Load ES with base of video memory
4     mov es,ax
5     mov al,'a' ; ASCII code for the character to display
6     mov ah,0x07 ; Color code 7 = gray text on black background
7     mov word es:[si],ax
8     hlt
9 times 510-($-$$) db 0 ; Skip to end of boot sector
10 db 0x55 0xaa ; Magic Numbers to make disk bootable

```

Listing 6: Writing an a to the top left of the screen.

registers were long enough to support pointers to any location in memory, we wouldn't need memory segmentation. Next we discuss the i386's protected mode, which does away with memory segmentation and makes programming much easier.

1.3.2 32-Bit Protected Mode

The x86 family's 32-bit protected mode is an extension of 16-bit real mode that provides wider registers and support for virtual memory. The programmer's model of the i386 CPU is shown in Figure 1.7. All of the 8086 registers are defined in the same way in the i386: **AX** is a 16-bit integer register with upper and lower halves **AH** and **AL**. The four integer registers have extended versions **EAX**, **EBX**, **ECX**, **EDX** in the i386, each 32 bits long. Assembly instruction opcodes (**mov**, **add**, **push**, etc) in protected mode are all the same as real (16-bit) mode. Since C compilers for 32-bit protected mode are widely available, we won't discuss the assembly language in much detail. But in order to write an operating system for the i386, you do need to understand a bit of the CPU's architecture.

Global Descriptor Tables

Memory segmentation lives on in a weird way on the i386's protected mode.

1.4 Exercises

Exercise 1.1 Write a program that clears the screen by writing the space character to every position on the terminal.

Exercise 1.2 Write a function that prints one character to an arbitrary (x, y) location on the screen. Your function should take three inputs: (1) the character to print, (2) the x location, and (3) the y location. *Hint: you may need to use the 8086 **mul** instruction to compute the offset into video memory where the character should be written.*

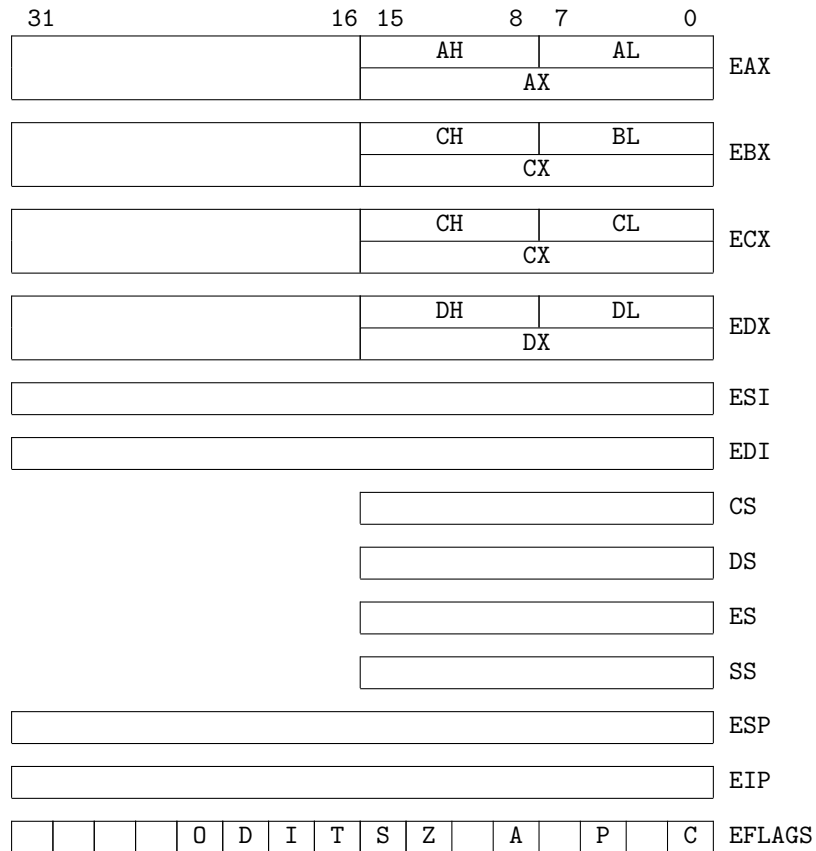


Figure 1.7: Registers of the i386.

Chapter 2

x86 Boot Process

The process by which a computer loads its operating system after it powers up is called booting—a reference to “pulling itself up by its bootstraps.” After the computer powers on, its memory and peripherals are uninitialized. We need some software in place to load the operating system from the hard disk into memory.

The IBM PC has been produced in many variants: different CPUs, different disk controllers, different display adapters, all made by different manufacturers. We cannot expect every program that runs on the PC platform to have internal driver support for every possible hardware device—particularly not the simple low-level bootloader. Instead, PC manufacturers provide drivers for their hardware burned into ROM that all present a common interface to the user’s software. These drivers, called the BIOS, allow programmers to access basic hardware features without integrating complex driver functionality into their programs.

The BIOS is also responsible for reading the bootloader from the first sector of the hard drive and passing control to it.

2.1 BIOS

The bootloader and other user code access BIOS functions using software interrupt instructions that are handled by the BIOS. Applications call the BIOS by placing commands into the CPU registers (*AX*, *BX*, etc) and then executing a *INT* instruction, which causes control to jump to BIOS code. The BIOS reads the commands in the CPU registers and executes the corresponding function.

The secret behind 8086 variant’s enduring success in the marketplace is compatibility. Even today, a brand-new Intel Core i7 system with all its modern features can run programs and operating systems that were written for IBM PCs manufactured in 1982. This might not seem impressive, but it is extremely valuable for businesses. Companies who invest in new PC-based hardware have never had to worry that hardware upgrades would make their software obsolete. They can always invest in new hardware and be certain that their existing software will continue to work. But how can Intel CPUs continue to support old software as the hardware architecture evolves over time?

On powerup, all x86 CPU variants—even modern 64-bit variants—begin executing code in 16-bit real mode starting at address `0xFFFF0`. 16-bit real mode is a compatibility mode that is supported in modern x86 CPUs to ensure that older software will continue to run on new machines. The instruction set architecture of real mode is the same as the original 8086 CPU. This is the instruction set that the BIOS and bootloader must use to load the operating system.

The BIOS is burned into ROM on the motherboard, and its first instruction lives at address `0xFFFF0` (see Table 2.1 for the BIOS's memory map). On powerup, the BIOS initializes the hardware and the DRAM controller. Then, using its internal drivers, it loads the first 512-byte sector from the boot disk into RAM at address `0x7C00`. This sector is called the master boot record, and it is the first piece of user-writable code that runs.

2.1.1 BIOS Driver Interface

The BIOS supports low-level drivers that can be called by programs to access hardware. The interface between applications and the BIOS happens through software interrupts—assembly language instructions that cause interrupts to occur. When the CPU encounters an `INT` instruction, it stops executing the application code and pushes the `IP` and `FLAGS` registers onto the stack. The CPU then reads the address of the corresponding interrupt handler from the interrupt vector table and jumps to that address. The interrupt vector table and corresponding interrupt handlers are installed at boot time by the BIOS, so the handler is part of the BIOS code.

Each software interrupt handler supported by the BIOS handles a different hardware driver. Generally, before the application executes a software interrupt, it will place some command code along with additional parameters in the CPU registers that tell the BIOS what to do. If the application needs to read data from a disk, for example, it will populate the CPU registers with information about where on the disk it wants to read and how much data needs to be read. When the BIOS finishes executing its interrupt handler, it uses an `RETI` instruction to pop the `FLAGS` register and return address of the stack and return to the application.

Below we discuss a few useful BIOS functions and give example code for calling them. For a more complete list of all BIOS functions, search Google for *Ralph Brown's Interrupt List*.

Getting a Keystroke from the Keyboard

Probably the simplest BIOS call is requesting a keystroke from the keyboard. This call just hangs until the user presses a key, then it returns with the ASCII code for the keypress in register `AL`. To call the Get Keystroke function in the BIOS, just put the value 0 in register `AH` then execute an `INT 0x16` instruction:

```
[BITS 16]

boot: ; This file starts executing here.

    mov ah,0
    int 0x16 ; This instruction will hang until you press a key.
            ; On key press, BIOS returns with AL = ASCII code
    hlt

times 510-($-$$) db 0 ; Skip to end of boot sector
```

Address Region	Contents
0x00000 - 0x003FF	Real Mode Interrupt Vector Table
0x00400 - 0x004FF	BIOS Data Area
0x00500 - 0x07BFF	Unused
0x07C00 - 0x07DFF	Bootloader
0x07E00 - 0x9FFFF	Unused
0xA0000 - 0xBFFFF	Video RAM (VRAM) Memory
0xB0000 - 0xB7777	Monochrome Video Memory
0xB8000 - 0xBFFFF	Color Video Memory
0xC0000 - 0xC7FFF	Video ROM BIOS
0xC8000 - 0xEFFFF	BIOS Shadow Area
0xF0000 - 0xFFFFF	System BIOS

Table 2.1: BIOS memory map.


```
db 0x55 ; Magic Numbers
db 0xaa ; To make disk bootable
```

You can compile and run the program above using `nasm` and `qemu` as discussed in Section 1.3.1. This program will not make any kind of visual display because we haven't told it to print anything. In the next example, we discuss how to print characters to the screen by calling the BIOS.

2.1.2 Printing Characters to the Screen

To print a character to the terminal, you have to pass a bit more information to the BIOS. The operation code `0x0E` in register `AH` tells the BIOS that we want to print a character. The ASCII code for the character to print goes in register `AL`. Registers `BH` and `BL` get the BIOS page number and foreground color respectively. I've written an example function that prints the character `A` to the terminal below.

Register	Meaning
AH	Command Code <code>0x0E</code>
AL	Character to write
BH	Page number (usually 0)
BL	Foreground color (7 for gray text)

Table 2.2: Parameters passed to the BIOS print command `INT 0x10`

```
[BITS 16]

boot: ; This file starts executing here.

    mov ah,0
    int 0x16 ; This instruction will hang until you press a key.
            ; On key press, BIOS returns with AL = ASCII code
    mov ah, $0xe
    mov bh,0 ; Page 0
    mov bl,7 ; Foreground color 7 (gray)
    int 0x10 ; Call the BIOS, print char in AL
    jmp boot ; Go back and get another character

times 510-($-$$) db 0 ; Skip to end of boot sector
db 0x55 ; Magic Numbers
db 0xaa ; To make disk bootable
```

Reading from the Disk

Before we talk about how to read from the disk, we need to discuss how data is organized on disks.

Sectors, Cylinders and Heads You're probably used to storing data on disks in files, but those are higher level abstractions provided by the operating system. At the most basic level, a disk is a *block device* that allows us to store information in blocks of 512 bytes. That's it.

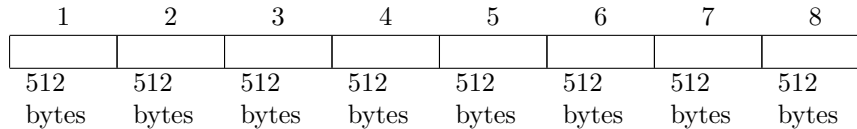
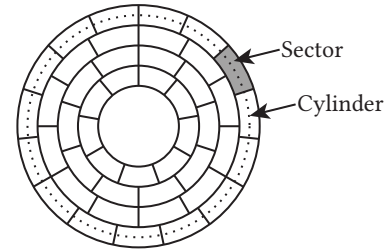


Figure 2.2: Layout of sectors on a disk.

Figure 2.1.2 shows the layout of 512 byte sectors on a disk. In each sector, we can write any 512-byte block of data we want. If we want to store larger blocks of data, they have to be split across multiple sectors. Later, when we want to read the data back, we just need to remember what sectors our data lives in.

Reading from the disk can be a little challenging. We need to tell the BIOS (1) the disk read command, (2) where we want to read from on the disk, (3) where we want the data to go in memory, and (4) how much data to read. For (1), the disk read command 0x02 goes in register AH. For (2), we need to specify a sector on the disk to read. The disk is addressed in cylinder, head, and sector numbers. Confusingly, cylinder and head indexes start at 0, and sector indexes start at 1. So the first sector on the disk in C:H:S notation is 0:0:1. The first hard disk number is 0x80. For (3), we need to give a complete 20-bit address (segment and offset) where we want the data from disk to be stored in memory. This is given to the BIOS in registers ES:BX. For (4), we need to tell the BIOS how many 512-byte sectors it should read from disk (passed in AL). Finally, call the BIOS with INT 0x13 to initiate the disk read.

Figure 2.1: Geometry of a magnetic hard disk. Sectors are arranged in concentric circles called *cylinders*.

Register	Meaning
----------	---------

AH	Command Code 0x02
AL	Number of sectors to read
ES	Segment address where data should be written
BX	16-bit offset into ES to write data to
CH	Low 8 bits of cylinder number
CL	Sector number
DH	Head number
DL	Drive number (0x80 is the first hard disk)

Table 2.3: Parameters passed to the BIOS disk read command INT 0x13.

The BIOS disk read functionality only allows us to read individual sectors from a disk into memory. It has no support for filesystems—that’s usually handled by the operating system.

2.1.3 DOS

In the early days of PCs, the commonly-used “operating system” was called *Disk Operating System* (DOS). DOS was not really an operating system in the usual sense. It was more like a set of extensions to the PC BIOS. The most useful support it provided above the BIOS was support for file systems and a rudimentary terminal interface. The BIOS, after all, can only read and write individual sectors from a disk—it doesn’t have support for creating files with names and locating those files on disk.

In those days, the most common form of storage was floppy disks, which could hold 720 kbytes or 1.44 mbytes of data depending on the disk type. Programs had to be small to fit on a disk, and computers could only run one program at a time. A program would be loaded off of its floppy disk into main memory, and it would use INT instructions to call BIOS or DOS functions.

For example, suppose a program calls the C library's `fopen` to open a file on disk. The C library would parse the file name and generate a *file control block* data structure from the file path passed to `fopen`. It would then call DOS's INT 0x21 function, which is a DOS-specific extension to the BIOS. DOS would examine the file control block and use its file system driver to locate the data sectors of the file on the disk. The DOS file system driver would make calls to the BIOS disk read function using calls to INT 0x13.

2.2 MBR

The MBR contains two critical components for the boot process: the bootloader and the partition table. The MBR is small—only 512 bytes. This doesn't leave much room to write a complicated program. Pretty much the only thing it can do is to find a real program that we actually want to run somewhere else on the boot disk and load that program into memory. But actually the MBR is so small that there isn't really even enough space to implement a proper hard disk driver, which of course is needed to read data from the disk. Fortunately, there is a disk driver available inside the BIOS—the MBR just needs to call that driver with the right parameters to load a real program from the disk. The MBR bootloader is a small program that calls the BIOS to load a larger second-stage bootloader from the hard drive.

2.3 GRUB

We only have one disk sector—512 bytes—for the MBR. There is not enough space in the MBR to store a complete bootloader that would be capable of loading a full kernel into memory. The bootloader usually needs lots of software to read the kernel from the disk: disk drivers, filesystem drivers, terminal drivers at least. PCs use a larger stage 1 bootloader to load the kernel from the disk. For the most part, each operating system has a custom bootloader. Linux uses a bootloader called GRUB (the GRand Unified Bootloader) which we will also use to load our kernel.

Base Address	Contents
0x000	Bootloader code
0x1be	Partition entry 1
0x1ce	Partition entry 2
0x1de	Partition entry 3
0x1ee	Partition entry 4
0x1fe	Boot signature 0x55
0x1ff	Boot signature 0xaa

GRUB does a lot of dirty work for us: it configures hardware, loads our kernel from disk into memory, and most importantly it puts the CPU into protected mode. Putting the CPU into protected mode is not a fun job: it requires setting up a bunch of messy tables in memory

2.3.1 Partitions

2.3.2 Binary File Formats

2.3.3 GRUB Hello World

Let's write a *Hello World* program that can be loaded and run by GRUB. The overall procedure we will follow is to create a disk image with an empty filesystem, install GRUB on the image, then copy our program to the filesystem. A disk image is a file that contains



Figure 2.3: Layout of the disk image we will create for the GRUB hello world.

a byte-for-byte copy of the contents of some disk. Disk images are useful for replicating OS installations—fresh installations of operating systems are usually done from disk images. Virtual machines also use disk images to store the contents of their virtual disks. Throughout this book, we will be creating disk images from scratch which we will use to boot our operating system.

We said before that hard disks are laid out as sectors of 512 bytes. The disk image that we create will also be organized in 512 byte blocks. A diagram of the disk image is shown below. The first 2048 sectors of the disk image ($2048 \times 512\text{bytes} = 1\text{Mbyte}$) will be assigned to the MBR and bootloader (GRUB). The first 2048 sectors are not part of any filesystem, so they can't store files in the way we are used to. Instead, they store the raw code of the bootloader loaded into memory by GRUB at boot time.

The first partition starts at sector 2048 and contains a FAT filesystem. FAT is a simple filesystem format developed for MS-DOS in the 1980s. We will be using it later in the book to store files for our operating system. For now, we will use the FAT filesystem to store the binary file for our hello world program.

The binary file that contains our program needs to have some special format to be recognized by GRUB. Its header is different from the standard ELF header used by Linux programs. GRUB expects our binary file to have a *multiboot* header which has fields that are similar to the ELF header. Other than the *multiboot* header, the rest of the binary file is the same. Fields in the Multiboot 2 header are explained in Table 2.3.3.

Field Name	Meaning
Magic Number	The integer 0xe85250d6 that tells GRUB that our binary is valid.
Flags	Binary-valued flags that tell GRUB what kind of data is contained in this multiboot2 entry.
Length	Length in bytes of this multiboot2 entry
Checksum	Integer to ensure that this multiboot2 entry has not been corrupted. The sum of all entries

Table 2.4: Fields in the Multiboot 2 header.

First, we will compile and link the C program. If you are compiling on an Intel or AMD x86 system, use the native compiler called `gcc`. If you are cross-compiling on a non-x86 system like an Apple M1 or other ARM device, use an x86 cross-compiler.

The compiler switches we use to build the kernel in this book tell the compiler to generate code for the i386 CPU using no external libraries. Most of the time when we compile programs, we want the compiler to bring in code from external libraries, which allow us to make calls to functions like `printf`. But those libraries assume that the program is running in an operating system. Since we are writing the operating system ourselves, we don't want to include those external functions. If we need them, we will have to write our own versions that don't depend on Linux syscalls.

```
user@system ~ $ gcc -c -ffreestanding -mgeneral-regs-only -mno-mmx -m32 -march=i386 -fno-pie
user@system ~ $ ld --section-start=.text=100000 --section-start=.rodata=0 -e main -melf_i386
```

Now we will create a 32 mbyte disk image called `rootfs.img` filled with zeros. The `dd` tool copies data in blocks from one file to another. It's kind of like `cat` for binary files. Here

```

1 #define MULTIBOOT2_HEADER_MAGIC 0xe85250d6
2
3 const unsigned int multiboot_header[] = {MULTIBOOT2_HEADER_MAGIC, // Magic number
4                                         0, // Flags
5                                         16, // Length
6                                         -(16+MULTIBOOT2_HEADER_MAGIC), // Checksum
7                                         0, // Type
8                                         12}; // Size
9
10 void main() {
11     unsigned short *vram = (unsigned short*)0xb8000; // Base address of video mem
12     const unsigned char color = 7; // gray text on black background
13     vram[0] = (((unsigned short)color)<<8) | (char)'H';
14     vram[1] = (((unsigned short)color)<<8) | (char)'E';
15     vram[2] = (((unsigned short)color)<<8) | (char)'L';
16     vram[3] = (((unsigned short)color)<<8) | (char)'L';
17     vram[4] = (((unsigned short)color)<<8) | (char)'O';
18
19     while(1);
20 }

```

Listing 7: The code for GRUB Hello World.

we are using it to copy from `/dev/zero`, which is a fake file that just reads the binary value 0. The output will be a new disk image file that we are creating called `rootfs.img` filled with 32 mbytes of binary 0.

```
user@system ~ $ dd if=/dev/zero of=rootfs.img bs=1M count=32
```

```

set timeout=5
set default=0 # Set the default menu entry

menuentry "Hello Grub" {
    set root=(hd0,msdos1)
    multiboot2 /hello # Load our hello grub program into memory
    boot             # Run hello grub
}

```

Next we will create the GRUB image, a file that contains the GRUB code to be copied to the first megabyte of our disk image. In Figure 2.3.3, we drew the GRUB code as starting at sector 1 and occupying up to sector 2047. The GRUB image does not include the MBR (that's in a separate file). The command below creates `grub.img`, which will be copied to our disk image starting at sector 1.

```
user@system ~ $ grub-mkimage -p "(hd0,msdos1)/boot" -o grub.img -0 i386-pc normal biosdisk multiboot
```

Now that we have a GRUB image, we can copy it to our disk image using `dd`. The first command below copies the MBR to sector 0 of our disk image. The MBR is stored in a 512-byte binary file on our Linux system called `boot.img`. The second command copies the grub image, which we created in the previous setp, to the disk image starting at sector 1.

```
user@system ~ $ dd if=/usr/lib/grub/i386-pc/boot.img of=rootfs.img conv=notrunc
user@system ~ $ dd if=grub.img of=rootfs.img conv=notrunc seek=1
```

Before we install our Hello Grub program onto our disk image, we need to format the image. Formatting is the process of setting up a partition table and a filesystem on a disk

or disk image. We will only create one partition on the disk image: starting at sector 2048 and occupying all sectors to the end of the disk.

```
user@system ~ $ echo 'start=2048, type=83, bootable' | sfdisk rootfs.img
```

Create a FAT16 filesystem on the first partition (starting at sector 2048) and copy our hello world binary to that filesystem.

```
user@system ~ $ mkfs.vfat --offset 2048 -F16 rootfs.img
user@system ~ $ mcopy -i rootfs.img@01M hello ::/
user@system ~ $ mmd -i rootfs.img@01M boot
user@system ~ $ mcopy -i rootfs.img@01M grub.cfg ::/boot
```

If all these commands worked, you should have a bootable disk image that can be used as a virtual disk for the `qemu` emulator. Use the command below to start the emulator and run your program in GRUB.

```
user@system ~ $ qemu-system-i386 -hda rootfs.img
```

2.3.4 Linker Scripts

In the previous section, we designed Grub Hello World to be a simple single-file C program which contained all of the components that Grub needs to boot an OS. We used a cheap hack to locate the multiboot header at address 0: we declared it as a `const unsigned int` array, and we told the linker to place all `const` variables (which live in the `rodata` section of the ELF file) at address 0. Our trick is fine for a small kernel, but we need a more scalable solution that will work for larger codebases.

Linker scripts allow us to specify how our binary is assembled and what address is assigned to its various components. As we discussed, every ELF file is composed of segments: the `.text` segment stores code, the `.data` segment stores initialized values of global variables, the `.rodata` segment stores constant data, and so on. Table 2.3.4 lists common default section names and their uses.

Compiling Hello Grub was clunky because we had to specify the address of each segment on the command line. Linker scripts let us specify the address of each section in a separate file, which gives us more control over the compilation. Listing 8 shows a simple linker script that is equivalent to the command line switches we used to compile Hello Grub.

Section Name	Contents
<code>.text</code>	Program code
<code>.bss</code>	Uninitialized global variables
<code>.data</code>	Initialized global variables
<code>.rodata</code>	Constant data

2.4 Exercises

Exercise 2.1 Write a `Makefile` that builds your Hello Grub program.

Exercise 2.2 Write a program that prints the character `a` to the screen repetitively. Fill up the whole screen with `as`.

Exercise 2.3 Write a delay function that causes your program to hang for a short amount of time. You can use a `for` loop. Modify your program from Exercise 2.2 to call your delay function and hang each time it prints the character `a`.

```
1 ENTRY(main) /* Name of the first function to run */
2 OUTPUT_FORMAT(elf32-i386) /* Specify ELF file format */
3
4 SECTIONS
5 {
6     . = 0; /* Set current location to address 0 */
7     .rodata : { *(.rodata) } /* Place .rodata section at address 0 */
8     /* Begin putting sections at 1 MiB, a conventional place for kernels to be
9        loaded at by the bootloader. */
10    . = 1M;
11    .text : { *(.text) } /* Place .text starting at 1M (0x100000) */
12    .data : { *(.data) } /* Place .data immediately after .text */
13    .bss : { *(.bss) } /* Place .bss section immediately after .data */
14 }
```

Listing 8: Simple linker script for GRUB Hello World.

Exercise 2.4 Write a function that can print a string to the screen by writing to video memory. The inputs to your function should be the string to print and the (x, y) location on the screen to start printing:

```
void print_string(char *string, unsigned int x, unsigned int y)
```

Exercise 2.5 Modify your `print_string` function from Exercise 2.4 to implement a complete terminal output driver. Your driver should keep track of the screen position where it last wrote a character. Each successive call to `print_string` should display characters to the terminal immediately after the last character that was printed without the need for the programmer to specify a screen location.

```
void print_string(char *string)
```


Index

8086

- Base Pointer(BP), 17
- Instruction Pointer (IP), 9
- real mode, 23
- registers, 9
- segment registers, 19

BIOS, 23

- drivers, 24

DOS, 26

epilogue, 15

i386

- protected mode, 21

nasm, 11

prologue, 15

qemu, 11

stack frame, 14

strlen, 13

Video Memory, 20