

Chapter 1

Bootstrap

Hardware

A computer's CPU (central processing unit, or processor) runs a conceptually simple loop: it inspects the value of a register called the program counter, reads a machine instruction from that address in memory, advances the program counter past the instruction, and executes the instruction. Repeat. If the execution of the instruction does not modify the program counter, this loop will interpret the memory pointed at by the program counter as a sequence of machine instructions to run one after the other. Instructions that do change the program counter implement conditional branches, unconditional branches, and function calls.

The execution engine is useless without the ability to store and modify program data. The fastest storage for data is provided by the processor's register set. A register is a storage cell inside the processor itself, capable of holding a machine word-sized value (typically 16, 32, or 64 bits). Data stored in registers can typically be read or written quickly, in a single CPU cycle. The x86 provides eight general purpose 32-bit registers—`%eax`, `%ebx`, `%ecx`, `%edx`, `%edi`, `%esi`, `%ebp`, and `%esp`—and a program counter `%eip` (the “instruction pointer”). The common `e` prefix stands for extended, as these are 32-bit extensions of the 16-bit registers `%ax`, `%bx`, `%cx`, `%dx`, `%di`, `%si`, `%bp`, `%sp`, and `%ip`. The two register sets are aliased so that, for example, `%ax` is the bottom half of `%eax`: writing to `%ax` changes the value stored in `%eax` and vice versa. The first four registers also have names for the bottom two 8-bit bytes: `%al` and `%ah` denote the low and high 8 bits of `%ax`; `%bl`, `%bh`, `%cl`, `%ch`, `%dl`, and `%dh` continue the pattern. In addition to these registers, the x86 has eight 80-bit floating-point registers as well as a handful of special-purpose registers like the control registers `%cr0`, `%cr2`, `%cr3`, and `%cr4`; the debug registers `%dr0`, `%dr1`, `%dr2`, and `%dr3`; the segment registers `%cs`, `%ds`, `%es`, `%fs`, `%gs`, and `%ss`; and the global and local descriptor table pseudo-registers `%gdtr` and `%ldtr`. The control, segment selector, and descriptor table registers are important to any operating system, as we will see in this chapter. The floating-point and debug registers are less interesting and not used by `xv6`.

Registers are fast but expensive. Most processors provide at most a few tens of general-purpose registers. The next conceptual level of storage is the main random-access memory (RAM). Main memory is 10-100x slower than a register, but it is much cheaper, so there can be more of it. An x86 processor has a few dozen registers, but a typical PC today has gigabytes of main memory. Because of the enormous differences in both access speed and size between registers and main memory, most processors,

including the x86, store copies of recently-accessed sections of main memory in on-chip cache memory. The cache memory serves as a middle ground between registers and memory both in access time and in size. Today's x86 processors typically have two levels of cache, a small first-level cache with access times relatively close to the processor's clock rate and a larger second-level cache with access times in between the first-level cache and main memory. This table shows actual numbers for an Intel Core 2 Duo system:

Intel Core 2 Duo E7200 at 2.53 GHz
TODO: Plug in non-made-up numbers!

storage	access time	size
register	0.6 ns	64 bytes
L1 cache	0.5 ns	64 kilobytes
L2 cache	10 ns	4 megabytes
main memory	100 ns	4 gigabytes

For the most part, x86 processors hide the cache from the operating system, so we can think of the processor as having just two kinds of storage—registers and memory—and not worry about the distinctions between the different levels of the memory hierarchy. The exceptions—the only reasons an x86 operating system needs to worry about the memory cache—are concurrency (Chapter 4) and device drivers (Chapter 6).

One reason memory access is so much slower than register access is that the memory is a set of chips physically separate from the processor chip. To allow the processor to communicate with the memory, there is a collection of wires, called a bus, running between the two. A simple mental model is that some of the wires, called lines, carry address bits; some carry data bits. To read a value from main memory, the processor sends high or low voltages representing 1 or 0 bits on the address lines and a 1 on the “read” line for a prescribed amount of time and then reads back the value by interpreting the voltages on the data lines. To write a value to main memory, the processor sends appropriate bits on the address and data lines and a 1 on the “write” line for a prescribed amount of time. This model is an accurate description of the earliest x86 chips, but it is a drastic oversimplification of a modern system. Even so, thanks to the processor-centric view the operating system has of the rest of the computer, this simple model suffices to understand a modern operating system. The details of modern I/O buses are the province of computer architecture textbooks.

Processors must communicate not just with memory but with hardware devices too. The x86 processor provides special in and out instructions that read and write values from device addresses called I/O ports. The hardware implementation of these instructions is essentially the same as reading and writing memory. Early x86 processors had an extra address line: 0 meant read/write from an I/O port and 1 meant read/write from main memory. Each hardware device monitors these lines for reads and writes to its assigned range of I/O ports. A device's ports let the software configure the device, examine its status, and cause the device to take actions; for example, software can use I/O port reads and writes to cause the disk interface hardware to read and write sectors on the disk.

Many computer architectures have no separate device access instructions. Instead the devices have fixed memory addresses and the processor communicates with the device (at the operating system's behest) by reading and writing values at those addresses. In fact, modern x86 architectures use this technique, called memory-mapped I/O, for most high-speed devices such as network, disk, and graphics controllers. For reasons of backwards compatibility, though, the old `in` and `out` instructions linger, as do legacy hardware devices that use them, such as the IDE disk controller, which we will see shortly.

Bootstrap

When an x86 PC boots, it starts executing a program called the BIOS, which is stored in flash memory on the motherboard. The BIOS's job is to prepare the hardware and then transfer control to the operating system. Specifically, it transfers control to code loaded from the boot sector, the first 512-byte sector of the boot disk. The BIOS loads a copy of that sector into memory at `0x7c00` and then jumps (sets the processor's `%ip`) to that address. When the boot sector begins executing, the processor is simulating an Intel 8088, the CPU chip in the original IBM PC released in 1981. The xv6 boot sector's job is to put the processor in a more modern operating mode, to load the xv6 kernel from disk into memory, and then to transfer control to the kernel. In xv6, the boot sector comprises two source files, one written in a combination of 16-bit and 32-bit x86 assembly (`bootasm.S`; (1000)) and one written in C (`bootmain.c`; (1200)). This chapter examines the operation of the xv6 boot sector, from the time the BIOS starts it to the time it transfers control to the kernel proper. The boot sector is a microcosm of the kernel itself: it contains low-level assembly and C code, it manages its own memory, and it even has a device driver, all in under 512 bytes of machine code.

Code: Assembly bootstrap

The first instruction in the boot sector is `ccli` (1015), which disables processor interrupts. Interrupts are a way for hardware devices to invoke operating system functions called interrupt handlers. The BIOS is a tiny operating system, and it might have set up its own interrupt handlers as part of the initializing the hardware. But the BIOS isn't running anymore—xv6 is, or will be—so it is no longer appropriate or safe to handle interrupts from hardware devices. When xv6 is ready (in Chapter 3), it will re-enable interrupts.

The processor is in "real mode," in which it simulates an Intel 8088. In real mode there are eight 16-bit general-purpose registers, but the processor sends 20 bits of address to memory. The segment registers `%cs`, `%ds`, `%es`, and `%ss` provide the additional bits necessary to generate 20-bit memory addresses from 16-bit registers. When a program refers to a memory address, the processor automatically adds 16 times the value of one of the segment registers; these registers are 16 bits wide. Which segment register is usually implicit in the kind of memory reference: instruction fetches use `%cs`,

data reads and writes use `%ds`, and stack reads and writes use `%ss`. We'll call the addresses the processor chip sends to memory "physical addresses," and the addresses that programs directly manipulate "virtual addresses." People often write a full real-mode virtual memory reference as *segment:offset*, indicating the value of the relevant segment register and the address supplied by the program.

The BIOS does not guarantee anything about the contents of `%ds`, `%es`, `%ss`, so first order of business after disabling interrupts is to set `%ax` to zero and then copy that zero into `%ds`, `%es`, and `%ss` (1018-1021).

A virtual *segment:offset* can yield a 21-bit physical address, but the Intel 8088 could only address 20 bits of memory, so it discarded the top bit: `0xffff0+0xffff = 0x10ffef`, but virtual address `0xffff:0xffff` on the 8088 referred to physical address `0x0ffef`. Some early software relied on the hardware ignoring the 21st address bit, so when Intel introduced processors with more than 20 bits of physical address, IBM provided a compatibility hack that is a requirement for PC-compatible hardware. If the second bit of the keyboard controller's output port is low, the 21st physical address bit is always cleared; if high, the 21st bit acts normally. The boot sector must enable the 21st address bit using I/O to the keyboard controller on ports `0x64` and `0x60` (1023-1041).

Real mode's 16-bit general-purpose and segment registers make it awkward for a program to use more than 65,536 bytes of memory, and impossible to use more than a megabyte. Modern x86 processors have a "protected mode" which allows physical addresses to have many more bits, and a "32-bit" mode that causes registers, virtual addresses, and most integer arithmetic to be carried out with 32 bits rather than 16. The xv6 boot sequence enables both modes as follows.

In protected mode, a segment register is an index into a segment descriptor table. Each table entry specifies a base physical address, a maximum virtual address called the limit, and permission bits for the segment. These permissions are the protection in protected mode: they can be used to make sure that one program cannot access memory belonging to another program.

xv6 makes almost no use of segments (it uses the paging hardware instead, as the next chapter describes). The boot code sets up the segment descriptor table `gdt` (1092-1095) so that all segments have a base address of zero and the maximum possible limit (four gigabytes). The table has a null entry, one entry for executable code, and one entry to data. The code segment descriptor has a flag set that indicates that the code should run in 32-bit mode. The boot code executes an `lgdt` instruction (1054) to set the processor's global descriptor table (GDT) register with the value `gdt_desc` (1097-1099), which in turns points at the table `gdt`.

Once it has loaded the GDT register, the boot sector enables protected mode by setting the 1 bit (`CR0_PE`) in register `%cr0` (1055-1057). Enabling protected mode does not immediately change how the processor translates virtual to physical addresses or whether it is in 32-bit mode; it is only when one loads a new value into a segment register that the processor reads the GDT and changes its internal segmentation settings. Thus the processor continues to execute in 16-bit mode with the same segment translations as before. The switch to 32-bit mode happens when the code executes a far jump (`ljmp`) instruction (1059). The jump continues execution at the next line (1066)

but in doing so sets `%cs` to refer to the code descriptor entry in `gdt`. The entry describes a 32-bit code segment, so the processor switches into 32-bit mode. The boot sector code has nursed the processor through an evolution from 8088 through 80286 to 80386.

The boot sector's first action in 32-bit mode is to initialize the data segment registers with `SEG_KDATA` (1068-1071). Virtual address now map directly to physical addresses. The only step left before executing C code is to set up a stack in an unused region of memory. The memory from `0xa0000` to `0x100000` is typically littered with device memory regions, and the `xv6` kernel expects to be placed at `0x100000`. The boot sector itself is at `0x7c00` through `0x7d00`. Essentially any other section of memory would be a fine location for the stack. The boot sector chooses `0x7c00` (known in this file as `$start`) as the top of the stack; the stack will grow down from there, toward `0x0000`, away from the boot sector code.

Finally the boot sector calls the C function `bootmain` (1078). `Bootmain`'s job is to load and run the kernel. It only returns if something has gone wrong. In that case, the code sends a few output words on port `0x8a00` (1080-1086). On real hardware, there is no device connected to that port, so this code does nothing. If the boot sector is running inside the PC simulator `Bochs`, port `0x8a00` is connected to `Bochs` itself; the code sequence triggers a `Bochs` debugger breakpoint. `Bochs` or not, the code then executes an infinite loop (1087-1088). A real boot sector might attempt to print an error message first.

Code: C bootstrap

The C part of the boot sector, `bootmain.c` (1200), loads a kernel from an IDE disk into memory and then starts executing it. The kernel is an ELF format binary, defined in `elf.h`. An ELF binary is an ELF file header, `struct elfhdr` (0955), followed by a sequence of program section headers, `struct proghdr` (0974). Each `proghdr` describes a section of the kernel that must be loaded into memory. These headers typically take up the first hundred or so bytes of the binary. To get access to the headers, `bootmain` loads the first 4096 bytes of the file, a gross overestimation of the amount needed (1213). It places the in-memory copy at address `0x10000`, another out-of-the-way memory address.

`bootmain` casts freely between pointers and integers (1223, 1226, and so on). Programming languages distinguish the two to catch errors, but the underlying processor sees no difference. An operating system must work at the processor's level; occasionally it will need to treat a pointer as an integer or vice versa. C allows these conversions, in contrast to languages like Pascal and Java, precisely because one of the first uses of C was to write an operating system: Unix.

Back in the boot sector, what should be an ELF binary header has been loaded into memory at address `0x10000` (1213). The next step is to check that the first four bytes of the header, the so-called magic number, are the bytes `0x7F`, `'E'`, `'L'`, `'F'`, or `ELF_MAGIC` (0952). All ELF binary headers are required to begin with this magic number as identification. If the ELF header has the right magic number, the boot sector assumes that the binary is well-formed. There are many other sanity checks that a prop-

er ELF loader would do, as we will see in Chapter 9, but the boot sector doesn't have the code space. Checking the magic number guards against simply forgetting to write a kernel to the disk, not against malicious binaries.

An ELF header points at a small number of program headers (`proghdrs`) describing the sections that make up the running kernel image. Each `proghdr` gives a virtual address (`va`), the location where the section's content lies on the disk relative to the start of the ELF header (`offset`), the number of bytes to load from the file (`filesz`), and the number of bytes to allocate in memory (`memsz`). If `memsz` is larger than `filesz`, the bytes not loaded from the file are to be zeroed. This is more efficient, both in space and I/O, than storing the zeroed bytes directly in the binary. As an example, the xv6 kernel has two loadable program sections, code and data:

```
# objdump -p kernel

kernel:      file format elf32-i386

Program Header:
LOAD off    0x00001000 vaddr 0x00100000 paddr 0x00100000 align 2**12
      filesz 0x000063ca memsz 0x000063ca flags r-x
LOAD off    0x000073e0 vaddr 0x001073e0 paddr 0x001073e0 align 2**12
      filesz 0x0000079e memsz 0x000067e4 flags rw-
STACK off   0x00000000 vaddr 0x00000000 paddr 0x00000000 align 2**2
      filesz 0x00000000 memsz 0x00000000 flags rwx
```

Notice that the second section, the data section, has a `memsz` larger than its `filesz`: the first 0x79e bytes are loaded from the kernel binary and the remaining 0x6046 bytes are zeroed.

`Bootmain` uses the addresses in the `proghdr` to direct the loading of the kernel. It reads each section's content starting from the disk location `offset` bytes after the start of the ELF header, and writes to memory starting at address `va`. `Bootmain` calls `readseg` to load data from disk (1237) and calls `stosb` to zero the remainder of the segment (1239). `Stosb` (0442) uses the x86 instruction `rep stosb` to initialize every byte of a block of memory.

`Readseg` (1279) reads at least `count` bytes from the disk `offset` into memory at `va`. The x86 IDE disk interface operates in terms of 512-byte chunks called sectors, so `readseg` may read not only the desired section of memory but also some bytes before and after, depending on alignment. For the second program segment in the example above, the boot sector will call `readseg((uchar*)0x1073e0, 0x73e0, 0x79e)`. Due to sector granularity, this call is equivalent to `readseg((uchar*)0x107200, 0x7200, 0xa00)`: it reads 0x1e0 bytes before the desired memory region and 0x82 bytes afterward. In practice, this sloppy behavior turns out not to be a problem (see exercise XXX). `Readseg` begins by computing the ending virtual address, the first memory address above `va` that doesn't need to be loaded from disk (1283), and rounding `va` down to a sector-aligned disk offset. Then it converts the offset from a byte offset to a sector offset; it adds 1 because the kernel starts at disk sector 1 (disk sector 0 is the boot sector). Finally, it calls `readsect` to read each sector into memory.

`Readsect` (1260) reads a single disk sector. It is our first example of a device driver, albeit a tiny one. `Readsect` begins by calling `waitdisk` to wait until the disk sig-

nals that it is ready to accept a command. The disk does so by setting the top two bits of its status byte (connected to input port 0x1f7) to 01. `waitdisk` (1251) reads the status byte until the bits are set that way. Chapter 6 will examine more efficient ways to wait for hardware status changes, but busy waiting like this (also called polling) is fine for the boot sector.

Once the disk is ready, `readsect` issues a read command. It first writes command arguments—the sector count and the sector number (offset)—to the disk registers on output ports 0x1f2-0x1f6 (1264-1268). The bits 0xe0 in the write to port 0x1f6 signal to the disk that 0x1f3-0x1f6 contain a sector number (a so-called linear block address), in contrast to a more complicated cylinder/head/sector address used in early PC disks. After writing the arguments, `readsect` writes to the command register to trigger the read (1254). The command 0x20 is “read sectors.” Now the disk will read the data stored in the specified sectors and make it available in 32-bit pieces on input port 0x1f0. `waitdisk` (1251) waits until the disk signals that the data is ready, and then the call to `insl` reads the 128 (SECTSIZE/4) 32-bit pieces into memory starting at `dst` (1273).

`inb`, `outb`, and `insl` are not ordinary C functions. They are inlined functions whose bodies are assembly language fragments (0403, 0421, 0412). When `gcc` sees the call to `inb` (1254), the inlined assembly causes it to emit a single `inb` instruction. This style allows the use of low-level instructions like `inb` and `outb` while still writing the control logic in C instead of assembly.

The implementation of `insl` (0412) is worth looking at more closely. `rep insl` is actually a tight loop masquerading as a single instruction. The `rep` prefix executes the following instruction `%ecx` times, decrementing `%ecx` after each iteration. The `insl` instruction reads a 32-bit value from port `%dx` into memory at address `%edi` and then increments `%edi` by 4. Thus `rep insl` copies $4 \times \%ecx$ bytes, in 32-bit chunks, from port `%dx` into memory starting at address `%edi`. The register annotations tell GCC to prepare for the assembly sequence by storing `dst` in `%edi`, `cnt` in `%ecx`, and port in `%dx`. Thus the `insl` function copies $4 \times cnt$ bytes from the 32-bit port `port` into memory starting at `dst`. The `cld` instruction clears the processor’s direction flag, so that the `insl` instruction increments `%edi`; when the flag is set, `insl` decrements `%edi` instead. The x86 calling convention does not define the state of the direction flag on entry to a function, so each use of an instruction like `insl` must initialize it to the desired value.

The boot loader is almost done. `bootmain` loops calling `readseg`, which loops calling `readsect` (1235-1240). At the end of the loop, `bootmain` has loaded the kernel into memory. Now it is time to run the kernel. The ELF header specifies the kernel entry point, the `%eip` where the kernel expects to be started (just as the boot loader expected to be started at 0x7c00). `bootmain` casts the entry point integer to a function pointer and calls that function, essentially jumping to the kernel’s entry point (1244-1245). The kernel should not return, but if it does, `bootmain` will return, and then `bootasm.S` will attempt a Bochs breakpoint and then loop forever.

Where is the kernel in memory? `bootmain` does not directly decide; it just follows the directions in the ELF headers. The “linker” creates the ELF headers, and the `vx6` Makefile that calls the linker tells it that the kernel should start at 0x100000.

Assuming all has gone well, the kernel entry pointer will refer to the kernel’s `main`

function (see `main.c`). The next chapter continues there.

Real world

The boot sector described in this chapter compiles to around 470 bytes of machine code, depending on the optimizations used when compiling the C code. In order to fit in that small amount of space, the xv6 boot sector makes a major simplifying assumption, that the kernel has been written to the boot disk contiguously starting at sector 1. More commonly, kernels are stored in ordinary file systems, where they may not be contiguous, or are loaded over a network. These complications require the boot loader to be able to drive a variety of disk and network controllers and understand various file systems and network protocols. In other words, the boot loader itself must be a small operating system. Since such complicated boot loaders certainly won't fit in 512 bytes, most PC operating systems use a two-step boot process. First, a simple boot sector like the one in this chapter loads a full-featured boot-loader from a known disk location, often relying on the less space-constrained BIOS for disk access rather than trying to drive the disk itself. Then the full loader, relieved of the 512-byte limit, can implement the complexity needed to locate, load, and execute the desired kernel.

TODO: Also, x86 does not imply BIOS: Macs use EFI. I wonder if the Mac has an A20 line.

Exercises

1. Look at the kernel load addresses; why doesn't the sloppy readsect cause problems?
2. something about BIOS lasting longer + security problems
3. Suppose you wanted `bootmain()` to load the kernel at `0x200000` instead of `0x100000`, and you did so by modifying `bootmain()` to add `0x100000` to the `va` of each ELF section. Something would go wrong. What?