

xv6: a simple, Unix-like teaching operating system

Russ Cox

Frans Kaashoek

Robert Morris

August 31, 2020

Contents

1	Operating system interfaces	9
1.1	Processes and memory	10
1.2	I/O and File descriptors	13
1.3	Pipes	15
1.4	File system	17
1.5	Real world	19
1.6	Exercises	20
2	Operating system organization	21
2.1	Abstracting physical resources	22
2.2	User mode, supervisor mode, and system calls	22
2.3	Kernel organization	23
2.4	Code: xv6 organization	24
2.5	Process overview	24
2.6	Code: starting xv6 and the first process	27
2.7	Real world	28
2.8	Exercises	28
3	Page tables	29
3.1	Paging hardware	29
3.2	Kernel address space	31
3.3	Code: creating an address space	33
3.4	Physical memory allocation	34
3.5	Code: Physical memory allocator	34
3.6	Process address space	35
3.7	Code: sbrk	36
3.8	Code: exec	37
3.9	Real world	38
3.10	Exercises	39
4	Traps and system calls	41
4.1	RISC-V trap machinery	42
4.2	Traps from user space	43

4.3	Code: Calling system calls	44
4.4	Code: System call arguments	45
4.5	Traps from kernel space	46
4.6	Page-fault exceptions	46
4.7	Real world	48
4.8	Exercises	48
5	Interrupts and device drivers	49
5.1	Code: Console input	49
5.2	Code: Console output	50
5.3	Concurrency in drivers	51
5.4	Timer interrupts	51
5.5	Real world	52
5.6	Exercises	53
6	Locking	55
6.1	Race conditions	56
6.2	Code: Locks	58
6.3	Code: Using locks	60
6.4	Deadlock and lock ordering	60
6.5	Locks and interrupt handlers	62
6.6	Instruction and memory ordering	62
6.7	Sleep locks	63
6.8	Real world	64
6.9	Exercises	64
7	Scheduling	67
7.1	Multiplexing	67
7.2	Code: Context switching	68
7.3	Code: Scheduling	69
7.4	Code: mycpu and myproc	70
7.5	Sleep and wakeup	71
7.6	Code: Sleep and wakeup	74
7.7	Code: Pipes	75
7.8	Code: Wait, exit, and kill	76
7.9	Real world	77
7.10	Exercises	79
8	File system	81
8.1	Overview	81
8.2	Buffer cache layer	82
8.3	Code: Buffer cache	83
8.4	Logging layer	84

8.5	Log design	85
8.6	Code: logging	86
8.7	Code: Block allocator	87
8.8	Inode layer	87
8.9	Code: Inodes	89
8.10	Code: Inode content	90
8.11	Code: directory layer	91
8.12	Code: Path names	92
8.13	File descriptor layer	93
8.14	Code: System calls	94
8.15	Real world	95
8.16	Exercises	96
9	Concurrency revisited	99
9.1	Locking patterns	99
9.2	Lock-like patterns	100
9.3	No locks at all	100
9.4	Parallelism	101
9.5	Exercises	102
10	Summary	103

Foreword and acknowledgments

This is a draft text intended for a class on operating systems. It explains the main concepts of operating systems by studying an example kernel, named xv6. xv6 is modeled on Dennis Ritchie's and Ken Thompson's Unix Version 6 (v6) [14]. xv6 loosely follows the structure and style of v6, but is implemented in ANSI C [6] for a multi-core RISC-V [12].

This text should be read along with the source code for xv6, an approach inspired by John Lions' Commentary on UNIX 6th Edition [9]. See <https://pdos.csail.mit.edu/6.S081> for pointers to on-line resources for v6 and xv6, including several lab assignments using xv6.

We have used this text in 6.828 and 6.S081, the operating systems classes at MIT. We thank the faculty, teaching assistants, and students of those classes who have all directly or indirectly contributed to xv6. In particular, we would like to thank Adam Belay, Austin Clements, and Nickolai Zeldovich. Finally, we would like to thank people who emailed us bugs in the text or suggestions for improvements: Abutalib Aghayev, Sebastian Boehm, Anton Burtsev, Raphael Carvalho, Tej Chajed, Rasit Eskicioglu, Color Fuzzy, Giuseppe, Tao Guo, Naoki Hayama, Robert Hilder-
man, Wolfgang Keller, Austin Liew, Pavan Maddamsetti, Jacek Masiulaniec, Michael McConville, m3hm00d, miguelgvieira, Mark Morrissey, Harry Pan, Askar Safin, Salman Shah, Adeodato Simó, Ruslan Savchenko, Pawel Szczurko, Warren Toomey, tyfkda, tzerbib, Xi Wang, and Zou Chang Wei.

If you spot errors or have suggestions for improvement, please send email to Frans Kaashoek and Robert Morris (kaashoek,rtm@csail.mit.edu).

Chapter 1

Operating system interfaces

The job of an operating system is to share a computer among multiple programs and to provide a more useful set of services than the hardware alone supports. An operating system manages and abstracts the low-level hardware, so that, for example, a word processor need not concern itself with which type of disk hardware is being used. An operating system shares the hardware among multiple programs so that they run (or appear to run) at the same time. Finally, operating systems provide controlled ways for programs to interact, so that they can share data or work together.

An operating system provides services to user programs through an interface. Designing a good interface turns out to be difficult. On the one hand, we would like the interface to be simple and narrow because that makes it easier to get the implementation right. On the other hand, we may be tempted to offer many sophisticated features to applications. The trick in resolving this tension is to design interfaces that rely on a few mechanisms that can be combined to provide much generality.

This book uses a single operating system as a concrete example to illustrate operating system concepts. That operating system, xv6, provides the basic interfaces introduced by Ken Thompson and Dennis Ritchie’s Unix operating system [14], as well as mimicking Unix’s internal design. Unix provides a narrow interface whose mechanisms combine well, offering a surprising degree of generality. This interface has been so successful that modern operating systems—BSD, Linux, Mac OS X, Solaris, and even, to a lesser extent, Microsoft Windows—have Unix-like interfaces. Understanding xv6 is a good start toward understanding any of these systems and many others.

As Figure 1.1 shows, xv6 takes the traditional form of a *kernel*, a special program that provides services to running programs. Each running program, called a *process*, has memory containing instructions, data, and a stack. The instructions implement the program’s computation. The data are the variables on which the computation acts. The stack organizes the program’s procedure calls. A given computer typically has many processes but only a single kernel.

When a process needs to invoke a kernel service, it invokes a *system call*, one of the calls in the operating system’s interface. The system call enters the kernel; the kernel performs the service and returns. Thus a process alternates between executing in *user space* and *kernel space*.

The kernel uses the hardware protection mechanisms provided by a CPU¹ to ensure that each

¹This text generally refers to the hardware element that executes a computation with the term *CPU*, an acronym for central processing unit. Other documentation (e.g., the RISC-V specification) also uses the words processor, core, and hart instead of CPU.

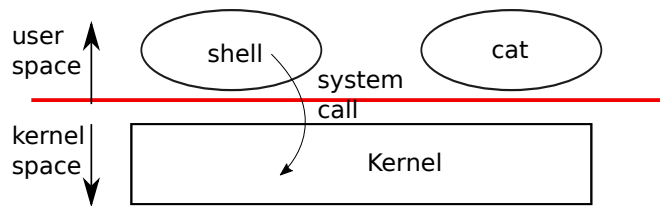


Figure 1.1: A kernel and two user processes.

process executing in user space can access only its own memory. The kernel executes with the hardware privileges required to implement these protections; user programs execute without those privileges. When a user program invokes a system call, the hardware raises the privilege level and starts executing a pre-arranged function in the kernel.

The collection of system calls that a kernel provides is the interface that user programs see. The xv6 kernel provides a subset of the services and system calls that Unix kernels traditionally offer. Figure 1.2 lists all of xv6’s system calls.

The rest of this chapter outlines xv6’s services—processes, memory, file descriptors, pipes, and a file system—and illustrates them with code snippets and discussions of how the *shell*, Unix’s command-line user interface, uses them. The shell’s use of system calls illustrates how carefully they have been designed.

The shell is an ordinary program that reads commands from the user and executes them. The fact that the shell is a user program, and not part of the kernel, illustrates the power of the system call interface: there is nothing special about the shell. It also means that the shell is easy to replace; as a result, modern Unix systems have a variety of shells to choose from, each with its own user interface and scripting features. The xv6 shell is a simple implementation of the essence of the Unix Bourne shell. Its implementation can be found at (user/sh.c:1).

1.1 Processes and memory

An xv6 process consists of user-space memory (instructions, data, and stack) and per-process state private to the kernel. Xv6 *time-shares* processes: it transparently switches the available CPUs among the set of processes waiting to execute. When a process is not executing, xv6 saves its CPU registers, restoring them when it next runs the process. The kernel associates a process identifier, or `PID`, with each process.

A process may create a new process using the `fork` system call. `Fork` creates a new process, called the *child process*, with exactly the same memory contents as the calling process, called the *parent process*. `Fork` returns in both the parent and the child. In the parent, `fork` returns the child’s `PID`; in the child, `fork` returns zero. For example, consider the following program fragment written in the C programming language [6]:

```
int pid = fork();
if(pid > 0){
    printf("parent: child=%d\n", pid);
```