

# Overview of Functions of an Operating System

Norman Matloff  
University of California, Davis  
©2001, N. Matloff

May 31, 2001

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	It's Just a Program! . . . . .	2
1.2	What Is an OS for, Anyway? . . . . .	3
1.3	A Bit More on System Calls . . . . .	4
1.4	Making These Concepts Concrete: Commands You Can Try Yourself . . . . .	5
<b>2</b>	<b>System Bootup</b>	<b>6</b>
<b>3</b>	<b>Application Program Loading</b>	<b>7</b>
<b>4</b>	<b>Timesharing</b>	<b>8</b>
4.1	Many Processes, Taking Turns . . . . .	8
4.2	Example of OS Code: Linux for Intel CPUs . . . . .	9
4.3	Process States . . . . .	10
4.4	How the Process Tree Is Built . . . . .	10
4.5	Making These Concepts Concrete: Commands You Can Try Yourself . . . . .	11
<b>5</b>	<b>Virtual Memory</b>	<b>12</b>
5.1	Make Sure You Understand the Goals . . . . .	12
5.2	Example of Virtual Nature of Addresses . . . . .	13
5.3	Overview of How the Goals Are Achieved . . . . .	13

5.4	Creation and Maintenance of the Page Table . . . . .	14
5.5	Details on Usage of the Page Table . . . . .	15
5.5.1	Virtual-to-Physical Address Translation, Page Table Lookup . . . . .	15
5.5.2	Page Faults . . . . .	16
5.5.3	Access Violations . . . . .	17
5.6	Improving Performance . . . . .	18
5.7	Intel Page Tables . . . . .	18
5.8	Making These Concepts Concrete: Commands You Can Try Yourself . . . . .	19
<b>A</b>	<b>Hardware Interrupts</b>	<b>19</b>
A.1	General Operation . . . . .	19
A.2	Some Details for Intel CPUs and PCs . . . . .	20

## 1 Introduction

### 1.1 It's Just a Program!

First and foremost, it is vital to understand that an **operating system** (OS) is just a program – a very large, very complex program, but still just a program. The OS provides support for the loading and execution of other programs (which we will refer to below as “application programs”), and the OS will set things up so that it has some special privileges which user programs don’t have, but in the end, the OS is simply a program.

For example, when your program, say **a.out**,<sup>1</sup> is running, the OS is *not* running. Thus the OS has no power to suspend your program while your program is running – since the OS isn’t running! This is a key concept, so let’s first make sure what the statement even means.

What does it mean for a program to be “running” anyway? Recall that the CPU is constantly performing its fetch/execute/fetch/execute/... cycle. For each fetch, it fetches whatever instruction the Program Counter (PC) is pointing to. If the PC is currently pointing to an instruction in your program, then your program is running! Each time an instruction of your program executes, the circuitry in the CPU will update the PC, having it point to either the next instruction (the usual case) or an instruction located elsewhere in your program (in the case of jumps).

The point is that the only way your program can stop running is if the PC is changed to point to another program, say the OS. How might this happen? Other than cases involving bugs in your program, there are only two ways this can occur:

---

<sup>1</sup>Or it could be a program which you didn’t write yourself, say **gcc**.

- Your program can *voluntarily* relinquish the CPU to the OS. It does this via a **system call**, which is a call to some function in the operating system which provides some useful service.

For example, suppose the C source file from which **a.out** was compiled had a call to **scanf()**. The **scanf()** function is a C library function, which was linked into **a.out** during compilation of **a.out**. But **scanf()** itself calls **read()**, a function within the OS. So, when **a.out** reaches the **scanf()** call, that will result in a call to the OS, but after the OS does the read from the keyboard, the OS will return to **a.out**.<sup>2</sup>

- The other possibility is that a hardware interrupt occurs. (See Appendix A for an introduction to hardware interrupts, and their implementation on Intel-based machines.) This is a signal – a physical pulse of current along an **interrupt-request** line in the bus – from some input/output (I/O) device such as the keyboard to the CPU. The circuitry in the CPU is designed to then jump to a place in memory which we designated upon bootup of the machine. This will be a place in the OS, so the OS will now run. The OS will attend to the I/O device, e.g. record the keystroke in the case of the keyboard, and then return to the interrupted program.

Note in our keystroke example that the keystroke may not have been made by you. While your program is running, some other user of the machine may hit a key. The interrupt will cause your program to be suspended; the OS will run the device driver for whichever device caused the interrupt – the keyboard, if the person was sitting at the console of the machine, or the network interface card, if the person was logged in remotely, say via **telnet** – which will record the keystroke in the buffer belonging to that other user; and the OS will do IRET, causing your program to resume.

So, when your program is running, it is king. The OS has no power to stop it. The only ways your program can stop running is if it voluntarily does so or is forced to stop by action occurring at an I/O device.

## 1.2 What Is an OS for, Anyway?

The major functions of a typical OS are to:

- load application programs for execution
- provide services, e.g. I/O, in the form of functions which the application programs can call
- enable **timesharing**, in which many application programs seem to be running simultaneously but are actually “taking turns,” *by coordinating with hardware operations*
- enable **virtual memory** operations for application programs, which both allows flexible use of memory and enforces security, again *by coordinating with hardware operations*
- maintain the file system
- manage I/O, including networking

---

<sup>2</sup>An exception is the system call **exit()**, which is called by your program when it is finished with execution. If you do not include this call in your C source file, the compiler will put one in for you.

Note carefully the interaction of the OS with the hardware. The OS manages the I/O devices, thus shielding the authors of application programs from having to deal directly with them. For example, suppose you wish to have your program read from a disk file. It would be a real nuisance if you had to deal with the minute details of the physical locations on disk of the various pieces of the file. Instead, you simply call `fopen()` and `fread()`, and the OS will worry about that for you. The OS knows these locations, by the way, because when the file was created, the creation was done through the OS too; at that time, the OS chose those locations, and recorded them.

In addition, if the hardware allows for it, and if the OS is designed to make use of it,<sup>3</sup> then the OS not only *relieves* the application programmer from having to perform the physical accesses of the disk, but also *forbids* him/her from doing so. This is for the purpose of security; we would not want the applications programmer to either inadvertently or maliciously trash someone else's disk file, for instance.

How the OS does these things is explained in the following sections. We will model the discussion after a UNIX system, but the description here applies to most modern OSs. It is assumed here that the reader is familiar with basic UNIX commands; a Unix tutorial is available at <http://heather.cs.ucdavis.edu/~matloff/unix.html>

### 1.3 A Bit More on System Calls

Recall that the OS makes available to application programs services such as I/O.<sup>4</sup> When you call `printf()`, for instance, it is just in the C library, not the OS, but it in turn calls `write()`, which *is* in the OS. The call to `write()` (which your program could also make directly) is a system call.

Most modern CPUs run in two or more **privilege levels**. As noted earlier, we for example would not want to give ordinary application programs direct access to I/O devices, e.g. disk drives, for security reasons. Thus the CPU is designed so that certain instructions, for example those which perform I/O, can be executed only at higher privilege levels, say Kernel Mode. (The term **kernel** refers to the OS.)

For this reason, one usually cannot implement a system call using an ordinary subroutine `CALL` instruction, because we need to have a mechanism that will change the machine to Kernel Mode. (Clearly, we cannot just have an instruction to do this, since ordinary user programs could execute this instruction and thus get into Kernel Mode themselves, wreaking all kinds of havoc!) Another problem is that the linker will not know where in the OS the desired subroutine resides.

Instead, system calls are implemented via an instruction type which is called a **software interrupt**. On Intel machines, this takes the form of the `INT` instruction, which has one operand.

We will assume Linux in the remainder of this subsection, and the operand is `0x80`. In other words, the call to `write()` in your C program (or in `printf()`) will be translated to

---

<sup>3</sup>For example, Pentium CPUs do have such capabilities, but the Windows OS does not make use of them. Linux and Windows NT do make use of them.

<sup>4</sup>You should not jump to the conclusion that all, or almost all, of the services deal with I/O. For example, the `execve()` service is used by one program to start the execution of another. Another non-I/O example is `getpid()`, which will return the process number of the program which calls it.

```
(code to put parameters values into designated registers)
int 0x80
```

The INT instruction works like a hardware interrupt, in the sense that it will force a jump to the OS, and change the privilege level to Kernel Mode, enabling the OS to execute the privileged instructions it needs. You should keep in mind, though, that here the “interrupt” is caused deliberately by the program which gets “interrupted,” via an INT instruction. This is much different from the case of a hardware interrupt, which is an action totally unrelated to the program which is interrupted.

The operand, 0x80 above, is the analog of the device number in the case of hardware interrupts. The CPU will jump to the location indicated by the vector at  $c(IDT)+8*0x80$ .

When the OS is done, it will execute an IRET instruction to return to the application program which made the system call. The IRET also makes a change back to User Mode.

As indicated above, a system call generally has parameters, just as ordinary subroutine calls do. One parameter is common to all the services – the service number, which is passed to the OS via the EAX register. Other registers may be used too, depending on the service.

As an example, the following Intel Linux assembly language program writes the string “ABC” to the screen:

```
hi:    .string "ABC"

.globl _start
_start:

    movl $4, %eax    # the write() system call, number 4 obtained
                    # from /usr/include/asm/unistd.h
    movl $1, %ebx    # 1 = file handle for stdout
    movl $hi, %ecx   # write from where
    movl $3, %edx    # write how many bytes
    int $0x80        # system call
```

For this particular OS service, the parameters are passed in the registers EBX, ECX and EDX (and, as mentioned before, with EAX specifying which service we want).

## 1.4 Making These Concepts Concrete: Commands You Can Try Yourself

The UNIX **strace** command will report which system calls your program makes. Place it before your program name on the command line, as with the **time** command above.

## 2 System Bootup

As will be explained later, when we wish to run an application program, the OS loads the program into memory. But how does the OS itself get loaded into memory and begin execution? The process by which this is done is called **bootup**.

The CPU hardware will be designed so that upon powerup the Program Counter (PC) is initialized to some specific value, say 0xffffffff in the case of Intel CPUs. And those who **fabricate** the computer (i.e. who put together the CPU, memory, bus, etc. to form a complete system) will include a small ROM at that same address, again say 0xffffffff. The contents of the ROM are the **boot loader** program. So, immediately after powerup, the boot loader program is running!

The goal of the boot loader is to load in the OS from disk to memory. In the simple form, the boot loader reads a specified area of the disk, copies the contents there (which will be the OS) to memory, and then finally executes a JMP instruction (or equivalent) to that section of memory—so that now the OS is running. In a more complicated form, the boot loader only reads part of the OS into memory, and then performs a JMP instruction to the OS; the OS then reads the rest of itself into memory.

For the sake of concreteness, let's look more closely at the Intel case. The program in ROM here is the BIOS, the Basic I/O System. It contains parts of the device drivers for that machine,<sup>5</sup> and also contains the boot loader program.

The boot loader program has been written to read from the first sector of the disk.<sup>6</sup> That sector is called the Master Boot Record (MBR). The boot loader in ROM will copy the code from the MBR to memory, starting at location 0x00007c00. It then does a jump to that address, so that that code is now running.

Now if the machine had originally been shipped with Windows installed, the code in the MBR was written to then load the Windows OS into memory. If on the other hand the machine had been shipped with Windows NT, OS/2, Linux or some other OS installed, the code in the MBR would have been written accordingly, and that OS would now be loaded into memory.

Many people who use Linux retain both Windows and Linux on their hard drives, and have a **dual-boot** setup. They start with a Windows machine, but then install Linux as well. As part of the process of installing Linux, a program named LILO (Linux Loader) is written into the MBR. So, when the boot loader in ROM loads the code from the MBR into memory and jumps to that part of memory, LILO will now run. LILO will then ask the user whether he/she wants to boot Linux or Windows, and then load the requested system into memory.

In any case, after the OS is loaded into memory by the code in the MBR, that code will perform a jump to the OS, so the OS is running. Typically, not all of the OS will be in memory yet, so the OS will now read the rest of itself into memory.

---

<sup>5</sup>These may or may not be used, depending on the OS. Windows uses them, but Linux doesn't.

<sup>6</sup>Disks are divided in blocks called **sectors**. The boot device could be something else instead of a hard drive, such as a floppy or a CD-ROM. This is set in the BIOS, with a priority ordering of which device to try to boot from first.

### 3 Application Program Loading

Suppose you have just compiled a program, producing, say for a Linux system, an executable file **a.out**. To run it, you type

```
% a.out
```

For the time being, assume a very simple machine/OS combination, with no **virtual memory**. Here is what will occur:

- During your typing of the command, the shell is running, say **tcsh** or **bash**. Again, the shell is just a program (one which could be assigned as a homework problem in a course), and it reads the command you type.
- The shell will then make a system call, `execve()`,<sup>7</sup> asking the OS to run **a.out**. The OS is now running.
- The OS will look in its disk directory, to determine where on disk the file **a.out** is. It will read the beginning section of **a.out**, which contains information on the size of the program, a list of the data segments used by the program, and so on.
- The OS will check its memory-allocation table (this is just an array in the OS) to find an unused region or regions of memory large enough for **a.out**. (Note that the OS has loaded all currently-active programs in memory up to now, just like it is loading **a.out** now, so it knows exactly which parts of memory are in use and which are free.) We need space both for **a.out**'s instructions (called the **text** portion of the program in UNIX terminology), and data – static data items (scalar and array variables, called the **data** segment in UNIX), as well as for stack space and the **heap** (used for calls to `calloc()` and `malloc()`).
- The OS will then load **a.out** (including text and data) into those regions of memory, and will update its memory-allocation table accordingly. Also, the OS will create a **page table** for **a.out**, which we will describe later.
- The OS will check a certain section of the **a.out** file, in which the linker previously recorded **a.out**'s **entry point**, i.e. the instruction in **a.out** at which execution is to begin.
- The OS is now ready to initiate execution of **a.out**. It will set the stack pointer to point to the place it had chosen earlier for **a.out**'s stack. (The OS will save its own register values, including stack pointer value, beforehand.) Then it will actually start **a.out**, say by executing a **JMP** instruction (or equivalent) to jump to **a.out**'s entry point.
- The **a.out** program is now running!

---

<sup>7</sup>It will also call another system call, `fork()`, but we will not go into that here.

Note that **a.out** itself may call `execve()` and start other programs running. For example, **gcc** does this.<sup>8</sup> It first runs the **cpp** C preprocessor (which translates `#include`, `#define` etc. from your C source file). Then it runs **cc1**, which is the “real” compiler (**gcc** itself is just a manager that runs the various components, as you can see now); this produces an assembly language file.<sup>9</sup> Then **gcc** runs **as**, the assembler, to produce a .o machine code file. The latter must be linked with some code, e.g. `/usr/lib/crt1.o`, which sets up the `main()` structure, e.g. access to the “argv” command-line arguments, and to the C library, e.g. `/lib/libc.so.6`; so, **gcc** runs the linker, **ld**. In all these cases, **gcc** starts these programs by calling `execve()`.

## 4 Timesharing

### 4.1 Many Processes, Taking Turns

**Timesharing** involves having several programs (or even several instances of the same program) running in what appears to be a simultaneous manner. Since the system has only one CPU (we will exclude the case of multiprocessor systems in this discussion), then this simultaneity is of course only an illusion, since only one program can run at any given time, but it is a worthwhile illusion, as we will see.

First of all, how is this illusion attained? The answer is that we have the programs all take turns running, with each turn – called a **quantum** or **timeslice** – being of very short duration, for example 50 milliseconds. Say we have four programs, u, v, x and y, running currently. What will happen is that first u runs for 50 milliseconds, then u is suspended and v runs for 50 milliseconds, then v is suspended and x runs for 50 milliseconds, and so on (after y gets its turn, then u gets a second turn, etc.). Since the turn-switching (formally known as **context-switching**) is happening so fast (every 50 milliseconds), it appears to us humans that each program is running continuously (though at one-fourth speed), rather than on and off, on and off, etc.

But how can the OS enforce these quanta? For example, how can the OS force the program u above to stop after 50 milliseconds? As discussed earlier, the answer is, “It can’t! The OS is dead while u is running.” Instead, the turns are implemented via a timing device, which emits a hardware interrupt at the proper time. For example, we could set the timer to emit an interrupt every 50 milliseconds. We would write a timer device driver, and incorporate it into the OS.<sup>10</sup>

The timer device driver saves all u’s current register values, including its PC value and the value in its Processor Status Register. Later, when u’s next turn comes, those values will be restored, and u will resume execution as if nothing ever happened. For now, though, the OS routine will restore v’s previously-saved register values, making sure to restore the PC value last of all. That last action forces a jump from the OS to v, right at the spot in v where v was suspended at the end of its last quantum. (Again, the CPU just “minds

---

<sup>8</sup>Remember, a compiler is just a program too. It is a very large and complex program, but in principle no different from the programs you write.

<sup>9</sup>You can view this file by running **gcc** with its `-S` option. This is often handy if you are going to write an assembly-language subroutine to be called by your C code, so that you can see how the compiler will deal with the parameters in the call to the subroutine.

<sup>10</sup>We will make such an assumption here. However, what is more common is to have the timer interrupt more frequently than the desired quantum size. The timer is set to interrupt every 10 milliseconds. If we want a quantum size of 50, we write the OS so that a program’s turn is ended after the fifth interrupt.



its own business,” and does not “know” that one program, the OS, has handed over control to another, *v*; the CPU just keeps performing its fetch/execute cycle, fetching whatever the PC points to.)

At any given time, there are many different **processes** in memory. These are instances of executions of programs. If for instance there are three users running the **gcc** C compiler right now on a given machine, here one program corresponds to three processes.

## 4.2 Example of OS Code: Linux for Intel CPUs

Here is a bit about how the context switch is done in Linux, in the version for Intel machines. Each process has a section of memory, called the Task State Segment (TSS), which stores various pieces of information about that process, such as the register values which the program had at the time its last turn ended. In the code below registers EBX and ECX point to the TSSs of the process whose turn just ended, say *u*, and the process to which we will give the next turn, say *v*.

As an example of the operations performed, and to show you concretely that the OS is indeed really a program with real code, here is a typical excerpt of code:<sup>11</sup>

```
pushl %esi
pushl %edi
pushl %ebp
movl %esp, 532(%ebx) ...
movl 532(%ecx), %esp
...
popl %ebp
popl %edi
popl %esi
...
iret
```

Here is what that code does. In the Linux source code, the TSS is accessed as a C **struct**, which has various fields to store things like register values. For example, `tss.esp` contains the previously-stored value of ESP, the stack pointer; this field happens to be located 532 bytes past the beginning of the TSS.

Now, upon entry to the above OS code, ESP is still pointing to *u*’s stack, so the three PUSH instructions save *u*’s values of the ESI, EDI and EBP registers on *u*’s own stack.<sup>12</sup> The other register values of *u* must be saved too, including its value of ESP. The latter is done by the MOV operation (“movl” in AT&T assembly language syntax), which copies the current ESP value, i.e. *u*’s ESP value, to `tss.esp` in *u*’s TSS. Other register saving is similar, though not shown here.

Now the OS must prepare to start *v*’s next turn. Thus *v*’s previously-saved register values must be restored to the registers. To understand how that is done, you must keep in mind that that same code above had been

<sup>11</sup>I have slightly modified some of this for the sake of simplicity.

<sup>12</sup>Note the need to write this in assembly language instead of C, since C would not give us direct access to the registers or the stack. Most of Linux is written in C, but machine-dependent operations like the one here must be done in assembly language.

executed when v's last turn ended. Thus v's value of ESP is in tss.esp of its TSS, and the second MOV we see above copies that value to ESP. So, now we are using v's stack.

Next, note similarly that at the end of v's last turn, its values of ESI, EDI and EBP were pushed onto its stack, and of course they are still there. So, we just pop them off, and back into the registers, which is what those three POP instructions do.

Finally, note that the mechanism which made v's last turn end was a hardware interrupt from the timer. At that time, the values of the Flags Register, CS and PC were pushed onto the stack. An IRET instruction pops all that stuff back into the corresponding registers. Note only does that restore registers, but since v's old PC value is restored, v is now running!

### 4.3 Process States

The OS maintains a **process table** which shows the state of each process in memory, mainly Run state versus Sleep state. A process which is in Run state means that it is ready to run but simply waiting for its next turn. The OS will repeatedly cycle through the process table, starting turns for processes in Run state but skipping over those in Sleep state. The processes in Sleep state are waiting for something, typically an I/O operation, and thus currently ineligible for turns. So, each time a turn ends, the OS will browse through its process table, looking for a process in Run state, and then choosing one for its next turn.

Say our application program u above contains a call to scanf() to read from the keyboard. Recall that scanf() calls the OS function read(). The latter will check to see whether there are any characters ready in the keyboard buffer. Typically there won't be any characters there yet, because the user has not started typing yet. In this case the OS will place this process in Sleep state, and then start a turn for another process.

How does a process get switched to run state from Sleep state? Say our application program u was in Sleep state because it was waiting for user input from the keyboard (say it was waiting for just a single character). As explained earlier, when the user hits a key, that causes a hardware interrupt from the keyboard, which forces a jump to the OS. Suppose at that time program v happened to be in the midst of a quantum. The CPU would temporarily suspend v and jump to the keyboard driver in the OS. The latter would notice that the program u had been in Sleep state, waiting for keyboard input, and would now move u to Run state.

Note, though, that that does not mean that the OS now starts u's next turn; u simply becomes eligible to run. Recall that each time one process' turn ends, the OS will select another process to run, from the set of all processes currently in Run state, and u will now be in that set.

### 4.4 How the Process Tree Is Built

On UNIX systems, the first thing the OS does after bootup is to start a process named **init**. That process then starts all other OS processes, such as the one handling user logins, **in.logind**. When a user logs in, the latter starts up a shell for the user, say **tcsh**, which in turn will start whatever processes the user commands, say **a.out**. As we have seen before, those may in turn spawn further processes.

## 4.5 Making These Concepts Concrete: Commands You Can Try Yourself

First, try the **ps** command. On UNIX systems, much information on current processes is given by the **ps** command, including:

- state (Run, Sleep, etc.)
- page usage (how many pages, number of page faults, etc.)
- ancestry (which process is the “parent” of the given process)

**The reader is urged to try this out. You will understand the concepts presented here much better after seeing some concrete information which the ps command can give you.** The format of this command’s options varies somewhat from machine to machine (on Linux, I recommend “ps ax”), so check the man page for details, but run it with enough options that you get the fullest output possible.

Another command to try is **w**. One of the pieces of information given by the **w** command is the average number of processes in Run state in the last few minutes. The larger this number is, the slower will be the response time of the machine as perceived by a user, as his program is now taking turns together with more programs run by other people.

On Linux (and some other UNIXx) systems, you can also try the **pstree** command, which graphically shows the “family tree” (ancestry relations) of each process. For example, here is the output I got by running it on one of our CSIF PCs:

```
% pstree
init--+-atd
|   -crond
|   -gpm
|   -inetd---in.rlogind---tcsh---pstree
|   -kdm--+-X
|       '-kdm---wmaker--+-gnome-terminal--+-gnome-pty-helpe
|                       |   '-tcsh--+-netscape-commun---netscape-+
|                               |   '-vi
|                               |   -2*[gnome-terminal--+-gnome-pty-helpe]
|                               |   '-tcsh]
|                               |   -gnome-terminal--+-gnome-pty-helpe
|                               |   '-tcsh---vi
|                               '-wmclock
|   -kerneld
|   -kflushd
|   -klogd
|   -kswapd
|   -lpd
|   -6*[mingetty]
```

```

| -2*[netscape-commun---netscape-commun]
| -4*[nfsiod]
| -portmap
| -rpc.rusersd
| -rwhod
| -sendmail
| -sshd
| -syslogd
| -update
| -xconsole
| -xntpd
| -ypbind---ypbind
%

```

Note how, for example, that some user is running **vi**. He gave the **vi** command to **tcsh**, which started a **vi** process for him. The **tcsh** had in turn been started by **gnome-terminal**, a version of **xterm** running the window this user had been running the shell in.

## 5 Virtual Memory

### 5.1 Make Sure You Understand the Goals

Now let us add in the effect of virtual memory (VM). VM has the following basic goals:

- Overcome limitations on memory size:  
We want to be able to run a program, or collectively several programs, whose memory needs are larger than the amount of physical memory available
- Relieve the compiler and linker of having to know what memory is free when a program is run:  
We want to facilitate **relocation** of programs, meaning that the compiler and linker, do not have to worry about where in memory a program will be loaded when it is run.
- Enable security:  
We want to ensure that one program will not accidentally (or intentionally) destroy another program's operation by writing to the latter's area of memory
- Enable sharing:  
We want to be able to have only one copy in memory of the "text" portion of a large program (e.g. a compiler) even though several users are each running instances of the program

## 5.2 Example of Virtual Nature of Addresses

The word *virtual* means “apparent.” It will appear that a program resides entirely in main memory, when in fact only part of it is there; it will appear (e.g. from the compiler’s point of view) that the program is loaded starting at location 0 in memory, but that will not be the case in actuality.

To make this more concrete, suppose our C source file from which we compiled **a.out** included a statement

```
int x;
```

and suppose the compiler and linker had assigned the address 200 to x. In other words, a statement in our C source file like

```
printf ("%d", &x);
```

would print out the value 200.

Then **a.out** might have instructions like

```
movl (200), %eax
```

on an Intel machine. This instruction copies the contents of word 200 of memory to the CPU register EAX.

At the time **a.out** is loaded by the OS into memory, the OS will divide both the text (instructions) and data portions of **a.out** into chunks, and find unused places in memory at which to place these chunks. The chunks are called **pages** of the program, and the same-sized places in memory in which the OS puts them are called pages of memory. The OS sets up a **page table**, which is an array in memory which is maintained by the OS, in which the OS records the correspondences, i.e. lists which page of the program is stored in which page of memory.

So, what appears to be in word 200 in memory from the program code above may actually be in, say, word 1204. At the time the CPU executes that instruction, the CPU will determine where “word 200” really is by doing a lookup in the page table. In our example here, the table will show that the item we want is actually in word 1204, and the CPU will then read from that location.

In this example, we say the **virtual address** is 200, and the **physical address** is 1204.

## 5.3 Overview of How the Goals Are Achieved

Let’s look at our stated goals in Section 5.1 above, and see how they are achieved:

- Overcome limitations on memory size:

To conserve memory space, the OS will initially load only part of **a.out** into memory, with the remainder being left back on disk. The pages of the program which are not loaded will be marked in the page table as currently being nonresident, and their locations on disk will be shown in the table. During execution of the program, if the program needs one of the nonresident pages, the CPU will notice that the page is nonresident (this is called a “page fault”), and cause an internal interrupt. That cause a jump to the OS, which will bring in that page from disk, and then jump back to the program, which will resume at the instruction which accessed the missing page.

Note that pages will often go back and forth between disk and memory in this manner. Each time the program needs a missing page, that page is brought in from disk and a page which had been resident is written back to disk in order to make room for the new one.

A big issue is the algorithm the OS uses to decide which page to move back to disk whenever it brings a page from disk after a page fault. This is beyond the scope of this document here, but one point to notice is that the algorithm will be chosen so as to work well on “most” programs. For some programs, that algorithm will result in a lot of page faults, due to it frequently being the case that right after a page is replaced, that same page is needed again.

- Relieve the compiler and linker of the burden of knowing what memory is free at the time the program executes:

This is clear from the example above, where the location “200” which the compiler and linker set up for *x* was in effect changed by the OS to 1204 at the time the program was loaded. The OS recorded this in the page table, and then during execution of the program, the VM hardware in the CPU does lookups in the page table to get the correct addresses.

- Enable security:

The page table will consist of one entry per page. That entry will, as noted earlier, include information as to where in memory that page of the program currently resides, of if currently nonresident, where on disk the page is stored. But in addition, the entry will also list the permissions the program has to access this particular page – read, write, execute – in a manner analogous to file-access permissions. If the program tries to access a page for which it does not have the proper permission, the VM hardware in the CPU will cause an internal interrupt, causing the OS to run. The OS will then kill the process, i.e. remove it from the process table.

- Enable sharing:

Suppose for example two users of a given machine wish to run **gcc** right now. This is a huge program, so conserving memory is quite important. An obvious strategy on the part of the OS would be to load only one copy of the text (i.e. instructions) part of **gcc**; of course, there must be separate copies of the data sections, as the data (.c source code files) for the two users are different.

VM allows us to accomplish this. Each user’s page table will list the same entries for the text, and thus they will share the text part of the program.

## 5.4 Creation and Maintenance of the Page Table

Note carefully the roles of the players here: It is the software, the OS, that creates and maintains the page table, but it is the hardware that actually uses the page table to generate addresses, check page residency and

check security.

When the OS creates a new process, it must find chunks (pages) of memory into which it will load part or all of the given program. It will create a page table for this process, and record in the page table the locations of these chunks (as well as record the locations on disk of the chunks which it did not load into memory).

The hardware will have a special Page Table Register (PTR) to point to the page table of the current process. When the OS starts a turn for a process, it will restore the previously-saved value of the PTR, and thus this process' page table will now be in effect.

## 5.5 Details on Usage of the Page Table

### 5.5.1 Virtual-to-Physical Address Translation, Page Table Lookup

Whenever the running program generates an address – either the address of an instruction, as will be the case for an instruction fetch, or the address of data, as will be the case during the execution of some types of instructions – this address is only virtual. It must be translated to the physical address at which the requested item actually resides. The circuitry in the CPU is designed to do this translation by performing a lookup in the page table.

The address space is broken into pages. For convenience, say the page size is 4096 bytes. For any virtual address, the virtual page number is equal to the address divided by the page size, 4096, and its offset within that byte is the address mod 4096. Since  $4096 = 2^{12}$ , that means that in a 32-bit virtual address, the upper 20 bits form the page number, and the lower 12 bits form the offset.

Consider for example the Intel instruction

```
movl $3, (0x735bca62)
```

This would copy the constant 3 to location 0x735bca62 (193539426 base-10). That means virtual page number 0x735bc (472508 base-10), offset 0xa62 (2658 base-10) within that page. In other words, the first byte of the word we will write to is byte 2658 within page 472508 in the virtual address space.

Suppose the entries in our page table are 32 bits wide, i.e. one word per entry.<sup>13</sup> Let's label the bits of an entry 31 to 0, where bit 31 is in the most-significant (i.e. leftmost) position and bit 0 is in the least significant (i.e. rightmost) place. Suppose the format of an entry is as follows:

- bits 31-12: physical page number if resident, disk location if not
- bit 11: 1 if page is resident, 0 if not
- bit 10: 1 if have read permission, 0 if not
- bit 9: 1 if have write permission, 0 if not

---

<sup>13</sup>If we were to look at the source code for the OS, we would probably see that the page table is stored as a very long array of type **unsigned int**, with each array element being one page table entry.

- bit 8: 1 if have execute permission, 0 if not
- bits 7-0: other information, not discussed here

Now, here is what will happen when the CPU executes the instruction

```
movl $3, (0x735bca62)
```

above:

- The CPU, seeing that this is virtual page number 0x735bc, will go to get that entry in the page table, as follows. Suppose the contents of the PTR is 0x256a1000. Then the table entry of interest here is at location  $0x735bc * 4 + 0x256a1000 = 0x2586e6f0$ . The CPU will read from that location, getting, say, 0xc2248eac.
- The CPU looks at bits 11-8 of that entry, getting 0xe, finding that the page is resident and that the program has read and write permission but not execute permission. The permission requested was write, so this is OK.
- The CPU looks at bits 31-12, getting 0xc2248. The virtual offset, which we found earlier to be 0xa62, is always retained, so the CPU now knows that the physical address of the virtual location 0x735bca62 is 0xc2248a62. The CPU puts this in the Memory Address Register (MAR), puts 3 in the Memory Data Register (MDR), and asserts the Write line in the bus. This writes 3 to memory location 0xc2248a62, and we are done.

### 5.5.2 Page Faults

Suppose in our example above bit 11 of the page table entry had been 0, indicating that the requested page was not in memory. This event is known as a **page fault**. If that occurs, the CPU will perform an internal interrupt,<sup>14</sup> which will force a jump to the OS. The OS will first decide which currently-resident page to replace, then write that page back to disk.<sup>15</sup> The OS would then bring in the requested page from disk. The OS would then update two entries in the page table: (a) it would change the entry for the page which was replaced, changing bit 11 to indicating the page is not resident, and changing bits 31-12; and (b) the OS would update the page table to indicate that the new item is resident now in memory, and show where it resides.

Since accessing the disk is far, far slower than accessing memory, a program will run quite slowly if it has too many page faults. If for example your PC at home does not have enough memory, you will find that you often have to wait while a large application program is loading, during which time you can hear the disk drive doing a lot of work, as the OS ejects many currently-resident pages to bring in the new application.

<sup>14</sup>The CPU will also record the PC value of the instruction which caused the page fault, so that that instruction can be restarted after the page fault is processed. In Pentium CPUs, the CR2 register is used to store this PC value.

<sup>15</sup>While we will not assume so here, most OSs will do this write-back only if it is necessary. One of the bits in our field of bits 7-0 above would be used as the **dirty bit** for this purpose. We will not pursue this aspect here.



### 5.5.3 Access Violations

If on the other hand an access violation occurs, the OS will announce an error – in UNIX, referred to as a **segmentation fault** – and kill the process, i.e. remove it from the process table.

For example, considering the following code:

```
int q[200];

main()
{
    int i;

    for (i = 0; i < 2000; i++) ~ {
        q[i] = i;
    }
}
```

Notice that the programmer has apparently made an error in the loop, setting up 2000 iterations instead of 200. The C compiler will not catch this at compile time, nor will the machine code generated by the compiler check that the array index is out of bounds at execution time.

If this program is run on a non-VM platform,<sup>16</sup> then it will merrily execute without any apparent error. It will simply write to the 180 words which follow the end of the array `q`. This may or may not be harmful, depending on what those words had been used for.

But on a VM platform, in our case UNIX, an error will indeed be reported, with a “Segmentation fault” message. However, as we look into how this comes about, the timing of the error may surprise you. The error is not likely to occur when `i = 200`; it is likely to be much later than that.

To illustrate this, I ran this program on a PC under Linux. I first added some code to aid in investigating what exactly occurs:

```
int q[200];

main()
{
    int i;

    printf( "%x %x\n", &q[0], &q[199])
```

---

<sup>16</sup>Recall that “VM platform” requires both that our CPU has VM capability, and that our OS uses this capability.

```

for (i = 0; i < 2000; i++)~ {
    printf("%d\n", i);
    q[i] = i;
}
}

```

After running this program, I found that the seg fault occurs not at  $i = 200$ , but actually at  $i = 616$ . Let's see why.

The first call to `printf()` reveals that `q` begins at virtual address `0x8049640` (and ends at `0x804995c`). On Intel machines, the page size is 4096 bytes, so a virtual address breaks down into a 20-bit page number and a 12-bit offset, just as in Section 5.5.1 above. In our case here, `q` begins in virtual page number `0x8049`, offset `0x640`. Converting these to base 10 for convenience, we find that `q` begins in virtual page number 32841, at byte 1600 of that byte.

Now it is the latter number which is of interest here. Remember, the page size is 4096 bytes, and we now know that `q` begins at the 1601<sup>st</sup> byte of a page.<sup>17</sup> Then `q` ends 200 words (800 bytes) later, at the 2400<sup>th</sup> byte of the page. *Now, here is the point: The first few nonexistent “elements” of `q` which follow the end of `q` are still in that same page – and since the program has write permission for the entire page, there will be no seg fault for trying to write to “`q[200]`”, “`q[201]`”, and so on.*

Eventually, for large enough  $i$ , “`q[i]`” will not be in that page, and its virtual page number will correspond to an entry in the page table which does not have write permission. Then the page fault will occur. It turned out to be for  $i = 616$ .

## 5.6 Improving Performance

Virtual memory comes at a big cost, in the form of overhead incurred by accessing the page tables. For this reason, the hardware will also typically include a **translation lookaside buffer** (TLB). This is a special cache to keep a copy of part of the page table in the CPU, to reduce the number of times one must access memory, where the page table resides.

## 5.7 Intel Page Tables

On Intel machines, each process actually has many tables, not just one. And some of them may even be nonresident currently; in other words, even the page tables are paged!

There is a “page table table,” which serves as directory of page tables for the current process. A special register, CR3, points to the page table table.

A virtual address is broken down into three fields, rather than the two in our descriptions above:

---

<sup>17</sup>Recall that the offset is retained after the virtual-to-physical address translation. So, even though we do not know what physical page number `q` was in when I ran the program, it does not matter in terms of offset. No matter which physical page it was, we know that `q` began at the 1601<sup>st</sup> byte in that page.

- bits 31-22:  $i$ , the index into the page table
- bits 21-12:  $j$ , index into some page table
- bits 11-0: the offset

So, in processing a virtual address, the CPU will first fetch from the address  $c+4*i$ , where  $c$  is the current contents of CR3. That fetched word  $w$  will tell the CPU whether the page table for this virtual address is currently resident in memory. If that table is resident,  $w$  will show the address of the table,  $d$ , while if it is nonresident, it will show where on disk to get it. The CPU will bring in the table from disk in the latter case. Once the CPU has  $d$ , it will then fetch from the address  $d+4*j$ , yielding the page table entry for the item our running program has requested. At this point, the operations are the same as described earlier.

## 5.8 Making These Concepts Concrete: Commands You Can Try Yourself

The UNIX **time** command will report how much time your program takes to run, how many page faults it generated, etc. Place it just before your program's name on the command line. (This program could be either one you wrote, or something like, say, **gcc**.) For example, if you have a program **x** with argument 12, type

```
time x 12
```

instead of

```
x 12
```

Also, the **top** program is very good, displaying lots of good information on the memory usage of each process.

# A Hardware Interrupts

## A.1 General Operation

A hardware **interrupt** is an electrical signal sent from an I/O device along an interrupt-request line IRQ in the system bus to the CPU (the line will be connected to an IRQ pin in the CPU). The CPU is designed to finish whatever instruction  $i$  and then jump to another part of memory where an **interrupt service routine** (ISR), or **device driver**, has been stored. Keep in mind that the device drivers are part of the OS.

The ISR now runs (first saving on the stack the values of any registers it will use, as any subroutine would do), performing whatever actions are needed to read or write the I/O data. When it is done, the ISR pops

from the stack any previously-saved register values, and executes an IRET (“interrupt return”) instruction. The latter is similar to an ordinary RET instruction used to return from a subroutine call, but in this case we are returning to whatever program had been running at the time the CPU received the interrupt signal. Since the hardware had saved the program state at that time, which is now restored by the hardware, the other program now resumes execution at the point at which it had been interrupted.

Most machines use **vector**ed interrupts. This means that there is a table stored in memory (initialized by the OS upon bootup) which stores one “interrupt vector” for each possible interrupting I/O device. When an interrupt is received from device  $i$ , the CPU will look up the vector for the  $i^{th}$  device. The vector will state the location of the driver for that device, and the CPU will then do a jump to that location, causing the driver to begin executing.

But how does the CPU know which device caused the interrupt? One way of handling this would be to have a different IRQ bus line for each I/O device. Some systems do use this approach, but generally it is infeasible, as a typical machine has many I/O devices, and we do not want to design the CPU to have so many pins. (Pins take up precious space on the periphery of a chip.)

## A.2 Some Details for Intel CPUs and PCs

Linux and similar systems run the Intel CPU in **protected mode**, which enables the hardware to provide various types of security features needed by a modern OS, such as virtual memory. Following is a simplified description of how these systems work.

Again for security reasons, we want I/O to be performed only by the OS. The Intel CPU has several modes of operation, which we will simplify here to just User Mode and Kernel (i.e. OS) Mode. The hardware is set up so that I/O instructions such as IN and OUT can be done only in Kernel Mode.

The Intel CPU contains an Interrupt Descriptor Table, IDT, which points to the beginning of the interrupt vector table in memory. Each vector is 8 bytes long, so the vector for the  $i^{th}$  I/O device is located at  $c(IDT)+8*i$ , where  $c()$  means “contents of.”

A PC also includes another piece of Intel hardware, the 8259A interrupt controller. The I/O devices are actually connected to the 8259A, which in turn is connected to the IRQ line, instead of the devices being connected directly to the IRQ. The 8259A has many input pins, one for each I/O device.<sup>18</sup>

When an interrupt from device  $i$  occurs, the 8259A will record the value of  $i$ , and then assert the IRQ line in the bus. The Intel CPU will then push the values of the Flags, CS and PC registers onto the stack,<sup>19</sup> and then proceed to determine the value of  $i$ . It does this by asserting the INTA (“interrupt acknowledge”) line in the bus. The 8259A then sends the value of  $i$  along the data bus, where it is received by the CPU. The CPU then does the lookup in the vector table, finds the information for device  $i$ , and jumps to the driver for device  $i$ . The information in the vector for device  $i$  will also result in the CPU changing to Kernel Mode, which is important since the driver must execute I/O instructions such as IN and OUT.

The driver will then execute, reading and/or writing data at the ports of the I/O device. When the driver is done, it executes IRET, everything is restored, and the interrupted program resumes execution where it left

<sup>18</sup>If we run out of such pins, two or more 8259A devices may be **cascaded** together.

<sup>19</sup>Don’t worry about the CS register.

off, as if nothing had ever happened.

If several I/O devices cause interrupts at about the same time, the 8259A can “queue” them, so that all will be processed.