

# 11

## CASE STUDY 2: WINDOWS 2000

Windows 2000 is a modern operating system that runs on high-end desktop PCs and servers. In this chapter we will examine various aspects of it, starting with a brief history, then moving on to its architecture. After this we will look at processes, memory management, I/O, the file system, and finally, security. One thing we will not look at is networking as that could easily fill a chapter (or an entire book) all by itself.

### 11.1 HISTORY OF WINDOWS 2000

Microsoft operating systems for desktop and laptop PCs can be divided into three families: MS-DOS, Consumer Windows (Windows 95/98/Me), and Windows NT. Below we will briefly sketch each of these families.

#### 11.1.1 MS-DOS

In 1981, IBM, at the time the biggest and most powerful computer company in the world, produced the 8088-based IBM PC. The PC came equipped with a 16-bit real-mode, single-user, command-line oriented operating system called MS-DOS 1.0. The operating system was provided by Microsoft, a tiny startup, mostly known at that time for its BASIC interpreter used on 8080 and Z-80 systems. This operating system consisted of 8 KB of memory resident code and was

closely modeled on CP/M, a tiny operating system for the 8-bit 8080 and Z80 CPUs. Two years later, a much more powerful 24-KB operating system, MS-DOS 2.0, was released. It contained a command line processor (shell), with a number of features borrowed from UNIX.

When Intel came out with the 286 chip, IBM built a new computer around it, the PC/AT, released in 1986. AT stood for Advanced Technology, because the 286 ran at a then-impressive 8 MHz and could address—with great difficulty—all of 16 MB of RAM. In practice, most systems had at most 1 MB or 2 MB, due to the great expense of so much memory. The PC/AT came equipped with Microsoft's MS-DOS 3.0, by now 36 KB. Over the years, MS-DOS continued to acquire new features, but it was still a command-line oriented system.

### 11.1.2 Windows 95/98/Me

Inspired by the user interface of the Apple Lisa, the forerunner to the Apple Macintosh, Microsoft decided to give MS-DOS a graphical user interface (shell) that it called **Windows**. Windows 1.0, released in 1985, was something of a dud. Windows 2.0, designed for the PC-AT and released in 1987, was not much better. Finally, Windows 3.0 for the 386 (released in 1990), and especially its successors 3.1 and 3.11, caught on and were huge commercial successes. None of these early versions of Windows were true operating systems, but more like graphical user interfaces on top of MS-DOS, which was still in control of the machine and the file system. All programs ran in the same address space and a bug in any one of them could bring the whole system to a grinding halt.

The release of Windows 95 in August 1995 still did not completely eliminate MS-DOS, although it transferred nearly all the features from the MS-DOS part to the Windows part. Together, Windows 95 and the new MS-DOS 7.0 contained most of the features of a full-blown operating system, including virtual memory, process management, and multiprogramming. However, Windows 95 was not a full 32-bit program. It contained large chunks of old 16-bit assembly code (as well as some 32-bit code) and still used the MS-DOS file system, with nearly all its limitations. The only major change to the file system was the addition of long file names in place of the 8 + 3 character file names allowed in MS-DOS.

Even with the release of Windows 98 in June 1998, MS-DOS was still there (now called version 7.1) and running 16-bit code. Although yet more functionality migrated from the MS-DOS part to the Windows part, and a disk layout suitable for larger disks was now standard, under the hood, Windows 98 was not very different from Windows 95. The main difference was the user interface, which integrated the desktop and the Internet more closely. It was precisely this integration that attracted the attention of the U.S. Dept. of Justice, which then sued Microsoft claiming that it was an illegal monopoly, an accusation Microsoft vigorously denied. In April 2000, a U.S. Federal court agreed with the government.

In addition to containing a large lump of old 16-bit assembly code in the kernel, Windows 98 had two other serious problems. First, although it was a multiprogramming system, the kernel itself was not reentrant. If a process was busy manipulating some kernel data structure and then suddenly its quantum ran out and another process started running, the new process might find the data structure in an inconsistent state. To prevent this type of problem, after entering the kernel, most processes first acquired a giant mutex covering the whole system before doing anything. While this approach eliminated potential inconsistencies, it also eliminated much of the value of multiprogramming since processes were frequently forced to wait for unrelated processes to leave the kernel before they could enter it.

Second, each Windows 98 process had a 4-GB virtual address space. Of this, 2 GB was completely private to the process. However, the next 1 GB was shared (writably) among all other processes in the system. The bottom 1 MB was also shared among all processes to allow all of them to access the MS-DOS interrupt vectors. This sharing facility was heavily used by most Windows 98 applications. As a consequence, a bug in one program could wipe out key data structures used by unrelated processes, leading to them all crashing. Worse yet, the last 1 GB was shared (writably) with the kernel and contained some critical kernel data structures. Any rogue program that overwrote these data structures with garbage could bring down the system. The obvious solution of not putting kernel data structures in user space was not possible because this feature was essential to making old MS-DOS programs work under Windows 98.

In the millennium year, 2000, Microsoft brought out a minor revision to Windows 98 called **Windows Me (Windows Millennium Edition)**. Although it fixed a few bugs and added a few features, under the covers it is essentially Windows 98. The new features included better ways to catalog and share images, music, and movies, more support for home networking and multiuser games, and more Internet-related features, such as support for instant messaging and broadband connections (cable modems and ADSL). One interesting new feature was the ability to restore the computer to its previous settings after a misconfiguration. If a user reconfigures the system (e.g., changing the screen from  $640 \times 480$  to  $1024 \times 768$ ) and it no longer works, this feature makes it possible to revert to the last known working configuration.

### 11.1.3 Windows NT

By the late 1980s, Microsoft realized that building a modern 32-bit operating system on top of the leaky 16-bit MS-DOS probably was not the best way to go. It recruited David Cutler, one of the key designers of DEC's VMS operating system, to work for Microsoft and gave him the job of leading a team to produce a brand-new 32-bit Windows compatible operating system from the ground up. This new

system, later called **Windows NT (Windows New Technology)**, was intended for mission-critical business applications as well as for home users. At the time, mainframes still ruled the (business) world, so designing an operating system on the assumption that companies would use personal computers for anything important was a visionary goal, but one that history has shown to be a very good one. Features such as security and high reliability, clearly lacking on the MS-DOS-based versions of Windows, were high on the agenda for (Windows) NT. Cutler's background with VMS clearly shows in various places, with there being more than a passing similarity between the design of NT and that of VMS.

The project succeeded and the first version, called Windows NT 3.1, was released in 1993. This initial release number was chosen to match the number of Microsoft's then-popular 16-bit Windows 3.1 system. Microsoft expected that NT would rapidly replace Windows 3.1 because it was technically a far superior system.

Much to its surprise, nearly all users preferred to stick with the old 16-bit system they knew, rather than upgrade to an unknown 32-bit system they did not know, however better it might be. Furthermore, NT required far more memory than Windows 3.1 and there were no 32-bit programs for it to run, so why bother? The failure of NT 3.1 to catch on in the marketplace was the reason Microsoft decided to build a 32-bit-ish version of Windows 3.1, namely Windows 95. The continued user resistance to NT then caused Microsoft to produce Windows 98 and finally Windows Me; each one claimed to be the very last release of the MS-DOS-based systems.

Despite the fact that nearly all consumers and most businesses ignored NT 3.1 for desktop systems, it did acquire a small following in the server market. A few new 3.x releases with small changes occurred in 1994 and 1995. These slowly began to acquire more following among desktop users as well.

The first major upgrade to NT came with NT 4.0 in 1996. This system had the power, security, and reliability of the new operating system, but also sported the same user interface as the by-then very popular Windows 95. This compatibility made it much easier for users to migrate from Windows 95 to NT, and many of them did so. Some of the differences between Windows 95/98 and Windows NT are summarized in Fig. 11-1.

From the beginning, NT was designed to be portable, so it was written almost entirely in C, with only a tiny bit of assembly code for low-level functions such as interrupt handling. The initial release consisted of 3.1 million lines of C for the operating system, libraries, and the environment subsystems (discussed below). When NT 4.0 came out, the code base had grown to 16 million lines of code, still mostly C, but with a small amount of C++ in the user interface part. By this time the system was highly portable, with versions running on the Pentium, Alpha, MIPS, and PowerPC, among other CPUs. Some of these have been dropped since then. The story of how NT was developed is given in the book *Showstopper* (Zachary, 1994). The book also tells a lot about the key people involved.

Item	Windows 95/98	Windows NT
Full 32-bit system?	No	Yes
Security?	No	Yes
Protected file mappings?	No	Yes
Private addr space for each MS-DOS prog?	No	Yes
Unicode?	No	Yes
Runs on	Intel 80x86	80x86, Alpha, MIPS, ...
Multiprocessor support?	No	Yes
Re-entrant code inside OS?	No	Yes
Plug and play?	Yes	No
Power management?	Yes	No
FAT-32 file system?	Yes	Optional
NTFS file system	No	Yes
Win32 API?	Yes	Yes
Run all old MS-DOS programs?	Yes	No
Some critical OS data writable by user?	Yes	No

**Figure 11-1.** Some differences between Windows 98 and Windows NT.

### 11.1.4 Windows 2000

The release of NT following NT 4.0 was originally going to be called NT 5.0. However, in 1999, Microsoft changed the name to Windows 2000, mostly in an attempt to have a neutral name that both Windows 98 users and NT users could see as a logical next step for them. To the extent that this approach succeeds, Microsoft will have a single main operating system built on reliable 32-bit technology but using the popular Windows 98 user interface.

Since Windows 2000 really is NT 5.0, it inherits many properties from NT 4.0. It is a true 32-bit (soon to be 64-bit) multiprogramming system with individually protected processes. Each process has a private 32-bit (soon 64-bit) demand-paged virtual address space. The operating system runs in kernel mode, whereas user processes run in user mode, providing complete protection (with none of the protection flaws of Windows 98). Processes can have one or more threads, which are visible to, and scheduled by, the operating system. It has Dept. of Defense C2 security for all files, directories, processes, and other shareable objects (at least, if the floppy disk is removed and the network is unplugged). Finally, it also has full support for running on symmetric multiprocessors with up to 32 CPUs.

The fact that Windows 2000 really is NT 5.0 is visible in many places. For example, the system directory is called `\winnt` and the operating system binary (in `\winnt\system32`) is called `ntoskrnl.exe`. Right clicking on this file to examine its

properties shows that its version number is 5.xxx.yyy.zzz, where the 5 stands for NT 5, xxx is the release number, yyy is the build (compilation) number, and zzz the minor variant. Also, many of the files in `\winnt` and its subdirectories have *nt* in their names, such as *ntvdm*, NT's virtual MS-DOS emulator.

Windows 2000 is more than just a better NT 4.0 with the Windows 98 user interface. To start with, it contains a number of other features previously found only in Windows 98. These include complete support for plug-and-play devices, the USB bus, IEEE 1394 (FireWire), IrDA (the infrared link between portable computers and printers), and power management, among others. In addition, a number of new features not present in any other Microsoft operating system have been added, including active directory directory service, security using Kerberos, support for smart cards, system monitoring tools, better integration of laptop computers with desktop computers, a system management infrastructure, and job objects. Also, the main file system, NTFS, has been extended to support encrypted files, quotas, linked files, mounted volumes, and content indexing, for example. Another novel NTFS feature is the single instance store, which is a kind of copy-on-write link in which two users can share a linked file until one of them writes on it, at which time a copy is made automatically.

One other major improvement is internationalization. NT 4.0 came in separate versions for different languages with the text strings embedded in the code. Installing an English software package on a Dutch computer often caused parts of the operating system to stop using Dutch and start using English because certain files containing code and text strings were overwritten. This problem has been eliminated. Windows 2000 has a single binary that runs everywhere in the world. An installation, or even an individual user, can choose the language to use at run time because all the menu items, dialog strings, error reports, and other text strings have been removed from the operating system and put in separate directories, one per installed language. Like all previous versions of NT, Windows 2000 uses Unicode throughout the system to support languages not using the Latin alphabet, such as Russian, Greek, Hebrew, and Japanese.

One thing that Windows 2000 does not have is MS-DOS. It is simply not there in any form (nor was it there in NT). There is a command line interface, but this is a new 32-bit program that includes the old MS-DOS functionality and considerably new functionality as well.

Despite many portability features with regard to the code, hardware, language, etc., in one respect Windows 2000 is less portable than NT 4.0: it runs on only two platforms, the Pentium and the Intel IA-64. Originally NT supported additional platforms, including the PowerPC, MIPS, and Alpha, but over the years, Microsoft dropped one after another for commercial reasons.

Like previous versions of NT, Windows 2000 comes in several product levels, this time: Professional, Server, Advanced server, and Datacenter server. The differences between all these versions are minor however, with the same executable binary used for all versions. When the system is installed, the product type is

recorded in an internal database (the registry). At boot time, the operating system checks the registry to see which version it is. The differences are shown in Fig. 11-2.

Version	Max RAM	CPUs	Max clients	Cluster size	Optimized for
Professional	4 GB	2	10	0	Response time
Server	4 GB	4	Unlimited	0	Throughput
Advanced server	8 GB	8	Unlimited	2	Throughput
Datacenter server	64 GB	32	Unlimited	4	Throughput

Figure 11-2. The different versions of Windows 2000.

As can be seen from the figure, the differences include the maximum memory supported, the maximum number of CPUs (for a multiprocessor configuration), and the maximum number of clients that can be served. The cluster size relates to the ability of Windows 2000 to make two or four machines look like a single server to the outside world, a useful feature for Web servers, for example. Finally, the default parameters are tuned differently on Professional, to favor interactive programs over batch work, although these can easily be changed if desired. One last difference is that some extra software is provided on the servers and some extra tools are provided on Datacenter server for managing large jobs.

The reason for having multiple versions is simply marketing: this allows Microsoft to charge big companies more than they charge individuals for what is essentially the same product. This idea is not new, however, and hardly unique to Microsoft. For years, airlines have been charging business passengers much more, not only for Business Class, but also for Cattle Class if they want the luxury of buying the ticket a day before the flight instead of a month before the flight.

Technically, the way the version differences are maintained is that in a few places in the code, two variables are read from the registry, *ProductType* and *ProductSuite*. Depending on their values, slightly different code is executed. Changing these variables is in violation of the license. In addition, the system traps any attempt to change them and records the attempt at tampering in an indelible way so it can be detected later.

In addition to the basic operating system, Microsoft has also developed several tool kits for advanced users. These include the Support Tools, the Software Development Kit, the Driver Development Kit, and the Resource Kit. These include a large number of utilities and tools for tweaking and monitoring the system. The support tools are on the Windows 2000 CD-ROM in the directory `\support\tools`. The standard installation procedure does not install them, but there is a file *setup.exe* in that directory that does. The SDK and DDK are available to developers at [msdn.microsoft.com](http://msdn.microsoft.com). The Resource Kit is a Microsoft product in a box. There are also various third-party tools available for snooping on the Windows 2000 internals, including a nice set available for free at the Web site

[www.sysinternals.com](http://www.sysinternals.com). Some of these even provide more information than the corresponding Microsoft tools.

Windows 2000 is an immensely complex system, now consisting of over 29 million lines of C code. If printed 50 lines per page and 1000 pages per bound book, the full code would occupy 580 volumes. This opus would occupy 23 running meters of shelf space (for the paperback version). If arranged in bookcases 1 m wide with 6 shelves per bookcase, the set would occupy a wall 4 m wide.

Just for fun, a comparison of a few operating system source code sizes is given in Fig. 11-3. However, this table should be taken with a grain (or better yet, a metric ton) of salt because what constitutes the operating system is different for different systems. For example, the entire window system and GUI is part of the kernel in Windows, but not in any UNIX version. It is simply a user process there. Counting X Windows adds another 1.5 million lines of code to all the UNIX versions, and that does not even count the GUI code (Motif, GNOME, etc.), which is also not part of the operating system in the UNIX world. Additionally, some systems include code for multiple architectures (e.g., five for 4.4 BSD and nine for Linux), with each architecture adding 10,000 to 50,000 lines of code. The reason Free BSD 1.0 has only 235,000 lines of code whereas 4BSD Lite, from which it is derived, has 743,000 lines is that support for all the obsolete architectures (e.g., the VAX) was dropped in Free BSD.

Also, the number of file systems, devices drivers, and libraries supplied varies greatly from system to system. In addition, Windows contains large amounts of test code that UNIX does not contain as well as some utilities and support for numerous languages besides English. Finally, the measurements were made by different people, which introduces considerable variance (e.g., did makefiles, headers, configuration files and documentation count and how much was there?). This is not like comparing apples with oranges; it is like comparing apples with telephones. However, all the counts within a single family came from the same source, so intrafamily counts are somewhat meaningful.

Despite all these disclaimers, two conclusions are fairly clear:

1. System bloat seems to be as inevitable as death and taxes.
2. Windows is much bigger than UNIX

Whether small is beautiful or big is beautiful is a matter of heated controversy. The argument for the former is that small size and a lean-and-mean mentality produces a manageable, reliable system that users can understand. The argument for the latter is that many users want lots of features. In any event, it should also be clear that any students planning to write a full-blown, state-of-the-art operating system from scratch have their work cut out for them.

Although Windows 2000 is already the world heavyweight champion in terms of pure mass, it is still growing, with bugs being fixed and new features being added. The way Microsoft manages its development is worth noting. Hundreds



Year	AT&T	BSD	MINIX	Linux	Solaris	Win NT
1976	V6 9K					
1979	V7 21K					
1980		4.1 38K				
1982	Sys III 58K					
1984		4.2 98K				
1986		4.3 179K				
1987	SVR3 92K		1.0 13K			
1989	SVR4 280K					
1991				0.01 10K		
1993		Free 1.0 235K			5.3 850K	3.1 6M
1994		4.4 Lite 743K		1.0 165K		3.5 10M
1996				2.0 470K		4.0 16M
1997			2.0 62K		5.6 1.4M	
1999				2.2 1M		
2000		Free 4.0 1.4M			5.8 2.0M	2000 29M

**Figure 11-3.** A comparison of some operating system sizes. The first string in each box is the version; the second is the size measured in lines of source code, where K = 1000 and M = 1,000,000. Comparisons within a column have real meaning; comparisons across columns do not, as discussed in the text.

of programmers work on various aspects of Windows 2000 all day. Whenever a piece of code is finished, the programmer submits it electronically to the build team. At 6 P.M. every day, the door is closed and the system is rebuilt (i.e., recompiled and linked). Each build gets a unique sequence number, which can be seen by examining the version number of *ntoskrnl.exe* (the first public release of Windows 2000 was build 2195).

The new operating system is electronically distributed to thousands of machines around the Microsoft campus in Redmond, WA, where it is subjected to intense stress tests all night. Early the next morning, the results of all the tests are sent to the relevant groups, so they can see if their new code works. Each team then decides which code they want to work on that day. During the day, the programmers work on their code and at 6 P.M. the build-and-test-cycle begins anew.

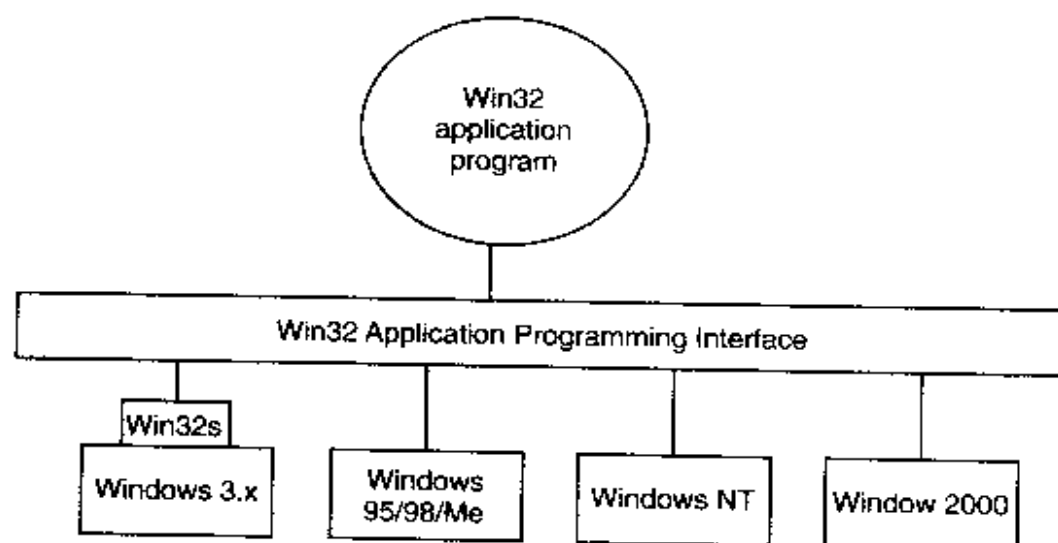
## 11.2 PROGRAMMING WINDOWS 2000

It is now time to start our technical study of Windows 2000. However, before getting into the details of the internal structure, we will first take a look at the programming interface and the registry, a small in-memory data base.

### 11.2.1 The Win32 Application Programming Interface

Like all other operating systems, Windows 2000 has a set of system calls it can perform. However, Microsoft has never made the list of Windows system calls public, and it also changes them from release to release. Instead, what Microsoft has done is define a set of function calls called the **Win32 API (Win32 Application Programming Interface)** that are publicly known and fully documented. These are library procedures that either make system calls to get the work done, or, in some cases, do the work right in user space. The existing Win32 API calls do not change with new releases of Windows, although new API calls are added frequently.

Binary programs for the Intel x86 that adhere exactly to the Win32 API interface will run unmodified on all versions of Windows since Windows 95. As shown in Fig. 11-4, an extra library is needed for Windows 3.x to match a subset of the 32-bit API calls to the 16-bit operating system, but for the other systems no adaptation is needed. It should be noted that Windows 2000 adds substantial new functionality to Win32, so it has additional API calls not included on older versions of Win32 and which will not work on older versions of Windows.



**Figure 11-4.** The Win32 API allows programs to run on almost all versions of Windows.

The Win32 API philosophy is completely different from the UNIX philosophy. In the latter, the system calls are all publicly known and form a minimal interface: removing even one of them would reduce the functionality of the operating system. The Win32 philosophy is to provide a very comprehensive interface, often with three or four ways of doing the same thing, and including many functions (i.e., procedures) that clearly should not be (and are not) system calls, such as an API call to copy an entire file.

Many Win32 API calls create kernel objects of one kind or another, including files, processes, threads, pipes, and so on. Every call creating an object returns a

result called a **handle** to the caller. This handle can subsequently be used to perform operations on the object. Handles are specific to the process that created the object referred to by the handle. They cannot be passed directly to another process and used there (just as UNIX file descriptors cannot be passed to other processes and used there). However, under certain circumstances, it is possible to duplicate a handle and pass it to other processes in a protected way, allowing them controlled access to objects belonging to other processes. Every object also has a security descriptor associated with it, telling in detail who may and may not perform what kinds of operations on the object.

Not all system-created data structures are objects and not all objects are kernel objects. The only ones that are true kernel objects are those that need to be named, protected, or shared in some way. Every kernel object has a system-defined type, has well-defined operations on it, and occupies storage in kernel memory. Although users can perform the operations (by making Win32 calls), they cannot get at the data directly.

The operating system itself can also create and use objects and does so heavily. Most of these objects are created to allow one component of the system to store some information for a substantial period of time or to pass some data structure to another component. For example, when a device driver is loaded, an object is created for it holding its properties and pointers to the functions it contains. Within the operating system, the driver is then referred to by using its object.

Windows 2000 is sometimes said to be object-oriented because the only way to manipulate objects is by invoking operations on their handles by making Win32 API calls. On the other hand, it lacks some of the most basic properties of object-oriented systems such as inheritance and polymorphism.

The Win32 API calls cover every conceivable area an operating system could deal with, and quite a few it arguably should not deal with. Naturally, there are calls for creating and managing processes and threads. There are also many calls that relate to interprocess (actually, interthread) communication, such as creating, destroying, and using mutexes, semaphores, events, and other IPC objects.

Although much of the memory management system is invisible to programmers (fundamentally, it is just demand paging), one important feature is visible: namely the ability of a process to map a file onto a region of its virtual memory. This allows the process the ability to read and write parts of the file as though they were memory words.

An important area for many programs is file I/O. In the Win32 view, a file is just a linear sequence of bytes. Win32 provides over 60 calls for creating and destroying files and directories, opening and closing files, reading and writing them, requesting and setting file attributes, and much more.

Another area for which Win32 provides calls is security. Every process has an ID telling who it is and every object can have an access control list telling in great detail precisely which users may access it and which operations they may

perform on it. This approach provides for a fine-grained security in which specific individuals can be allowed or denied specific access to every object.

Processes, threads, synchronization, memory management, file I/O, and security system calls are nothing new. Other operating systems have them too, although generally not hundreds of them, as Win32 does. But what really distinguishes Win32 are the thousands upon thousands of calls for the graphical interface. There are calls for creating, destroying, managing and using windows, menus, tool bars, status bars, scroll bars, dialog boxes, icons, and many more items that appear on the screen. There are calls for drawing geometric figures, filling them in, managing the color palettes they use, dealing with fonts, and placing icons on the screen. Finally, there are calls for dealing with the keyboard, mouse and other input devices as well as audio, printing, and other output devices. In short, the Win32 API (especially the GUI part) is immense and we could not even begin to describe it in any detail in this chapter, so we will not try. Interested readers should consult one of the many books on Win32 (e.g., Petzold, 1999; Simon, 1997; and Rector and Newcomer, 1997).

Although the Win32 API is available on Windows 98 (as well as on the consumer electronics operating system, Windows CE), not every version of Windows implements every call and sometimes there are minor differences as well. For example, Windows 98 does not have any security, so those API calls that relate to security just return error codes on Windows 98. Also, Windows 2000 file names use the Unicode character set, which is not available on Windows 98 and Windows 98 file names are not case sensitive, whereas Windows 2000 file names are case sensitive (although some kinds of searches on file names are not case sensitive). There are also differences in parameters to some API function calls. On Windows 2000, for example, all the screen coordinates given in the graphics functions are true 32-bit numbers; on Windows 98, only the low-order 16 bits are used because much of the graphics subsystem is still 16-bit code. The existence of the Win32 API on several different operating systems makes it easier to port programs between them, but since these minor variations exist, some care must be taken to achieve portability.

### 11.2.2 The Registry

Windows needs to keep track of a great deal of information about hardware, software, and users. In Window 3.x, this information was stored in hundreds of *.ini* (initialization) files spread all over the disk. Starting with Windows 95, nearly all the information needed for booting and configuring the system and tailoring it to the current user was gathered in a big central database called the **registry**. In this section we will give an overview of the Windows 2000 registry.

To start with, it is worth noting that although many parts of Windows 2000 are complicated and messy, the registry is one of the worst, and the cryptic nomenclature does not make it any better. Fortunately, entire books have been

written describing it (Born, 1998; Hipson, 2000; and Ivens, 1998). That said, the *idea* behind the registry is very simple. It consists of a collection of directories, each of which contains either subdirectories or entries. In this respect it is a kind of file system for very small files. It has directories and entries (the files).

The confusion starts with the fact that Microsoft calls a directory a **key**, which it definitely is not. Furthermore, all the top-level directories start with the string *HKEY*, which means handle to key. Subdirectories tend to have somewhat better chosen names, although not always.

At the bottom of the hierarchy are the entries, called **values**, which contain the information. Each value has three parts: a name, a type, and the data. The name is just a Unicode string, often **default** if the directory contains only one value. The type is one of 11 standard types. The most common ones are Unicode string, a list of Unicode strings, a 32-bit integer, an arbitrary length binary number, and a symbolic link to a directory or entry elsewhere in the registry. Symbolic names are completely analogous to symbolic links in file systems or shortcuts on the Windows desktop: they allow one entry to point to another entry or directory. A symbolic link can also be used as a key, meaning that something that appears to be a directory is just a pointer to a different directory.

At the top level, the Windows 2000 registry has six keys, called **root keys**, as listed in Fig. 11-5. Some interesting **subkeys** (subdirectories) are also shown here. To see this list on your system, use one of the registry editors, either *regedit* or *regedt32*, which unfortunately display different information and use different formats. They can also change registry values. Amateurs should not change keys or values on any system they plan to boot again. Just looking is safe, though. You have been warned.

The first key (i.e., directory), *HKEY\_LOCAL\_MACHINE*, is probably the most important as it contains all the information about the local system. It has five subkeys (i.e., subdirectories). The *HARDWARE* subkey contains many subkeys telling all about the hardware and which driver controls which piece of hardware. This information is constructed on the fly by the plug-and-play manager as the system boots. Unlike the other subkeys, it is not stored on disk.

The *SAM* (Security Account Manager) subkey contains the user names, groups, passwords, and other account and security information needed for logging in. The *SECURITY* subkey contains general security policy information, such as minimum length for passwords, how many failed login attempts are tolerated before the police are called, etc.

The *SOFTWARE* subkey is where software manufacturers store preferences etc. For example, if a user has *Adobe Acrobat*, *Photoshop* and *Premiere* installed, there will be a subkey *Adobe* here, and below that further subkeys for *Acrobat*, *Photoshop*, *Premiere*, and any other Adobe products. The entries in these subdirectories can store anything the Adobe programmers want to put there, generally system-wide properties such as the version and build number, how to uninstall the package, drivers to use, and so forth. The registry saves them the trouble of hav-

Key	Description
HKEY_LOCAL_MACHINE	Properties of the hardware and software
HARDWARE	Hardware description and mapping of hardware to drivers
SAM	Security and account information for users
SECURITY	System-wide security policies
SOFTWARE	Generic information about installed application programs
SYSTEM	Information for booting the system
HKEY_USERS	Information about the users; one subkey per user
USER-AST-ID	User AST's profile
AppEvents	Which sound to make when (incoming email/fax, error, etc.)
Console	Command prompt settings (colors, fonts, history, etc.)
Control Panel	Desktop appearance, screensaver, mouse sensitivity, etc.
Environment	Environment variables
Keyboard Layout	Which keyboard: 102-key US, AZERTY, Dvorak, etc.
Printers	Information about installed printers
Software	User preferences for Microsoft and third party software
HKEY_PERFORMANCE_DATA	Hundreds of counters monitoring system performance
HKEY_CLASSES_ROOT	Link to HKEY_LOCAL_MACHINE\SOFTWARE\CLASSES
HKEY_CURRENT_CONFIG	Link to the current hardware profile
HKEY_CURRENT_USER	Link to the current user profile

**Figure 11-5.** The root keys registry keys and selected subkeys. The capitalization has no meaning but follows the Microsoft practice here.

ing to invent their own method for storing this information. User-specific information also goes in the registry, but under HKEY\_USERS.

The SYSTEM subkey holds mostly information about booting the system, for example, the list of drivers that must be loaded. It also holds the list of services (daemons) that must be started after booting up and the configuration information for all of them.

The next top-level key is HKEY\_USERS, which contains the profiles of all the users. All the user-specific preferences in a number of areas are stored here. When a user changes a preference using the control panel, for example, the desktop color scheme, the new settings are recorded here. In fact, many of the programs on the control panel do little more than collect user information and change the registry accordingly. Some of the subkeys under HKEY\_USERS are shown in Fig. 11-5 and should need little additional comment. Some of the subkeys, such as Software, contain surprisingly large numbers of subkeys, even if no software packages are installed.

The next top-level key, HKEY\_PERFORMANCE\_DATA contains neither data read in from the disk nor data collected by the plug-and-play manager. Instead, it offers a window into the operating system. The system itself contains hundreds of counters for monitoring system performance. These counters are accessible via

this registry key. When a subkey is queried, a specified procedure is run to collect and return the information (possibly by reading one or more counters and combining them in some way). This key is not visible using *regedit* or *regedt32*. Instead one has to use the performance tools, such as *pfmon*, *perfmon*, and *pview*. There are many such tools, some on the Windows 2000 CD-ROM, some in the resource kits, and some from third parties.

The next three top-level keys do not actually exist. Each one is a symbolic link to some place elsewhere in the registry. The `HKEY_CLASSES_ROOT` key is the most interesting. It points to the directory that handles COM (Component Object Model) objects and also the associations between file extensions and programs. When a user double clicks on a file ending in, say, *.doc*, the program catching the mouse click looks here to see which program to run (probably *Microsoft Word*). The complete database of recognized extensions and which program each one is owned by is under this key.

The `HKEY_CURRENT_CONFIG` key links to the current hardware configuration. A user can construct multiple hardware configurations, for example by disabling various devices to see if they were the cause of strange system behavior. This key points to the current configuration. Similarly, `HKEY_CURRENT_USER` points to the current user so that user's preferences can be found quickly.

None of the last three keys really adds anything, since the information was available elsewhere anyway (although less conveniently accessible). Thus despite the fact that *regedit* and *regedt32* list five top-level keys, there are really only three top-level directories and one of those is not shown among the five displayed.

The registry is fully available to the Win32 programmer. There are calls to create and delete keys, look up values within keys, and more. Some of the more useful ones are listed in Fig. 11-6.

Win32 API function	Description
RegCreateKeyEx	Create a new registry key
RegDeleteKey	Delete a registry key
RegOpenKeyEx	Open a key to get a handle to it
RegEnumKeyEx	Enumerate the subkeys subordinate to the key of the handle
RegQueryValueEx	Look up the data for a value within a key

Figure 11-6. Some of the Win32 API calls for using the registry

When the system is turned off, most of the registry information (but not all, as discussed above) is stored on the disk in files called **hives**. Most of them are in `\winnt\system32\config`. Because their integrity is so critical to correct system functioning, when they are updated, backups are made automatically and writes are done using atomic transactions to prevent corruption in the event of system crash during the write. Loss of the registry requires reinstalling all software.

## 11.3 SYSTEM STRUCTURE

In the previous sections we examined Windows 2000 as seen by the programmer. Now we are going to look under the hood to see how the system is organized internally, what the various components do, and how they interact with each other and with user programs. Although there are many books on how to use Windows 2000, there are many fewer on how it works. Far and away the best place to look for additional information on this topic is *Inside Windows 2000*, 3rd ed. by Solomon and Russinovich (2000). Some of the material in this chapter is based on information in that book and information from the authors. Microsoft was also a key source.

### 11.3.1 Operating System Structure

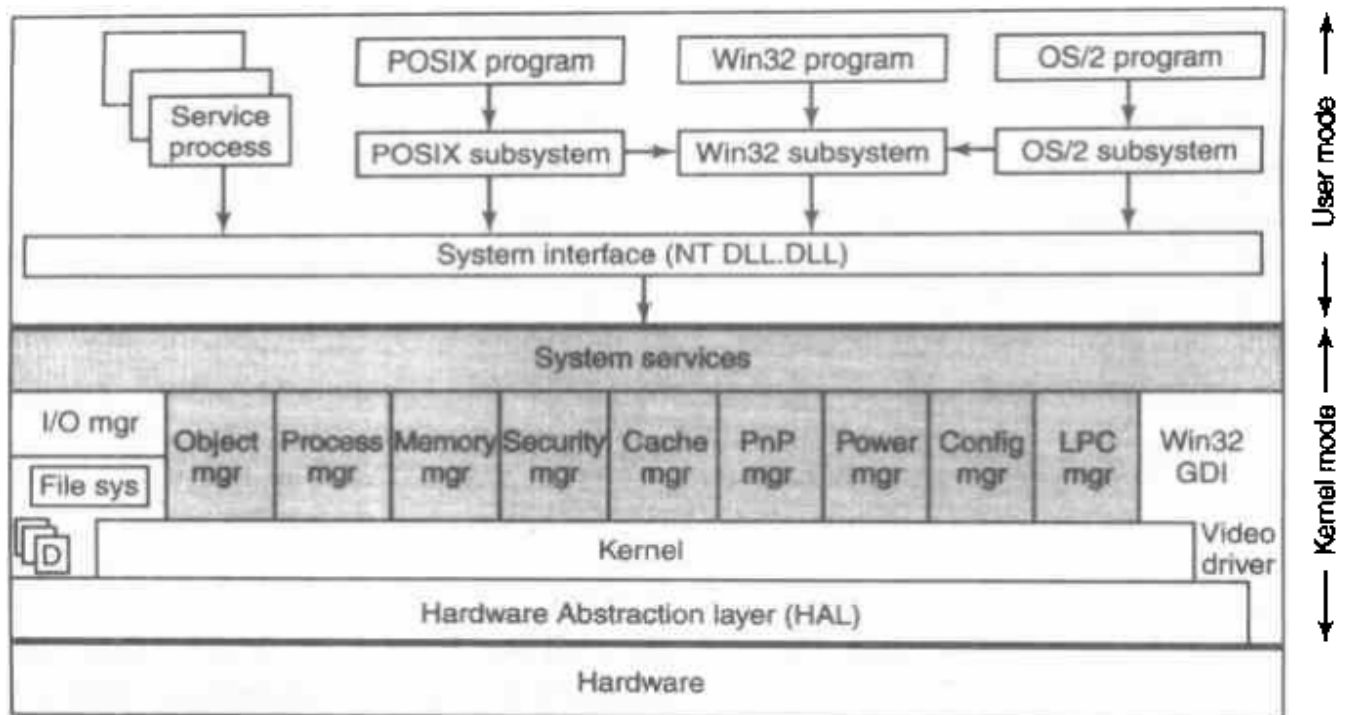
Windows 2000 consists of two major pieces: the operating system itself, which runs in kernel mode, and the environment subsystems, which run in user mode. The kernel is a traditional kernel in the sense that it handles process management, memory management, file systems, and so on. The environment subsystems are somewhat unusual in that they are separate processes that help user programs carry out certain system functions. In the following sections we will examine each of these parts in turn.

One of NT's many improvements over Windows 3.x was its modular structure. It consisted of a moderately small kernel that ran in kernel mode, plus some server processes that ran in user mode. User processes interacted with the server processes using the client-server model: a client sent a request message to a server, and the server did the work and returned the result to the client via a second message. This modular structure made it easier to port it to several computers besides the Intel line, including the DEC Alpha, IBM PowerPC, and SGI MIPS. It also protected the kernel from bugs in the server code. However, for performance reasons, starting with NT 4.0, pretty much all of the operating system (e.g., system call handling and all of the screen graphics) was put back into kernel mode. This design was carried over to Windows 2000.

Nevertheless, there is still some structure in Windows 2000. It is divided into several layers, each one using the services of the ones beneath it. The structure is illustrated in Fig. 11-7. One of the layers is divided horizontally into many modules. Each module has some particular function and a well-defined interface to the other modules.

The lowest two software layers, the HAL and the kernel, are written in C and in assembly language and are partly machine dependent. The upper ones are written entirely in C and are almost entirely machine independent. The drivers are written in C, or in a few cases C++. Below we will first examine the various components of the system starting at the bottom and working our way up.





**Figure 11-7.** The structure of Windows 2000 (slightly simplified). The shaded area is the executive. The boxes indicated by D are device drivers. The service processes are system daemons.

## The Hardware Abstraction Layer

One of the goals of Windows 2000 (and Windows NT before it) was to make the operating system portable across platforms. Ideally, when a new machine comes along, it should be possible to just recompile the operating system with the new machine's compiler and have it run the first time. Unfortunately, life is not like that. While the upper layers of the operating system can be made completely portable (because they mostly deal with internal data structures), the lower layers deal with device registers, interrupts, DMA, and other hardware features that differ appreciably from machine to machine. Even though most of the low-level code is written in C, it cannot just be scooped up from a Pentium, plopped down on, say, an Alpha, recompiled, and rebooted due to the many small hardware differences between the Pentium and the Alpha that have nothing to do with the different instruction sets and which cannot be hidden by the compiler.

Fully aware of this problem, Microsoft made a serious attempt to hide many of the machine dependencies in a thin layer at the bottom called the **HAL (Hardware Abstraction Layer)**. (The name HAL was no doubt inspired by the computer HAL in the late Stanley Kubrick's movie *2001: A Space Odyssey*. Rumor has it that Kubrick chose the name "HAL" by taking the name of the then-dominant computer company—IBM—and subtracting 1 from each letter.)

The job of the HAL is to present the rest of the operating system with abstract hardware devices, in particular, devoid of the warts and idiosyncracies with which

real hardware is so richly endowed. These devices are presented in the form of machine-independent services (procedure calls and macros) that the rest of the operating system and the drivers can use. By using the HAL services (which are identical on all Windows 2000 systems, no matter what the hardware is) and not addressing the hardware directly, drivers and the kernel require fewer changes when being ported to new hardware. Porting the HAL itself is straightforward because all the machine-dependent code is concentrated in one place and the goals of the port are well defined, namely, implement all of the HAL services.

The services chosen for inclusion in the HAL are those that relate to the chip set on the parentboard and which vary from machine to machine within reasonably predictable limits. In other words, it is designed to hide the differences between one vendor's parentboard and another one's, but not the differences between a Pentium and an Alpha. The HAL services include access to the device registers, bus-independent device addressing, interrupt handling and resetting, DMA transfers, control of the timers and real-time clock, low-level spin locks and multiprocessor synchronization, interfacing with the BIOS and its CMOS configuration memory. The HAL does not provide abstractions or services for specific I/O devices such as keyboards, mice, or disks or for the memory management unit.

As an example of what the hardware abstraction layer does, consider the issue of memory-mapped I/O versus I/O ports. Some machines have one and some have the other. How should a driver be programmed: to use memory-mapped I/O or not? Rather than forcing a choice, which would make the driver not portable to a machine that did it the other way, the hardware abstraction layer offers three procedures for driver writers to use for reading the device registers and another three for writing them:

```
uc = READ_PORT_UCHAR(port);      WRITE_PORT_UCHAR(port, uc);  
us = READ_PORT_USHORT(port);     WRITE_PORT_USHORT(port, us);  
ul = READ_PORT_ULONG(port);      WRITE_PORT_LONG(port, ul);
```

These procedures read and write unsigned 8-, 16-, and 32-bit integers, respectively, to the specified port. It is up to the hardware abstraction layer to decide whether memory-mapped I/O is needed here. In this way, a driver can be moved without modification between machines that differ in the way the device registers are implemented.

Drivers often need to access specific I/O devices for various purposes. At the hardware level, a device has one or more addresses on a certain bus. Since modern computers often have multiple buses (ISA, PCI, SCSI, USB, 1394, etc.), it can happen that two or more devices have the same bus address, so some way is needed to distinguish them. The HAL provides a service for identifying devices by mapping bus-relative device addresses onto system-wide logical addresses. In this way, drivers are not required to keep track of which device is on which bus.

These logical addresses are analogous to the handles the operating system gives user programs to refer to files and other system resources. This mechanism also shields higher layers from properties of alternative bus structures and addressing conventions.

Interrupts have a similar problem—they are also bus dependent. Here, too, the HAL provides services to name interrupts in a system-wide way and also provides services to allow drivers to attach interrupt service routines to interrupts in a portable way, without having to know anything about which interrupt vector is for which bus. Interrupt request level management is also handled in the HAL.

Another HAL service is setting up and managing DMA transfers in a device-independent way. Both the system-wide DMA engine and DMA engines on specific I/O cards can be handled. Devices are referred to by their logical addresses. The HAL also implements software scatter/gather (writing or reading from noncontiguous blocks of physical memory).

The HAL also manages clocks and timers in a portable way. Time is kept track of in units of 100 nsec starting at 1 January 1601, which is far more precise than MS-DOS's keeping track of time in units of 2 sec since 1 January 1980 and provides support for the many computer-related activities in the 17th, 18th, and 19th centuries. The time services decouple the drivers from the actual frequencies at which the clocks run.

Kernel components sometimes need to synchronize at a very low level, especially to prevent race conditions in multiprocessor systems. The HAL provides some primitives to manage this synchronization, such as spin locks, in which one CPU simply waits for a resource held by another CPU to be released, particularly in situations where the resource is typically only held for a few machine instructions.

Finally, after the system has been booted, the HAL talks to the BIOS and inspects the CMOS configuration memory, if any, to find out which buses and I/O devices the system contains and how they have been configured. This information is then put into the registry so other system components can look it up without having to understand how the BIOS or configuration memory work. A summary of some of the things the HAL does is given in Fig. 11-8.

Since the HAL is highly-machine dependent, it must match the system it is installed on perfectly, so a variety of HALs are provided on the Windows 2000 CD-ROM. At system installation time, the appropriate one is selected and copied to the system directory `\winnt\system32` on the hard disk as *hal.dll*. All subsequent boots use this version of the HAL. Removing this file will make the system unbootable.

Although the HAL is reasonably efficient, for multimedia applications, it may not be fast enough. For this reason, Microsoft also produced a software package called **DirectX**, which augments the HAL with additional procedures and allows user processes much more direct access to the hardware. DirectX is somewhat specialized, so we will not discuss it further in this chapter.

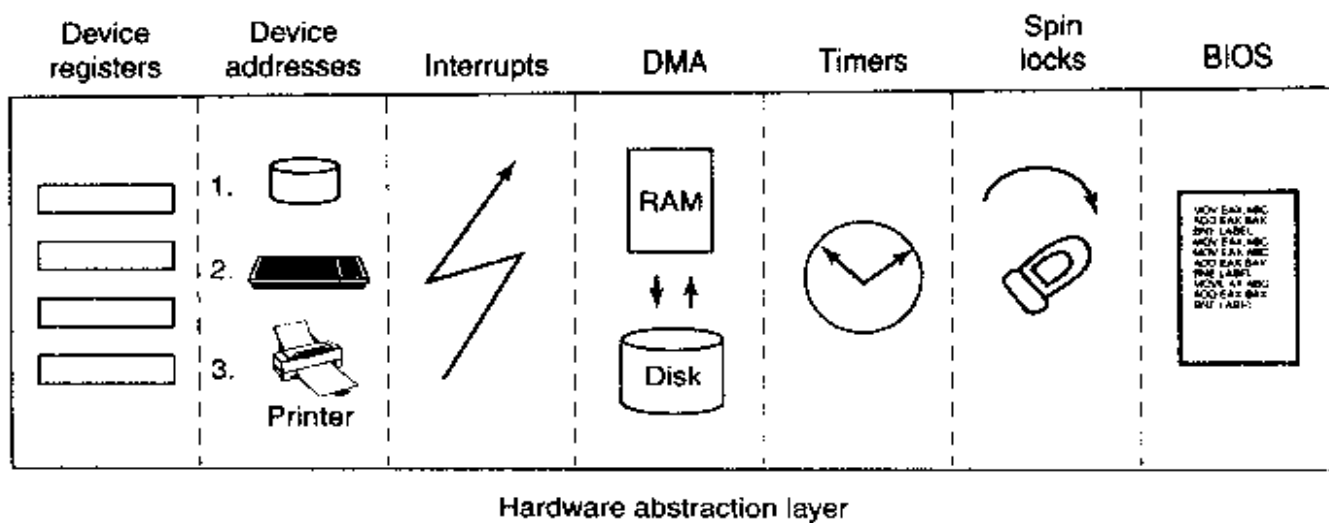


Figure 11-8. Some of the hardware functions the HAL manages.

## The Kernel Layer

Above the hardware abstraction layer is a layer that contains what Microsoft calls the **kernel**, as well as the device drivers. Some early documentation refers to the kernel as the “microkernel,” which it really never was because the memory manager, file system, and other major components resided in kernel space and ran in kernel mode from day one. The kernel is certainly not a microkernel now since virtually the entire operating system was put in kernel space starting with NT 4.0.

In the chapter on UNIX, we used the term “kernel” to mean everything running in kernel mode. In this chapter we will (very grudgingly) reserve the term “kernel” for the part labeled as such in Fig. 11-7 and call the totality of the code running in kernel mode the “operating system.” Part of kernel (and much of the HAL) is permanently resident in main memory (i.e., is not paged). By adjusting its priority, it can control whether it can tolerate being preempted by I/O interrupts or not. Although a substantial fraction of the kernel is machine specific, most of it is nevertheless written in C, except where top performance overshadows all other concerns.

The purpose of the kernel is to make the rest of the operating system completely independent of the hardware, and thus highly portable. It picks up where the HAL leaves off. It accesses the hardware via the HAL and builds upon the extremely low-level HAL services to construct higher-level abstractions. For example, the HAL has calls to associate interrupt service procedures with interrupts, and set their priorities, but does little else in this area. The kernel, in contrast, provides a complete mechanism for doing context switches. It properly saves all the CPU registers, changes the page tables, flushes the CPU cache, and so on, so that when it is done, the previously running thread has been saved in tables in memory. It then sets up the new thread’s memory map and loads its registers so the new thread can start running.

The code for thread scheduling is also in the kernel. When it is time to see if a new thread can run, for example, after a quantum runs out or after an I/O interrupt completes, the kernel chooses the thread and does the context switch necessary to run it. From the point of view of the rest of the operating system, thread switching is automatically handled by lower layers without any work on their part and in a portable way. The scheduling algorithm itself will be discussed later in this chapter when we come to processes and threads.

In addition to providing a higher-level abstraction of the hardware and handling thread switches, the kernel also has another key function: providing low-level support for two classes of objects: control objects and dispatcher objects. These objects are not the objects that user processes get handles to, but are internal objects upon which the executive builds the user objects.

**Control objects** are those objects that control the system, including primitive process objects, interrupt objects, and two somewhat strange objects called DPC and APC. A **DPC (Deferred Procedure Call)** object is used to split off the non-time-critical part of an interrupt service procedure from the time critical part. Generally, an interrupt service procedure saves a few volatile hardware registers associated with the interrupting I/O device so they do not get overwritten and reenables the hardware, but saves the bulk of the processing for later.

For example, after a key is struck, the keyboard interrupt service procedure reads the key code from a register and reenables the keyboard interrupt, but does not need to process the key immediately, especially if something more important (i.e., higher priority) is currently going on. As long as the key is processed within about 100 msec, the user will be none the wiser. DPCs are also used for timer expirations and other activities whose actual processing need not be instantaneous. The DPC queue is the mechanism for remembering that there is more work to do later.

Another kernel control object is the **APC (Asynchronous Procedure Call)**. APCs are like DPCs except that they execute in the context of a specific process. When processing a key press, it does not matter whose context the DPC runs in because all that is going to happen is that the key code will be inspected and probably put in a kernel buffer. However, if an interrupt requires copying a buffer from kernel space to a buffer in some user process' address space (e.g., as it may on completion of a read from the modem), then the copying procedure needs to run in the receiver's context. The receiver's context is needed so the page table will contain both the kernel buffer and the user buffer (all processes contain the entire kernel in their address spaces, as we will see later). For this reason, the kernel distinguishes between DPCs and APCs.

The other kind of kernel objects are **dispatcher objects**. These include semaphores, mutexes, events, waitable timers, and other objects that threads can wait on. The reason that these have to be handled (in part) in the kernel is that they are intimately intertwined with thread scheduling, which is a kernel task. As a little aside, mutexes are called "mutants" in the code because they were required to

implement the OS/2 semantics of not automatically unlocking themselves when a thread holding one exited, something the Windows 2000 designers considered bizarre. (The OS/2 semantics are relevant because NT was originally conceived of as a replacement for OS/2, the operating system shipped on IBM's PC/2.)

## The Executive

Above the kernel and device drivers is the upper portion of the operating system, called the **executive**, shown as the shaded area in Fig. 11-7. The executive is written in C, is architecture independent, and can be ported to new machines with relatively little effort. It consists of 10 components, each of which is just a collection of procedures that work together to accomplish some goal. There are no hard boundaries between the pieces and different authors describing the executive might even group the procedures differently into components. It should be noted that components on the same level can (and do) call each other extensively.

The **object manager** manages all objects known to the operating system. These include processes, threads, files, directories, semaphores, I/O devices, timers, and many others. The object manager allocates a block of virtual memory from kernel address space when an object is created and returns it to the free list when the object is deallocated. Its job is to keep track of all the objects.

To avoid any confusion, most of the executive components labeled "manager" in Fig. 11-7 are not processes or threads, but merely collections of procedures that other threads can execute when in kernel mode. A few of them, such as the power manager and plug-and-play manager, really are independent threads though.

The object manager also manages a name space in which newly created objects may be placed so they can be referred to later. All other components of the executive use objects heavily to do their work. Objects are so central to the functioning of Windows 2000 that they will be discussed in detail in the next section.

The **I/O manager** provides a framework for managing I/O devices and provides generic I/O services. It provides the rest of the system with device-independent I/O, calling the appropriate driver to perform physical I/O. It is also home to all the device drivers (indicated by *D* in Fig. 11-7). The file systems are technically device drivers under control of the I/O manager. Two different ones are present for the FAT and NTFS file systems, each one independent of the others and controlling different disk partitions. All the FAT file systems are managed by a single driver. We will study I/O further in Sec. 11.6 and one of the file systems, NTFS, in Sec. 11.7.

The **process manager** handles processes and threads, including their creation and termination. It deals with the mechanisms used to manage them, rather than policies about how they are used. It builds upon the kernel process and thread objects and adds extra functionality to them. It is the key to multiprogramming in Windows 2000. We will study process and thread management in Sec. 11.4.

The **memory manager** implements Windows 2000's demand-paged virtual memory architecture. It manages the mapping of virtual pages onto physical page frames. It thereby enforces the protection rules that restrict each process to only access those pages belonging to its address space and not to other processes' address spaces (except under special circumstances). It also handles certain system calls that relate to virtual memory. We will study memory management in Sec. 11.5.

The **security manager** enforces Windows 2000's elaborate security mechanism, which meets the U.S. Dept. of Defense's Orange Book C2 requirements. The Orange Book specifies a large number of rules that a conforming system must meet, starting with authenticated login through how access control is handled, to the fact that virtual pages must be zeroed out before being reused. We will study the security manager in Sec. 11.8.

The **cache manager** keeps the most recently used disk blocks in memory to speed up access to them in the (likely) event that they are needed again. Its job is to figure out which blocks are probably going to be needed again and which ones are not. It is possible to configure Windows 2000 with multiple file systems, in which case the cache manager works for all of them, so each one does not have to do its own cache management. When a block is needed, the cache manager is asked to supply it. If it does not have the block, the cache manager calls upon the appropriate file system to get it. Since files can be mapped into processes' address spaces, the cache manager must interact with the virtual memory manager to provide the necessary consistency. The amount of space devoted to caching is dynamic and can increase or decrease as demands on it change. We will study the cache manager in Sec. 11.9.

The **plug-and-play manager** is sent all notifications of newly attached devices. For some devices, a check is made at boot time and not thereafter. Other devices, for example, USB devices, can be attached at any time and their attachment triggers a message to the plug-and-play manager, which then locates and loads the appropriate driver.

The **power manager** rides herd on power usage. This consists of turning off the monitor and disks after they have been idle for a while. On laptops, the power manager monitors battery usage and takes action when the battery is about to run dry. Such action typically tells programs to save their files and prepare for a graceful shutdown. And pronto.

The **configuration manager** is in charge of the registry. It adds new entries and looks up keys when asked to.

The **local procedure call manager** provides for a highly-efficient interprocess communication used between processes and their subsystems. Since this path is needed to carry out some system calls, efficiency is critical here, which is why the standard interprocess communication mechanisms are not used.

The Win32 GDI executive module handles certain system calls (but not all of them). It was originally in user space but was moved to kernel space in NT 4.0 to



improve performance. The **GDI (Graphics Device Interface)** handles image management for the monitor and printers. It provides system calls to allow user programs to write on the monitor and printers in a device-independent way. It also contains the window manager and display driver. Prior to NT 4.0, it, too, was in user space but the performance was disappointing, so Microsoft moved it into the kernel to speed it up. It is worth mentioning that Fig. 11-7 is not at all to scale. For example, the Win32 and graphics device interface module is larger than the rest of the executive combined.

At the top of the executive is a thin layer called **system services**. Its function is to provide an interface to the executive. It accepts the true Windows 2000 system calls and calls other parts of the executive to have them executed.

At boot time, Windows 2000 is loaded into memory as a collection of files. The main part of the operating system, consisting of the kernel and executive, is located in the file *ntoskrnl.exe*. The HAL is a shared library located in a separate file, *hal.dll*. The Win32 and graphics device interface are together in a third file, *win32k.sys*. Finally, many device drivers are also loaded. Most of these have extension *.sys*.

Actually, things are not quite that simple. The *ntoskrnl.exe* file comes in uniprocessor and multiprocessor versions. Also, there are versions for the Xeon processor, which can have more than 4 GB of physical memory and the Pentium, which cannot. Finally, versions can consist of a free build (sold in stores and preinstalled by computer manufacturers) or a checked build (for debugging purposes). Together there could be eight combinations, although two pairs were combined leaving only six. One of these is copied to *ntoskrnl.exe* when the system is installed.

The checked builds are worth a few words. When a new I/O device is installed on a PC, there is invariably a manufacturer-supplied driver that has to be installed to make it work. Suppose that an IEEE 1394 card is installed on a computer and appears to work fine. Two weeks later the system suddenly crashes. Who does the owner blame? Microsoft.

The bug may indeed be Microsoft's, but some bugs are actually due to flakey drivers, over which Microsoft has no control and which are installed in kernel memory and have full access to all kernel tables as well as the entire hardware. In an attempt to reduce the number of irate customers on the phone, Microsoft tries to help driver writers debug their code by putting statements of the form

**ASSERT(some condition)**

throughout the code. These statements make sanity checks on all parameters to internal kernel procedures (which may be freely called by drivers) and make many other checks as well. The free builds have *ASSERT* defined as a macro that does nothing, removing all the checks. The checked builds have it defined as

```
#define ASSERT(a) if (! (a)) error( ... )
```



causing all the checks to appear in the *ntoskrnl.exe* executable code and be carried out at run time. While this slows down the system enormously, it helps driver writers debug their drivers before they ship them to customers. The checked builds also have numerous other debugging features turned on.

### The Device Drivers

The last part of Fig. 11-7 consists of the **device drivers**. Each device driver can control one or more I/O devices, but a device driver can also do things not related to a specific device, such as encrypting a data stream or even just providing access to kernel data structures. Device drivers are not part of the *ntoskrnl.exe* binary. The advantage of this approach is that once a driver has been installed on a system, it is added to a list in the registry and is loaded dynamically when the system boots. In this way, *ntoskrnl.exe* is the same for everyone, but every system is configured precisely for those devices it contains.

There are device drivers for macroscopically visible I/O devices such as disks and printers, but also for many internal devices and chips that practically no one has ever heard of. In addition, the file systems are also present as device drivers, as mentioned above. The largest device driver, the one for Win32, GDI, and video, is shown on the far right of Fig. 11-7. It handles many system calls and most of the graphics. Since customers can install new device drivers, they have the power to affect the kernel and corrupt the system. For this reason, drivers must be written with great care.

### 11.3.2 Implementation of Objects

Objects are probably the single most important concept in Windows 2000. They provide a uniform and consistent interface to all system resources and data structures such as processes, threads, semaphores, etc. This uniformity has various facets. First, all objects are named and accessed the same way, using object handles. Second, because all accesses to objects go through the object manager, it is possible to put all the security checks in one place and ensure that no process can make an end run around them. Third, sharing of objects among processes can be handled in a uniform way. Fourth, since all object opens and closes go through the object manager, it is easy to keep track of which objects are still in use and which can be safely deleted. Fifth, this uniform model for object management makes it easy to manage resource quotas in a straightforward way.

A key to understanding objects is to realize that an (executive) object is just some number of consecutive words in memory (i.e., in kernel virtual address space). An object is a data structure in RAM, no more and no less. A file on disk is not an object, although an object (i.e., a data structure in kernel virtual address space) is created for a file when it is opened. A consequence of the fact that objects are just kernel data structures is that when the system is rebooted (or

crashes) all objects are lost. In fact, when the system boots, there are no objects present at all (except for the idle and system processes, whose objects are hardwired into the *ntoskrnl.exe* file). All other objects are created on the fly as the system boots up and various initialization (and later user) programs run.

Objects have a structure, as shown in Fig. 11-9. Each object contains a header with certain information common to all objects of all types. The fields in this header include the object's name, the object directory in which it lives in object space, security information (so a check can be made when an object is opened), and a list of processes with open handles to the object (if a certain debugging flag is enabled).

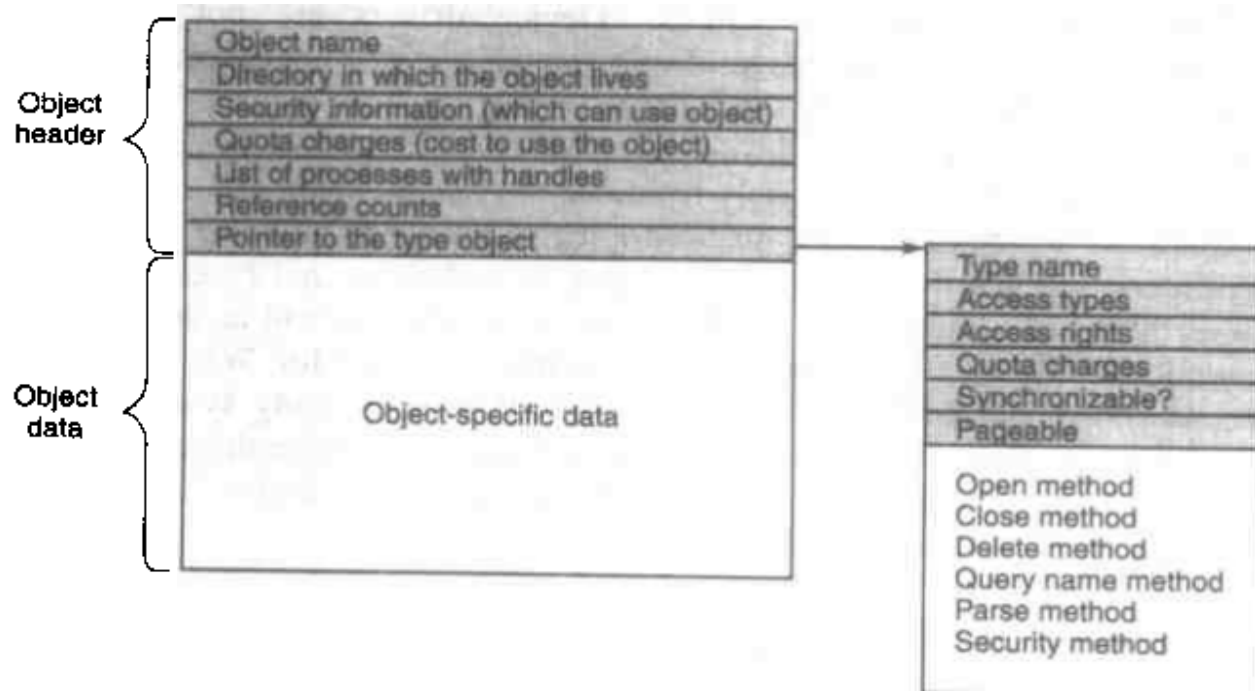


Figure 11-9. The structure of an object.

Each object header also contains a quota charge field, which is the charge levied against a process for opening the object. If a file object costs 1 point and a process belongs to a job that has 10 file points worth of quota, the processes in that job can only open 10 files in total. In this way resource limits can be enforced for each object type separately.

Objects occupy valuable real estate—pieces of kernel virtual address space—so when an object is no longer needed it should be removed and its address space reclaimed. The mechanism for reclamation is to have a reference counter in each object's header. It counts the number of open handles held by processes. This counter is incremented every time the object is opened and decremented every time it is closed. When it hits 0, no more users hold handles to the object. When an object is acquired or released by an executive component, a second counter is incremented or decremented, even though no actual handle is issued. When both counters hit 0, no user process is using the object and no executive process is using the object, so the object can be removed and its memory freed.

The object manager needs to maintain dynamic data structures (its objects), but it is not the only part of the executive with this need. Other pieces also need to allocate and release chunks of kernel memory dynamically. To meet these needs, the executive maintains two page pools in kernel address space: for objects and for other dynamic data structures. Such pools operate as heaps, similar to the C language calls *malloc* and *free* for managing dynamic data. One pool is paged and the other is nonpaged (pinned in memory). Objects that are needed often are kept in the nonpaged pool; objects that are rarely accessed, such as registry keys and some security information, are kept in the paged pool. When memory is tight, the latter can be paged out and faulted back on demand. In fact, substantial portions of the operating system code and data structures are also pageable, to reduce memory consumption. Objects that may be needed when the system is running critical code (and when paging is not permitted) must go in the nonpaged pool. When a small amount of storage is needed, a page can be taken from either pool and then broken up into units as small as 8 bytes.

Objects are typed, which means each one has certain properties common to all objects of its type. The type is indicated by a pointer in the header to a type object, as shown in Fig. 11-9. The type object information includes items such as the type name, whether a thread can wait on the object (yes for mutexes, no for open files), and whether new objects of this type go on the paged or nonpaged pool. Each object points to its type object.

The last thing a type object has is also the most important: pointers to the code for certain standard operations such as open, close, and delete. Whenever one of these operations is invoked on an object, the pointer to the type object is followed and the relevant code located and executed. This mechanism gives the system the opportunity to initialize new objects, and recover storage when they are deleted.

Executive components can create new types dynamically. There is no definitive list of object types, but some of the more common ones are listed in Fig. 11-10. Let us briefly go over the object types in Fig. 11-10. Process and thread are obvious. There is one object for every process and every thread, which holds the main properties needed to manage the process or thread. The next three objects, semaphore, mutex, and event, all deal with interprocess synchronization. Semaphores and mutexes work as expected, but with various extra bells and whistles (e.g., maximum values and timeouts). Events can be in one of two states: signaled or nonsignaled. If a thread waits on an event that is in signaled state, the thread is released immediately. If the event is in nonsignaled state, it blocks until some other thread signals the event, which releases all blocked threads. An event can also be set up, so after a signal has been successfully waited for, it automatically reverts to nonsignaled state, rather than staying in signaled state.

Port, timer, and queue objects also relate to communication and synchronization. Ports are channels between processes for exchanging messages. Timers provide a way to block for a specific time interval. Queues are used to notify threads that a previously started asynchronous I/O operation has completed.

Type	Description
Process	User process
Thread	Thread within a process
Semaphore	Counting semaphore used for interprocess synchronization
Mutex	Binary semaphore used to enter a critical region
Event	Synchronization object with persistent state (signaled/not)
Port	Mechanism for interprocess message passing
Timer	Object allowing a thread to sleep for a fixed time interval
Queue	Object used for completion notification on asynchronous I/O
Open file	Object associated with an open file
Access token	Security descriptor for some object
Profile	Data structure used for profiling CPU usage
Section	Structure used for mapping files onto virtual address space
Key	Registry key
Object directory	Directory for grouping objects within the object manager
Symbolic link	Pointer to another object by name
Device	I/O device object
Device driver	Each loaded device driver has its own object

**Figure 11-10.** Some common executive object types managed by object manager.

Open file objects are created when a file is opened. Files that are not opened do not have objects managed by the object manager. Access tokens are security objects; they identify a user and tell what special privileges the user has, if any. Profiles are structures used for storing periodic samples of the program counter of a running thread to see where the program is spending its time.

Sections are objects used by the memory system for handling memory-mapped files. They record which file (or part thereof) is mapped onto which memory addresses. Keys are registry keys and are used to relate names to values. Object directories are entirely local to the object manager. They provide a way to collect related objects together in exactly the same way directories work in the file system. Symbolic links are also similar to their file system counterparts: they allow a name in one part of the object name space to refer to an object in a different part of the object name space. Each known device has a device object that contains information about it and is used to refer to the device within the system. Finally, each device driver that has been loaded has an object in the object space.

Users can create new objects or open existing objects by making Win32 calls such as `CreateSemaphore` or `OpenSemaphore`. These are calls to library procedures that ultimately result in the appropriate system calls being made. The result of any successful call that creates or opens an object is a 64-bit handle table

entry that is stored in the process' private handle table in kernel memory. The 32-bit index of the handle's position in the table is returned to the user to use on subsequent calls.

The 64-bit handle table entry in the kernel contains two 32-bit words. One word contains a 29-bit pointer to the object's header. The low-order 3 bits are used as flags (e.g., whether the handle is inherited by child processes). These bits are masked off before the pointer is followed. The other word contains a 32-bit rights mask. It is needed because permissions checking is done only at the time the object is created or opened. If a process has only read permission to an object, all the other rights bits in the mask will be 0s, giving the operating system the ability to reject any operation on the object other than reading it.

The handle tables for two processes and their relationships to some objects are illustrated in Fig. 11-11. In this example, process *A* has access to threads 1 and 2 and access to mutexes 1 and 2. Process *B* has access to thread 3 and mutexes 2 and 3. The corresponding entries in the handle tables hold the rights to each of these objects. For example, process *A* might have the rights to lock and unlock its mutexes, but not the right to destroy them. Note that mutex 2 is shared by both processes allowing threads in them to synchronize. The other mutexes are not shared, which might mean that the threads within process *A* use mutex 1 for their internal synchronization and the threads within process *B* use mutex 3 for their internal synchronization.

### The Object Name Space

As objects are created and deleted during execution, the object manager needs a way to keep track of them. To do this job, it maintains a name space, in which all objects in the system are located. The name space can be used by a process to locate and open a handle for some other process' object, provided it has been granted permission to do so. The object name space is one of three name spaces maintained by Windows 2000. The other ones are the file system name space and the registry name space. All three are hierarchical name spaces with multiple levels of directories for organizing entries. The directory objects listed in Fig. 11-10 provide the means to implement this hierarchical name space for objects.

Since executive objects are volatile (i.e., vanish when the computer is shut down, unlike file system and registry entries), when the system boots up, there are no objects in memory and the object name space is empty. During booting, various parts of the executive create directories and then fill them with objects. For example, as the plug-and-play manager discovers devices out there, it creates a device object for each one and enters this object into the name space. When the system is fully booted, all I/O devices, disk partitions, and other interesting discoveries are in the object name space.

Not all objects get entered by the Columbus method—just go look and see what you find. Some executive components look in the registry to see what to do. A key example here is device drivers. During bootup, the system looks in the

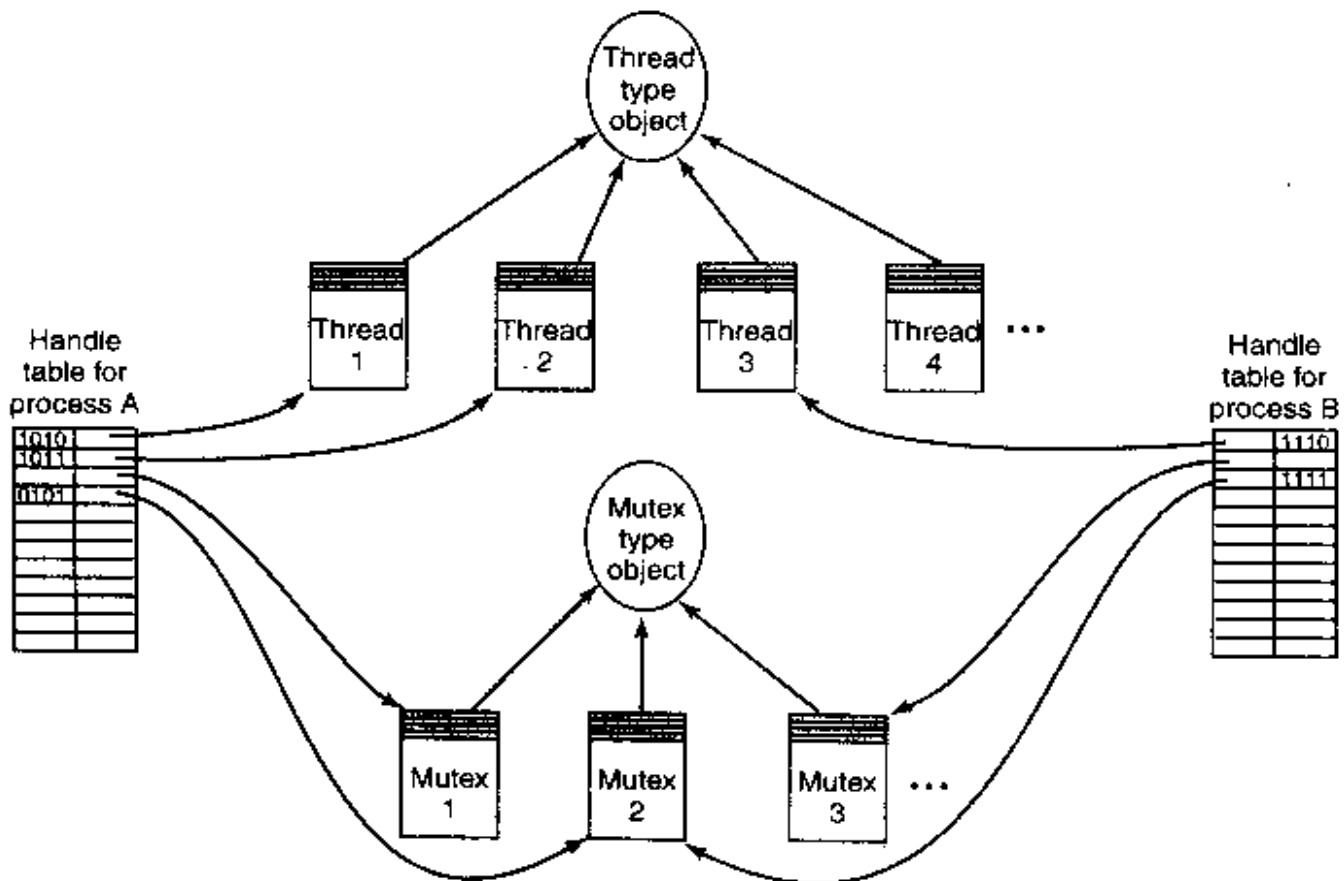


Figure 11-11. The relationship between handle tables, objects, and type objects.

registry to see which device drivers are needed. As they are loaded one by one, an object is created for each one and its name is inserted into the object space. Within the system, the driver is referred to by a pointer to its object.

Although the object name space is crucial to the entire operation of the system, few people know that it even exists because it is not visible to users without special viewing tools. One such viewing tool is *winobj*, available for free at [www.sysinternals.com](http://www.sysinternals.com). When run, this tool depicts an object name space that typically contains the object directories listed in Fig. 11-12 as well as a few others.

The somewhat strangely named directory `\??` contains the names of all the MS-DOS-style device names, such as `A:` for the floppy disk and `C:` for the first hard disk. These names are actually symbolic links to the directory `\Device` where the device objects live. The name `\??` was chosen to make it alphabetically first to speed up lookup of all path names beginning with a drive letter. The contents of the other object directories should be self explanatory.

### 11.3.3 Environment Subsystems

Going back to Fig. 11-7, we see that Windows 2000 consists of components in kernel mode and components in user mode. We have now completed our examination of the kernel mode components; now it is time to look at the user mode

Directory	Contents
??	Starting place for looking up MS-DOS devices like C:
Device	All discovered I/O devices
Driver	Objects corresponding to each loaded device driver
ObjectTypes	The type objects shown in Fig. 11-11
Windows	Objects for sending messages to all the windows
BaseNamedObjs	User-created objects such as semaphores, mutexes, etc.
Arcname	Partition names discovered by the boot loader
NLS	National language support objects
FileSystem	File system driver objects and file system recognizer objects
Security	Objects belonging to the security system
KnownDLLs	Key shared libraries that are opened early and held open

Figure 11-12. Some typical directories in the object name space.

components of which there are three kinds: DLLs, environment subsystems, and service processes. These components work together to provide each user process with an interface that is distinct from the Windows 2000 system call interface.

Windows 2000 supports three different documented APIs: Win32, POSIX, and OS/2. Each of these interfaces has a published list of library calls that programmers can use. The job of the DLLs (Dynamic Link Libraries) and environment subsystems is to implement the functionality of the published interface, thereby hiding the true system call interface from application programs. In particular, the Win32 interface is the official interface for Windows 2000, Windows NT, Windows 95/98/Me, and to a limited extent, Windows CE. By using the DLLs and Win32 environment subsystem, a program can be written to the Win32 specification and run unmodified on all these versions of Windows, even though the system calls are not the same on the various systems.

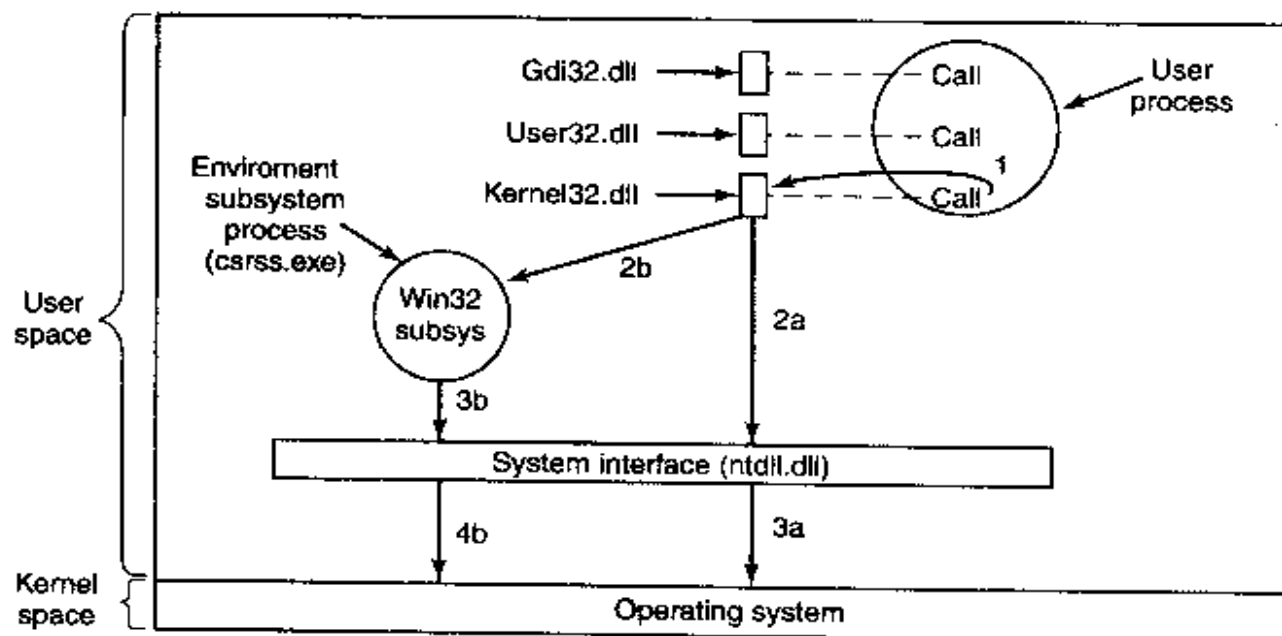
To see how these interfaces are implemented, let us look at Win32. A Win32 program normally contains many calls to Win32 API functions, for example, `CreateWindow`, `DrawMenuBar`, and `OpenSemaphore`. There are thousands of such calls, and most programs use a substantial number of them. One possible implementation would be to statically link every Win32 program with all the library procedures that it uses. If this were done, each binary program would contain one copy of each procedure it used in its executable binary.

The trouble with this approach is that it wastes memory if the user has multiple programs open at once and they use many of the same library procedures. For example, *Word*, *Excel*, and *Powerpoint* all use exactly the same procedures for opening dialog boxes, drawing windows, displaying menus, managing the clipboard, etc., so if a user had all of them open and active at once, there would be three (identical) copies of each of the libraries in memory.

To avoid this problem, all versions of Windows support shared libraries, called **DLLs (Dynamic Link Libraries)**. Each DLL collects together a set of closely related library procedures and their data structures into a single file, usually (but not always) with extension *.dll*. When an application is linked, the linker sees that some of the library procedures belong to DLLs and records this information in the executable's header. Calls to procedures in DLLs are made indirectly through a transfer vector in the caller's address space. Initially this vector is filled with 0s, since the addresses of the procedures to be called are not yet known.

When the application process is started, the DLLs that are needed are located (on disk or in memory) and mapped into the process' virtual address space. The transfer vector is then filled in with the correct addresses so that the procedures can then be called via the transfer vector with only a negligible loss of efficiency. The win here is that even though multiple application programs have the same DLL mapped in, only one copy of the DLL text is needed in physical memory (but each process gets its own copy of the private static data in the DLL). Windows 2000 uses DLLs extremely heavily for all aspects of the system.

Now we have enough background to see how the Win32 and other process interfaces are implemented. Each user process generally links with a number of DLLs that together implement the Win32 interface. To make an API call, one of the procedures in a DLL is called, shown as step 1 in Fig. 11-13. What happens next depends on the Win32 API call. Different ones are implemented in different ways.



**Figure 11-13.** Various routes taken to implement Win32 API function calls.

In some cases, the DLL calls another DLL (*ntdll.dll*) that actually traps to the operating system. This path is shown as steps 2a and 3a in Fig. 11-13. It is also possible that the DLL does all the work itself without making a system call at all.



For other Win32 API calls a different route is taken, namely, first a message is sent to the Win32 subsystem process (*csrss.exe*), which then does some work and then makes a system call (steps 2b, 3b, and 4b). Here, too, in some cases the environment subsystem does all the work in user space and just returns immediately. The message passing between the application process and the Win32 subsystem process has been carefully optimized for performance using a special local procedure call mechanism implemented by the executive and shown as LPC in Fig. 11-7.

In the first version of Windows NT, virtually all the Win32 API calls took route 2b, 3b, 4b, putting a large chunk of the operating system in user space (e.g., the graphics). However, starting with NT 4.0, much of the code was put into kernel mode (in the Win32/GDI driver in Fig. 11-7) for performance reasons. In Windows 2000, only a small number of Win32 API calls, for example process and thread creation, take the long route. The other ones take the direct route, bypassing the Win32 environment subsystem.

As an aside, the three most important DLLs are shown in Fig. 11-13, but they are not the only ones. There are over 800 separate DLLs in the *\winnt\system32* directory totalling 130 MB. To avoid any confusion, the number of DLL files is over 800; the number of API calls contained in them exceeds 13,000. (The 29 million lines of code had to compile into something, after all.) A few of the more important DLLs are listed in Fig. 11-14. The number of exported functions (i.e., those visible outside the file) in each one is given, but these tend to change (meaning increase) over time. The number of exported functions in the first public release of *ntdll.dll* in Windows 2000 is 1179. These are the real system calls. The 1209 calls exported by *ntoskrnl.exe* are the functions available to device drivers and other code linked with the kernel. The list of exported functions in any *.exe* or *.dll* file can be viewed using the *depends* program in the platform SDK Kit.

File	Mode	Fcns	Contents
hal.dll	Kernel	95	Low-level hardware management, e.g., port I/O
ntoskrnl.exe	Kernel	1209	Windows 2000 operating system (kernel + executive)
win32k.sys	Kernel	-	Many system calls including most of the graphics
ntdll.dll	User	1179	Dispatcher from user mode to kernel mode
csrss.exe	User	0	Win32 environment subsystem process
kernel32.dll	User	823	Most of the core (nongraphics) system calls
gdi32.dll	User	543	Font, text, color, brush, pen, bitmap, palette, drawing, etc. calls
user32.dll	User	695	Window, icon, menu, cursor, dialog, clipboard, etc. calls
advapi32.dll	User	557	Security, cryptography, registry, management calls

**Figure 11-14.** Some key Windows 2000 files, the mode they run in, the number of exported function calls, and the main contents of each file. The calls in *win32k.sys* are not formally exported since *win32k.sys* is not called directly.

Although the Win32 process interface is the most important one, there are also two other ones: POSIX and OS/2. The POSIX environment provides minimal support for UNIX applications. It supports only the P1003.1 functionality and little else. It does not have threads, windowing, or networking, for example. In practice, porting any real UNIX program to Windows 2000 using this subsystem is close to impossible. It was included only because parts of the U.S. government require operating systems for government computers to be P1003.1 compliant. This subsystem is not self-contained and uses the Win32 subsystem for much of its work, but without exporting the full Win32 interface to its user programs (which would have made it usable, at no extra cost to Microsoft).

To allow UNIX users to migrate to Windows 2000, Microsoft has a product called Interix that provides a better degree of UNIX compatibility than the POSIX subsystem.

The OS/2 subsystem is similarly limited in functionality and does not support any graphical applications. In practice, it, too, is completely useless. Thus the original idea of having multiple operating system interfaces implemented by different processes in user space is essentially gone. What is left is a full Win32 implementation in kernel mode and little else.

## 11.4 PROCESSES AND THREADS IN WINDOWS 2000

Windows 2000 has a number of concepts for managing the CPU and grouping resources together. In the following sections we will examine these, discussing some of the relevant Win32 API calls, and show how these concepts are implemented.

### 11.4.1 Fundamental Concepts

Windows 2000 supports traditional processes, which can communicate and synchronize with one another, just as they can in UNIX. Each process contains at least one thread, which in turn contains at least one fiber (lightweight thread). Furthermore, processes can be collected into jobs for certain resource management purposes. Together, jobs, processes, threads, and fibers provide a very general set of tools for managing parallelism and resources, both on uniprocessors (single-CPU machines) and on multiprocessors (multiCPU machines). A brief summary of these four concepts is given in Fig. 11-15.

Let us examine these concepts from the largest to the smallest. A **job** in Windows 2000 is a collection of one or more processes that are to be managed as a unit. In particular, there are quotas and resource limits associated with each job, stored in the corresponding job object. The quotas include items such as the maximum number of processes (prevents any process from generating an unbounded number of children), the total CPU time available to each process individually and

Name	Description
Job	Collection of processes that share quotas and limits
Process	Container for holding resources
Thread	Entity scheduled by the kernel
Fiber	Lightweight thread managed entirely in user space

**Figure 11-15.** Basic concepts used for CPU and resource management.

to all the processes combined, and the maximum memory usage, again, per process and total. Jobs can also impose security restrictions on the processes in the job, such as not being able to acquire administrator (superuser) power, even with the proper password.

Processes are more interesting than jobs and also more important. As in UNIX, processes are containers for resources. Every process has a 4-GB address space, with the user occupying the bottom 2 GB (optionally 3 GB on Advanced Server and Datacenter server) and the operating system occupying the rest. Thus the operating system is present in every process' address, although protected from tampering by the memory management unit hardware. A process has a process ID, one or more threads, a list of handles (managed in kernel mode), and an access token holding its security information. Processes are created using a Win32 call that takes as its input the name of an executable file, which defines the initial contents of the address space and creates the first thread.

Every process starts out with one thread, but new ones can be created dynamically. Threads forms the basis of CPU scheduling as the operating system always selects a thread to run, not a process. Consequently, every thread has a state (ready, running, blocked, etc.), whereas processes do not have states. Threads can be created dynamically by a Win32 call that specifies the address within the enclosing process' address space it is to start running at. Every thread has a thread ID, which is taken from the same space as the process IDs, so an ID can never be in use for both a process and a thread at the same time. Process and thread IDs are multiples of four so they can be used as byte indices into kernel tables, the same as other objects.

A thread normally runs in user mode, but when it makes a system call it switches to kernel mode and continues to run as the same thread with the same properties and limits it had in user mode. Each thread has two stacks, one for use when it is in user mode and one for use when it is in kernel mode. In addition to a state, an ID, and two stacks, every thread has a context (in which to save its registers when it is not running), a private area for its own local variables, and possibly its own access token. If it has its own access token, this one overrides the process access token in order to let client threads pass their access rights to server threads who are doing work for them. When a thread is finished executing, it can exit. When the last thread still active in a process exits, the process terminates.

It is important to realize that threads are a scheduling concept, not a resource ownership concept. Any thread is able to access all the objects that belong to its process. All it has to do is grab the handle and make the appropriate Win32 call. There is no restriction on a thread that it cannot access an object because a different thread created or opened it. The system does not even keep track of which thread created which object. Once an object handle has been put in a process' handle table, any thread in the process can use it.

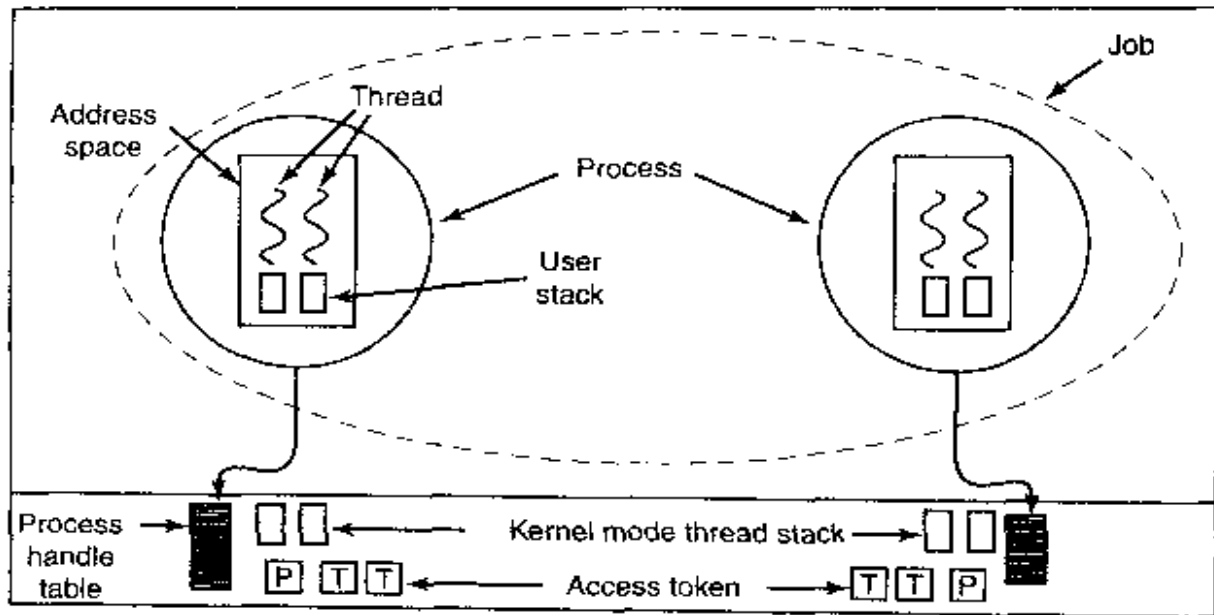
In addition to the normal threads that run within user processes, Windows 2000 has a number of daemon threads that run only in kernel space and are not associated with any user process (they are associated with the special system or idle processes). Some perform administrative tasks, such as writing dirty pages to the disk, while others form a pool that can be assigned to a component of the executive or a driver that needs to get some work done asynchronously in the background. We will study some of these threads later when we come to memory management.

Switching threads in Windows 2000 is relatively expensive because doing a thread switch requires entering and later leaving kernel mode. To provide very lightweight pseudoparallelism, Windows 2000 provides **fibers**, which are like threads, but are scheduled in user space by the program that created them (or its run-time system). Each thread can have multiple fibers, the same way a process can have multiple threads, except that when a fiber logically blocks, it puts itself on the queue of blocked fibers and selects another fiber to run in the context of its thread.

The operating system is not aware of this transition because the thread keeps running, even though it may be first running one fiber, then another. In fact, the operating system knows nothing at all about fibers, so there are no executive objects relating to fibers, as there are for jobs, processes, and threads. There are also no true system calls for managing fibers. However, there are Win32 API calls. These are among the Win32 API calls that do not make system calls, which we mentioned during the discussion of Fig. 11-13. The relationship between jobs, processes, and threads is illustrated in Fig. 11-16.

Although we will not discuss it in much detail, Windows 2000 is capable of running on a symmetric multiprocessor system. This requirement means that the operating system code must be fully reentrant, that is, every procedure must be written in such a way that two or more CPUs may be changing its variables at once, without causing problems. In many cases this means that code sections have to be protected by spin locks or mutexes to keep additional CPUs at bay until the first one is done (i.e., serialize access to critical regions). The number of CPUs the system can handle is governed by the licensing restrictions and listed in Fig. 11-2. There is no technical reason why Windows Professional cannot run on a 32-node multiprocessor—it is the same binary as Datacenter Server, after all.

The upper limit of 32 CPUs is a hard limit because word-length bitmaps are used to keep track of CPU usage in various ways. For example, one word-length



**Figure 11-16.** The relationship between jobs, processes, and threads. Several fibers can also be multiplexed on one thread (not shown).

bitmap keeps track of which of the (up to) 32 CPUs are currently idle, and another bitmap is used per process to list the CPUs this process is permitted to run on. The 64-bit version of Windows 2000 should be able to effortlessly support up to 64 CPUs; beyond that requires actually changing the code substantially (to use multiple words for the bitmaps).

### 11.4.2 Job, Process, Thread and Fiber Management API Calls

New processes are created using the Win32 API function `CreateProcess`. This function has 10 parameters, each of which has many options. This design is clearly a lot more complicated than the UNIX scheme, in which `fork` has no parameters, and `exec` has just three: pointers to the name of the file to execute, the (parsed) command line parameter array, and the environment strings. Roughly speaking, the 10 parameters to `CreateProcess` are as follows:

1. A pointer to the name of the executable file.
2. The command line itself (unparsed).
3. A pointer to a security descriptor for the process.
4. A pointer to a security descriptor for the initial thread.
5. A bit telling whether the new process inherits the creator's handles.
6. Miscellaneous flags (e.g., error mode, priority, debugging, consoles).
7. A pointer to the environment strings.

8. A pointer to the name of the new process' current working directory.
9. A pointer to a structure describing the initial window on the screen.
10. A pointer to a structure that returns 18 values to the caller.

Windows 2000 does not enforce any kind of parent-child or other hierarchy. All processes are created equal (no processes are created more equal). However, since 1 of the 18 parameters returned to the creating process is a handle to the new process (allowing considerable control over the new process), there is an implicit hierarchy in terms of who has a handle to whom. Although these handles cannot just be passed directly to other processes, there is a way for a process to make a duplicate handle suitable for another process and then give it the handle, so the implicit process hierarchy may not last long.

Each process in Windows 2000 is created with a single thread, but a process can create more threads later on. Thread creation is simpler than process creation: `CreateThread` has only six parameters instead of 10:

1. The optional security descriptor.
2. The initial stack size.
3. The starting address.
4. A user-defined parameter.
5. The initial state of the thread (ready or blocked).
6. The thread's ID.

The kernel does the thread creation, so it is clearly aware of threads (i.e., they are not implemented purely in user space as is the case in some other systems).

## Interprocess Communication

Threads can communicate in a wide variety of ways, including pipes, named pipes, mailslots, sockets, remote procedure calls, and shared files. Pipes have two modes: byte and message, selected at creation time. Byte-mode pipes work the same way as in UNIX. Message-mode pipes are somewhat similar but preserve message boundaries, so that four writes of 128 bytes will be read as four 128-byte messages, and not as one 512-byte message, as might happen with byte-mode pipes. Named pipes also exist and have the same two modes as regular pipes. Named pipes can also be used over a network; regular pipes cannot.

**Mailslots** are a feature of Windows 2000 not present in UNIX. They are similar to pipes in some ways, but not all. For one thing, they are one-way, whereas pipes are two-way. They can also be used over a network but do not provide guaranteed delivery. Finally, they allow the sending process to broadcast a message to many receivers, instead of to just one receiver.

**Sockets** are like pipes, except that they normally connect processes on different machines. For example, one process writes to a socket and another one on a remote machine reads from it. Sockets can also be used to connect processes on the same machine, but since they entail more overhead than pipes, they are generally only used in a networking context.

Remote procedure calls are a way for process *A* to have process *B* call a procedure in *B*'s address space on *A*'s behalf and return the result to *A*. Various restrictions on the parameters exist. For example, it makes no sense to pass a pointer to a different process.

Finally, processes can share memory by mapping onto the same file at the same time. All writes done by one process then appear in the address spaces of the other processes. Using this mechanism, the shared buffer used in producer-consumer problems can easily be implemented.

Just as Windows 2000 provides numerous interprocess communication mechanisms, it also provides numerous synchronization mechanisms, including semaphores, mutexes, critical regions, and events. All of these mechanisms work on threads, not processes, so that when a thread blocks on a semaphore, other threads in that process (if any) are not affected and can continue to run.

A semaphore is created using the `CreateSemaphore` API function, which can initialize it to a given value and define a maximum value as well. Semaphores are kernel objects and thus have security descriptors and handles. The handle for a semaphore can be duplicated using `DuplicateHandle` and passed to another process so that multiple processes can synchronize on the same semaphore. Calls for up and down are present, although they have the somewhat peculiar names of `ReleaseSemaphore` (up) and `WaitForSingleObject` (down). It is also possible to give `WaitForSingleObject` a timeout, so the calling thread can be released eventually, even if the semaphore remains at 0 (although timers reintroduce races).

Mutexes are also kernel objects used for synchronization, but simpler than semaphores because they do not have counters. They are essentially locks, with API functions for locking (`WaitForSingleObject`) and unlocking (`ReleaseMutex`). Like semaphore handles, mutex handles can be duplicated and passed between processes so that threads in different processes can access the same mutex.

The third synchronization mechanism is based on **critical sections**, (which we have called critical regions elsewhere in this book) which are similar to mutexes, except local to the address space of the creating thread. Because critical sections are not kernel objects, they do not have handles or security descriptors and cannot be passed between processes. Locking and unlocking is done with `EnterCriticalSection` and `LeaveCriticalSection`, respectively. Because these API functions are performed initially in user space and only make kernel calls when blocking is needed, they are faster than mutexes.

The last synchronization mechanism uses kernel objects called **events**, of which there are two kinds: **manual-reset events** and **auto-reset events**. Any event can be in one of two states: set and cleared. A thread can wait for an event

to occur with `WaitForSingleObject`. If another thread signals an event with `SetEvent`, what happens depends on the type of event. With a manual-reset event, all waiting threads are released and the event stays set until manually cleared with `ResetEvent`. With an auto-reset event, if one or more threads are waiting, exactly one thread is released and the event is cleared. An alternative operation is `PulseEvent`, which is like `SetEvent` except that if nobody is waiting, the pulse is lost and the event is cleared. In contrast, a `SetEvent` that occurs with no waiting threads is remembered by leaving the event in set state so a subsequent thread waiting on it is released immediately.

Events, mutexes, and semaphores can all be named and stored in the file system, like named pipes. Two or more processes can synchronize by opening the same event, mutex, or semaphore, rather than having one of them create the object and then make duplicate handles for the others, although the latter approach is certainly an option as well.

The number of Win32 API calls dealing with processes, threads, and fibers is nearly 100, a substantial number of which deal with IPC in one form or another. A summary of the ones discussed above as well as some other important ones are given in Fig. 11-17.

Most of the calls in Fig. 11-17 were either discussed above or should be self-explanatory. Again note that not all of these are system calls. As we mentioned earlier, Windows 2000 knows nothing about fibers. They are entirely implemented in user space. As a consequence, the `CreateFiber` call does its work entirely in user space without making any system calls (unless it has to allocate some memory). Many other Win32 calls have this property as well, including `EnterCriticalSection` and `LeaveCriticalSection` as we noted above.

### 11.4.3 Implementation of Processes and Threads

Processes and threads are more important and more elaborate than jobs and fibers, so we will concentrate on them here. A process is created when another process makes the Win32 `CreateProcess` call. This call invokes a (user-mode) procedure in *kernel32.dll* that creates the process in several steps using multiple system calls and other work.

1. The executable file given as a parameter is examined and opened. If it is a valid POSIX, OS/2, 16-bit Windows, or MS-DOS file, a special environment is set up for it. If it is a valid 32-bit Win32 *.exe* file, the registry is checked to see if it is special in some way (e.g., to be run under supervision of a debugger). All of this is done in user mode inside *kernel32.dll*.
2. A system call, `NtCreateProcess`, is made to create the empty process object and enter it into the object manager's name space. Both the kernel object and the executive object are created. In addition, the



Win32 API Function	Description
CreateProcess	Create a new process
CreateThread	Create a new thread in an existing process
CreateFiber	Create a new fiber
ExitProcess	Terminate current process and all its threads
ExitThread	Terminate this thread
ExitFiber	Terminate this fiber
SetPriorityClass	Set the priority class for a process
SetThreadPriority	Set the priority for one thread
CreateSemaphore	Create a new semaphore
CreateMutex	Create a new mutex
OpenSemaphore	Open an existing semaphore
OpenMutex	Open an existing mutex
WaitForSingleObject	Block on a single semaphore, mutex, etc.
WaitForMultipleObjects	Block on a set of objects whose handles are given
PulseEvent	Set an event to signaled then to nonsignaled
ReleaseMutex	Release a mutex to allow another thread to acquire it
ReleaseSemaphore	Increase the semaphore count by 1
EnterCriticalSection	Acquire the lock on a critical section
LeaveCriticalSection	Release the lock on a critical section

**Figure 11-17.** Some of the Win32 calls for managing processes, threads, and fibers.

process manager creates a process control block for the object and initializes it with the process ID, quotas, access token, and various other fields. A section object is also created to keep track of the process' address space.

- When *kernel32.dll* gets control back, it makes another system call, *NtCreateThread*, to create the initial thread. The thread's user and kernel stacks are also created. The stack size is given in the header of the executable file.
- Kernel32.dll* now sends a message to the Win32 environment subsystem telling it about the new process and passing it the process and thread handles. The process and threads are entered into the subsystems tables so it has a complete list of all processes and threads. The subsystem then displays a cursor containing a pointer with an hourglass to tell the user that something is going on but that the cursor can be used in the meanwhile. When the process makes its first GUI

call, usually to create a window, the cursor is removed (it times out after 2 seconds if no call is forthcoming).

5. At this point, the thread is able to run. It starts out by running a run-time system procedure to complete the initialization.
6. The run-time procedure sets the thread's priority, tells loaded DLLs that a new thread is present, and does other housekeeping chores. Finally, it begins running the code of the process' main program.

Thread creation also consists of a number of steps, but we will not go into them in much detail. It starts when the running process executes `CreateThread`, which calls a procedure inside *kernel32.dll*. This procedure allocates a user stack within the calling process and then makes the `NtCreateThread` call to create an executive thread object, initialize it, and also create and initialize a thread control block. Again, the Win32 subsystem is notified and enters the new thread in its tables. Then the thread starts running and completes its own initialization.

When a process or thread is created, a handle is returned for it. This handle can be used to start, stop, kill, and inspect the process or thread. It is possible for the owner of a handle to pass the handle to another process in a controlled and secure way. This technique is used to allow debuggers to have control over the processes they are debugging.

## Scheduling

Windows 2000 does not have a central scheduling thread. Instead, when a thread cannot run any more, the thread enters kernel mode and runs the scheduler itself to see which thread to switch to. The following conditions cause the currently running thread to execute the scheduler code:

1. The thread blocks on a semaphore, mutex, event, I/O, etc.
2. It signals an object (e.g., does an up on a semaphore).
3. The running thread's quantum expires.

In case 1, the thread is already running in kernel mode to carry out the operation on the dispatcher or I/O object. It cannot possibly continue, so it must save its own context, run the scheduler code to pick its successor, and load that thread's context to start it.

In case 2, the running thread is in the kernel, too. However, after signaling some object, it can definitely continue because signaling an object never blocks. Still, the thread is required to run the scheduler to see if the result of its action has released a higher priority thread that is now free to run. If so, a thread switch occurs because Windows 2000 is fully preemptive (i.e., thread switches can occur at any moment, not just at the end of the current thread's quantum).

In case 3, a trap to kernel mode occurs, at which time the thread executes the scheduler code to see who runs next. Depending on what other threads are waiting, the same thread may be selected, in which case it gets a new quantum and continues running. Otherwise a thread switch happens.

The scheduler is also called under two other conditions:

1. An I/O operation completes.
2. A timed wait expires.

In the first case, a thread may have been waiting on this I/O and is now released to run. A check has to be made to see if it should preempt the running thread since there is no guaranteed minimum run time. The scheduler is not run in the interrupt handler itself (since that may keep interrupts turned off too long). Instead a DPC is queued for slightly later, after the interrupt handler is done. In the second case, a thread has done a down on a semaphore or blocked on some other object, but with a timeout that has now expired. Again it is necessary for the interrupt handler to queue a DPC to avoid having it run during the clock interrupt handler. If a thread has been made ready by this timeout, the scheduler will be run and if nothing more important is available, the DPC will run next.

Now we come to the actual scheduling algorithm. The Win32 API provides two hooks for processes to influence thread scheduling. These hooks largely determine the algorithm. First, there is a call `SetPriorityClass` that sets the priority class of all the threads in the caller's process. The allowed values are: realtime, high, above normal, normal, below normal, and idle.

Second, there is a call `SetThreadPriority` that sets the relative priority of some thread (possibly, but not necessarily, the calling thread) compared to the other threads in its process. The allowed values are: time critical, highest, above normal, normal, below normal, lowest, and idle. With six process classes and seven thread classes, a thread can have any one of 42 combinations. This is the input to the scheduling algorithm.

The scheduler works as follows. The system has 32 priorities, numbered from 0 to 31. The 42 combinations are mapped onto the 32 priority classes according to the table of Fig. 11-18. The number in the table determines the thread's **base priority**. In addition, every thread has a **current priority**, which may be higher (but not lower) than the base priority and which we will discuss shortly.

To use these priorities for scheduling, the system maintains an array with 32 entries, corresponding to priorities 0 through 31 derived from the table of Fig. 11-18. Each array entry points to the head of a list of ready threads at the corresponding priority. The basic scheduling algorithm consists of searching the array from priority 31 down to priority 0. As soon as a nonempty slot is found, the thread at the head of the queue is selected and run for one quantum. If the quantum expires, the thread goes to the end of the queue at its priority level and the thread at the front is chosen next. In other words, when there are multiple

		Win32 process class priorities					
Win32 thread priorities		Realtime	High	Above Normal	Normal	Below Normal	Idle
	Time critical	31	15	15	15	15	15
	Highest	26	15	12	10	8	6
	Above normal	25	14	11	9	7	5
	Normal	24	13	10	8	6	4
	Below normal	23	12	9	7	5	3
	Lowest	22	11	8	6	4	2
	Idle	16	1	1	1	1	1

Figure 11-18. Mapping of Win32 priorities to Windows 2000 priorities.

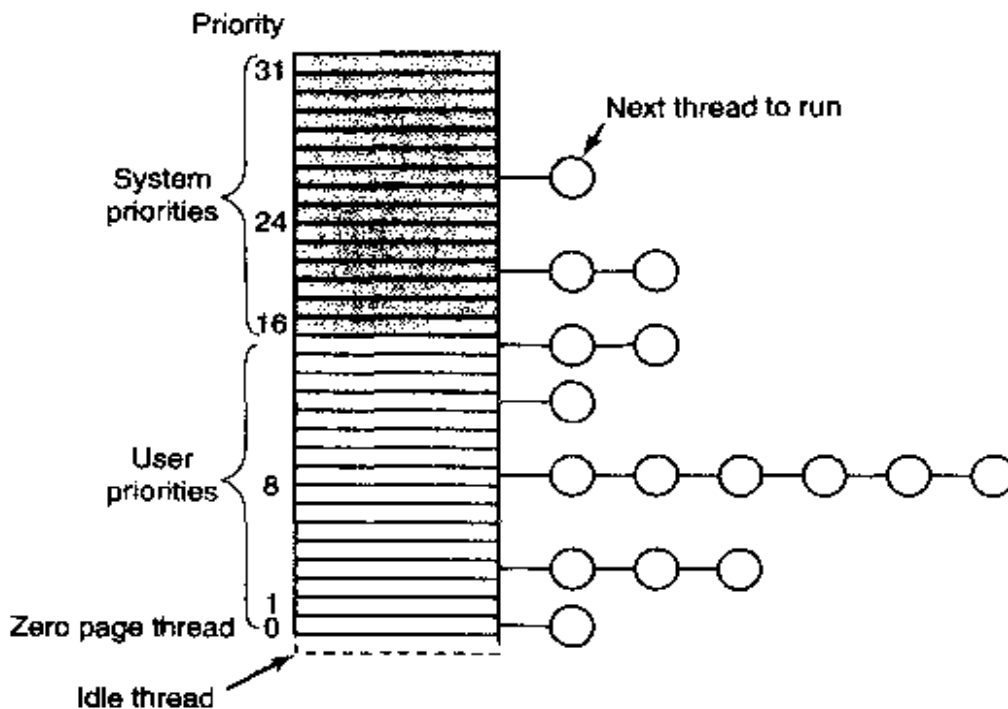
threads ready at the highest priority level, they run round robin for one quantum each. If no thread is ready, the idle thread is run.

It should be noted that scheduling is done by picking a thread without regard to which process that thread belongs. Thus the scheduler does *not* first pick a process and then pick a thread in that process. It only looks at the threads. It does not even know which thread belongs to which process. On a multiprocessor, each CPU schedules itself using the priority array. A spin lock is used to make sure that only one CPU at a time is inspecting the array.

The array of queue headers is shown in Fig. 11-19. The figure shows that there are actually four categories of priorities: realtime, user, zero, and idle, which is effectively  $-1$ . These deserve some comment. Priorities 16–31 are called real time, but they are not. There are no guarantees given and no deadlines are met. They are simply higher priority than 0–15. However, priorities 16 through 31 are reserved for the system itself and for threads explicitly assigned those priorities by the system administrator. Ordinary users may not run there for a good reason. If a user thread were to run at a higher priority than, say, the keyboard or mouse thread and get into a loop, the keyboard or mouse thread would never run, effectively hanging the system.

User threads run at priorities 1–15. By setting the process and thread priorities, a user can determine which threads get preference. The zero thread runs in the background and eats up whatever CPU time nobody else wants. Its job is to zero pages for the memory manager. We will discuss its role later. If there is absolutely nothing to do, not even zero pages, the idle thread runs. It is not really a full-blown thread though.

Over the course of time, some patches were made to the basic scheduling algorithm to improve system performance. Under certain specific conditions, the current priority of a user thread can be raised above the base priority (by the operating system), but never above priority 15. Since the array of Fig. 11-19 is



**Figure 11-19.** Windows 2000 supports 32 priorities for threads.

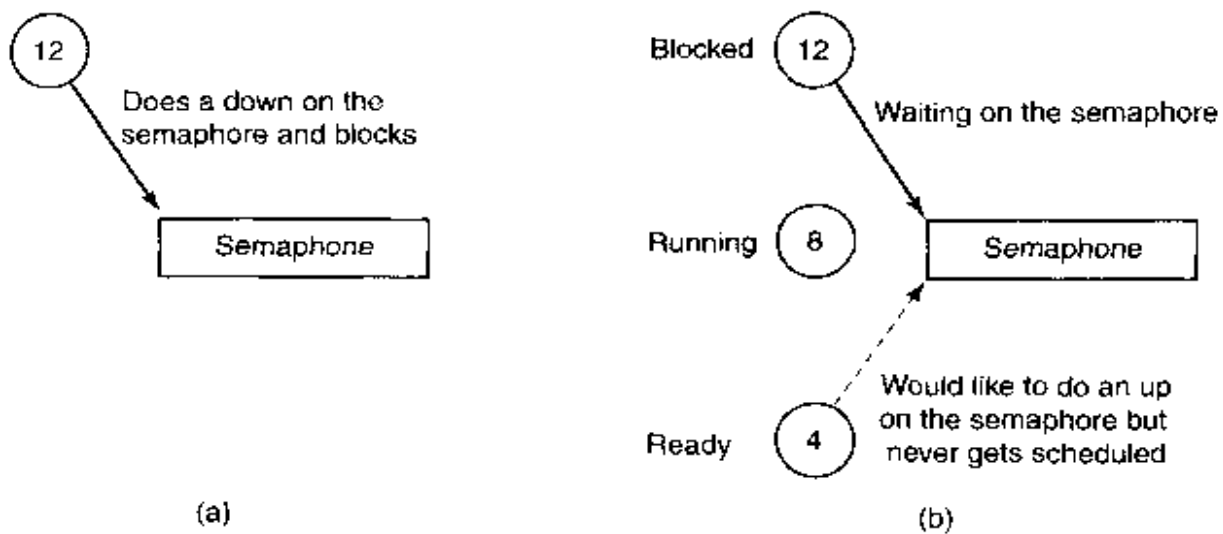
based on the current priority, changing this priority affects scheduling. No adjustments are ever made to threads running at priority 15 or higher.

Let us now see when a thread's priority is raised. First, when an I/O operation completes and releases a waiting thread, the priority is boosted to give it a chance to run again quickly and start more I/O. The idea here is to keep the I/O devices busy. The amount of boost depends on the I/O device, typically 1 for a disk, 2 for a serial line, 6 for the keyboard, and 8 for the sound card.

Second, if a thread was waiting on a semaphore, mutex, or other event, when it is released, it gets boosted by 2 units if it is in the foreground process (the process controlling the window to which keyboard input is sent) and 1 unit otherwise. This fix tends to raise interactive processes above the big crowd at level 8. Finally, if a GUI thread wakes up because window input is now available, it gets a boost for the same reason.

These boosts are not forever. They take effect immediately, but if a thread uses all of its next quantum, it loses one point and moves down one queue in the priority array. If it uses up another full quantum, it moves down another level, and so on until it hits its base level, where it remains until it is boosted again. Clearly, if a thread wants good service, it should play a lot of music.

There is one other case in which the system fiddles with the priorities. Imagine that two threads are working together on a producer-consumer type problem. The producer's work is harder, so it gets a high priority, say 12, compared to the consumer's 4. At a certain point, the producer has filled up a shared buffer and blocks on a semaphore, as illustrated in Fig. 11-20(a).



**Figure 11-20.** An example of priority inversion.

Before the consumer gets a chance to run again, an unrelated thread at priority 8 becomes ready and starts running, as shown in Fig. 11-20(b). As long as this thread wants to run, it will be able to, since it outguns the consumer and the producer, while higher, is blocked. Under these circumstances, the producer will never get to run again until the priority 8 thread gives up.

Windows 2000 solves this problem through what might be charitably called a big hack. The system keeps track of how long it has been since a ready thread ran last. If it exceeds a certain threshold, it is moved to priority 15 for two quanta. This may give it the opportunity to unblock the producer. After the two quanta are up, the boost is abruptly removed rather than decaying gradually. Probably a better solution would be to penalize threads that use up their quantum over and over by lowering their priority. After all, the problem was not caused by the starved thread, but by the greedy thread. This problem is well known under the name **priority inversion**.

An analogous problem happens if a priority 16 thread grabs a mutex and does not get a chance to run for a long time, starving more important system threads that are waiting for the mutex. This problem can be prevented within the operating system by having a thread that needs a mutex for a short time just disable scheduling while it is busy. On a multiprocessor, a spin lock should be used.

Before leaving the subject of scheduling, it is worth saying a couple of words about the quantum. On Windows 2000 Professional the default is 20 msec; on uniprocessor servers it is 120 msec; on multiprocessors various other values are used, depending on the clock frequency. The shorter quantum favors interactive users whereas the longer quantum reduces context switches and thus provides better efficiency. This is the meaning of the last column in Fig. 11-2. These defaults can be increased manually by 2x, 4x, or 6x if desired. As an aside, the size of the quantum was chosen a decade ago and not changed since although machines are now more than an order of magnitude faster. The numbers probably

could be reduced by a factor of 5 to 10 with no harm and possibly better response time for interactive threads in a heavily loaded system.

One last patch to the scheduling algorithm says that when a new window becomes the foreground window, all of its threads get a longer quantum by an amount taken from the registry. This change gives them more CPU time, which usually translates to better service for the window that just moved to the foreground.

#### 11.4.4 MS-DOS Emulation

One of the design goals of Windows 2000 was inherited from NT: try to run as many reasonable MS-DOS programs as possible. This goal is quite different from Windows 98's stated goal: run all old MS-DOS programs (to which we add: no matter how ill-behaved they may be).

The way Windows 2000 deals with ancient programs is to run them in a fully protected environment. When an MS-DOS program is started, a normal Win32 process is started and loaded with an MS-DOS emulation program, *ntvdm* (NT Virtual DOS Machine) that will monitor the MS-DOS program and carry out its system calls. Since MS-DOS only recognized memory up to 1 MB on the 8088 and only up to 16 MB with bank switching and other tricks on the 286, it is safe to put *ntvdm* high in the process' virtual address space where the program has no way to address it. This situation is shown in Fig. 11-21.

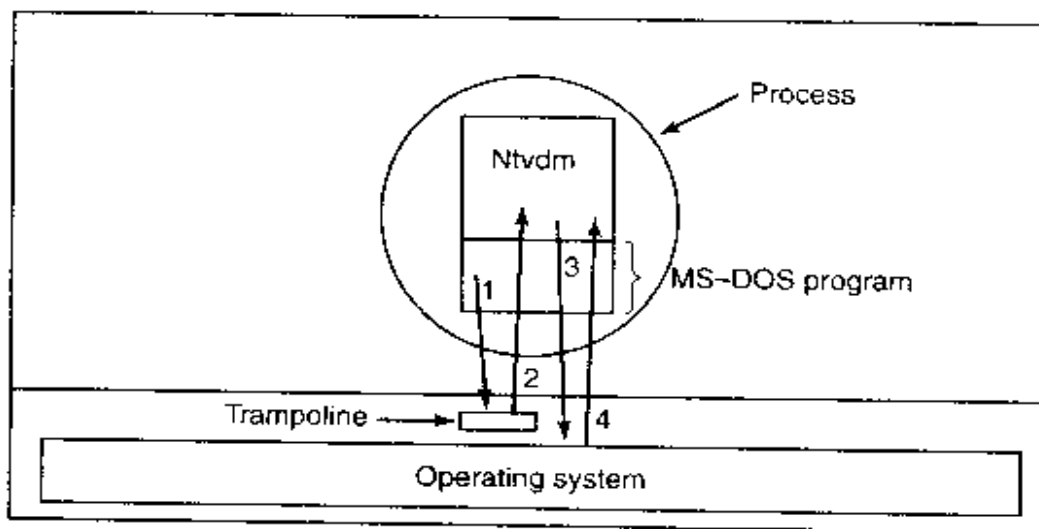


Figure 11-21. How old MS-DOS programs are run under Windows 2000.

When the MS-DOS program is just executing normal instructions, it can run on the bare hardware since the Pentium includes all the 8088 and 286 instructions as subsets. The interesting part is what happens when the MS-DOS program wants to do I/O or interact with the system. A well-behaved program just makes a system call. In expectation of this, *ntvdm* instructs Windows 2000 to reflect all MS-DOS system calls back to it. In effect, the system call just bounces off the operating

system and is caught by the emulator, as shown in steps 1 and 2 in Fig. 11-21. Sometimes this technique is referred to as using a **trampoline**.

Once it gets control, the emulator figures out what the program was trying to do and issues its own Win32 calls to get the work done (step 3 and 4 in Fig. 11-21). As long as the program is well behaved and just makes legal MS-DOS system calls, this technique works fine. The trouble is that some old MS-DOS programs bypassed the operating system and wrote directly to the video RAM, read directly from the keyboard, and so on, things that are impossible in a protected environment. To the extent that the illegal behavior causes a trap, there is some hope that the emulator can figure out what the program was trying to do and emulate it. If it does not know what the program wants, the program is just killed because 100 percent emulation was not a Windows 2000 requirement.

### 11.4.5 Booting Windows 2000

Before Windows 2000 can start up, it must be booted. The boot process creates the initial processes that bring up the system. In this section we will briefly discuss how the boot process works for Windows 2000. The hardware boot process consists of reading in the first sector of the first disk (the master boot record) and jumping to it, as we described in Sec. 5.4.2. This short assembly language program reads the partition table to see which partition contains the bootable operating system. When it finds the operating system partition, it reads in the first sector of that partition, called the **boot sector**, and jumps to it. The program in the boot sector reads its partition's root directory, searching for a file called *ntldr* (another piece of archaeological evidence that Windows 2000 is really NT). If it finds that file, it reads the file into memory and executes it. *Ntldr* loads Windows 2000. As an aside, there are several versions of the boot sector, depending on whether the partition is formatted as FAT-16, FAT-32, or NTFS. When Windows 2000 is installed, the correct version of the master boot record and boot sector are written to disk.

*Ntldr* now reads a file called *Boot.ini*, which is the only configuration information not in the registry. It lists all the versions of *hal.dll* and *ntoskrnl.exe* available for booting in this partition. The file also provides many parameters, such as how many CPUs and how much RAM to use, whether to give user processes 2 GB or 3 GB, and what rate to set the real-time clock to. *Ntldr* then selects and loads *hal.dll* and *ntoskrnl.exe* files as well as *bootvid.dll*, the default video driver for writing on the display during the boot process. *Ntldr* next reads the registry to find out which drivers are needed to complete the boot (e.g., the keyboard and mouse drivers, but also dozens more for controlling various chips on the parent-board). Finally, it reads in all these drivers and passes control to *ntoskrnl.exe*.

Once started, the operating system does some general initialization and then calls the executive components to do their own initialization. For example, the object manager prepares its name space to allow other components call it to insert



their objects into the name space. Many components also do specific things related to their function, such as the memory manager setting up the initial page tables and the plug-and-play manager finding out which I/O devices are present and loading their drivers. All in all, dozens of steps are involved, during which time the progress bar displayed on the screen is growing in length as steps are completed. The last step is creating the first true user process, the **session manager**, *smss.exe*. Once this process is up and running, booting is completed.

The session manager is a native Windows 2000 process. It makes true system calls and does not use the Win32 environment subsystem, which is not even running yet. In fact, one of its first duties is to start it (*csrss.exe*). It also reads the registry hives from disk and learns what else it is supposed to do. Typically its work includes entering many objects in the object manager's name space, creating any extra paging files needed, and opening important DLLs to have them around all the time. After it has done most of this work, it creates the login daemon, *winlogon.exe*.

At this point, the operating system is up and running. Now it is time to get the service processes (user space daemons) going and allow users to log in. *Winlogon.exe* first creates the authentication manager (*lsass.exe*), and then the parent process of all the services (*services.exe*). The latter looks in the registry to find out which user space daemon processes are needed and what files they are in. It then starts creating them. They are shown in Fig. 11-7. The fact that the disk is generally being heavily used after the first user has logged in (but has done nothing) is not the user's fault. The culprit is *services.exe* creating all the services. In addition, it also loads any remaining device drivers that have not yet been loaded. The hierarchy of initial processes and some typical services are shown in Fig. 11-22.

*Winlogon.exe* is also responsible for all user logins. The actual login dialog is handled by a separate program in *msgina.dll* to make it possible for third parties to replace the standard login with faceprint identification or something else other than name and password. After a successful login, *winlogon.exe* gets the user's profile from the registry and from it determines which shell to run. Many people do not realize it, but the standard Windows desktop is just *explorer.exe* with some options set. If desired, a user can select any other program as the shell, including the command prompt or even *Word*, by editing the registry. However, editing the registry is not for the faint of heart: a mistake here can make the system unusable.

## 11.5 MEMORY MANAGEMENT

Windows 2000 has an extremely sophisticated virtual memory system. It has a number of Win32 functions for using it and part of the executive plus six dedicated kernel threads for managing it. In the following sections we will look at the fundamental concepts, the Win32 API calls, and finally the implementation.

Process	Description
idle	Not really a process, but home to the idle thread
system	Creates smss.exe & paging files; reads registry; opens DLLs
smss.exe	First real proc; much initialization; creates csrss & winlogon
csrss.exe	Win32 subsystem process
winlogon.exe	Login daemon
lsass.exe	Authentication manager
services.exe	Looks in registry and starts services
Printer server	Allows remote jobs to use the printer
File server	Serves requests for local files
Telnet daemon	Allows remote logins
Incoming email handler	Accepts and stores inbound email
Incoming fax handler	Accepts and prints inbound faxes
DNS resolver	Internet domain name system server
Event logger	Logs various system events
Plug-and-play manager	Monitors hardware to see what is out there

**Figure 11-22.** The processes starting up during the boot phase. The ones above the line are always started. The ones below it are examples of services that could be started.

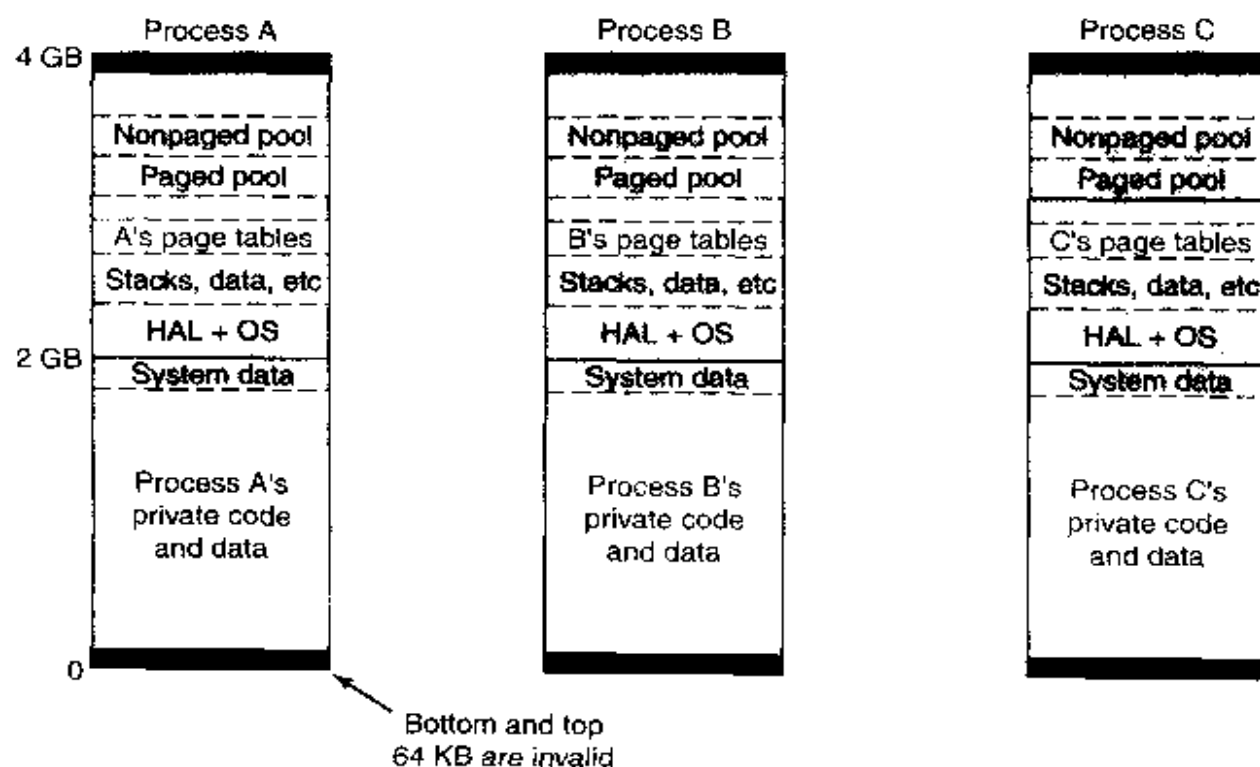
### 11.5.1 Fundamental Concepts

In Windows 2000, every user process has its own virtual address space. Virtual addresses are 32 bits long, so each process has 4 GB of virtual address space. The lower 2 GB minus about 256 MB are available for the process' code and data; the upper 2 GB map onto kernel memory in a protected way. The virtual address space is demand paged, with a fixed page size (4 KB on the Pentium).

The virtual address space layout for three user processes is shown in Fig. 11-23 in slightly simplified form. The bottom and top 64 KB of each process' virtual address space is normally unmapped. This choice was made intentionally to help catch programming errors. Invalid pointers are often 0 or -1, so attempts to use them on Windows 2000 will cause an immediate trap instead of reading garbage or, worse yet, writing to an incorrect memory location. However, when old MS-DOS programs are being run in emulation mode, they can be mapped in.

Starting at 64 KB comes the user's private code and data. This extends up to almost 2 GB. The last piece of the bottom 2 GB contains some system counters and timers that are shared among all users read only. Making them visible here allows processes to access them without the overhead of a system call.

The upper 2 GB contains the operating system, including the code, data, and the paged and nonpaged pools (used for objects, etc.). The upper 2 GB is shared among all user processes, except for the page tables, which are each process' own page tables. The upper 2 GB of memory is not writable and mostly not even



**Figure 11-23.** Virtual address space layout for three user processes. The white areas are private per process. The shaded areas are shared among all processes.

readable for user-mode processes. The reason for putting it here is that when a thread makes a system call, it traps into kernel mode and just keeps on running in the same thread. By making the whole operating system and all of its data structures (as well as the whole user process) visible within a thread's address space when it enters kernel mode, there is no need to change the memory map or flush the cache upon kernel entry. All that has to be done is switch to the thread's kernel stack. The trade-off here is less private address space per process in return for faster system calls. Large database servers already feel cramped, which is why the 3-GB user space option is available on Advanced Server and Datacenter Server.

Each virtual page can be in one of three states: free, reserved, or committed. A **free page** is not currently in use and a reference to it causes a page fault. When a process is started, all of its pages are in free state until the program and initial data are mapped into its address space. Once code or data is mapped onto a page, the page is said to be **committed**. A reference to a committed page is mapped using the virtual memory hardware and succeeds if the page is in main memory. If the page is not in main memory, a page fault occurs and the operating system finds and brings in the page from disk.

A virtual page can also be in **reserved** state, meaning it is not available for being mapped until the reservation is explicitly removed. For example, when a

new thread is created, 1 MB of stack space is reserved in the process' virtual address space, but only one page is committed. This technique means that the stack can eventually grow to 1 MB without fear that some other thread will allocate the needed contiguous piece of virtual address space out from under it. In addition to the free, reserved, and committed attributes, pages also have other attributes, such as being readable, writable, and executable.

An interesting trade-off occurs with assignment of backing store to committed pages. A simple strategy would be to assign a page in one of the paging files to back up each committed page at the time the page was committed. This would guarantee that there was always a known place to write out each committed page should it be necessary to evict it from memory. The downside of this strategy is that the paging file might have to be as large as the union of all processes' virtual memory. On a large system that rarely ran out of memory and thus rarely paged, this approach would waste disk space.

To avoid wasting disk space, Windows 2000 committed pages that have no natural home on the disk (e.g., stack pages) are not assigned a disk page until the moment that they have to be paged out. This design makes the system more complex because the paging files maps may have to be fetched during a page fault, and fetching them may cause one or more additional page faults inside the page fault handler. On the other hand, no disk space need be allocated for pages that are never paged out.

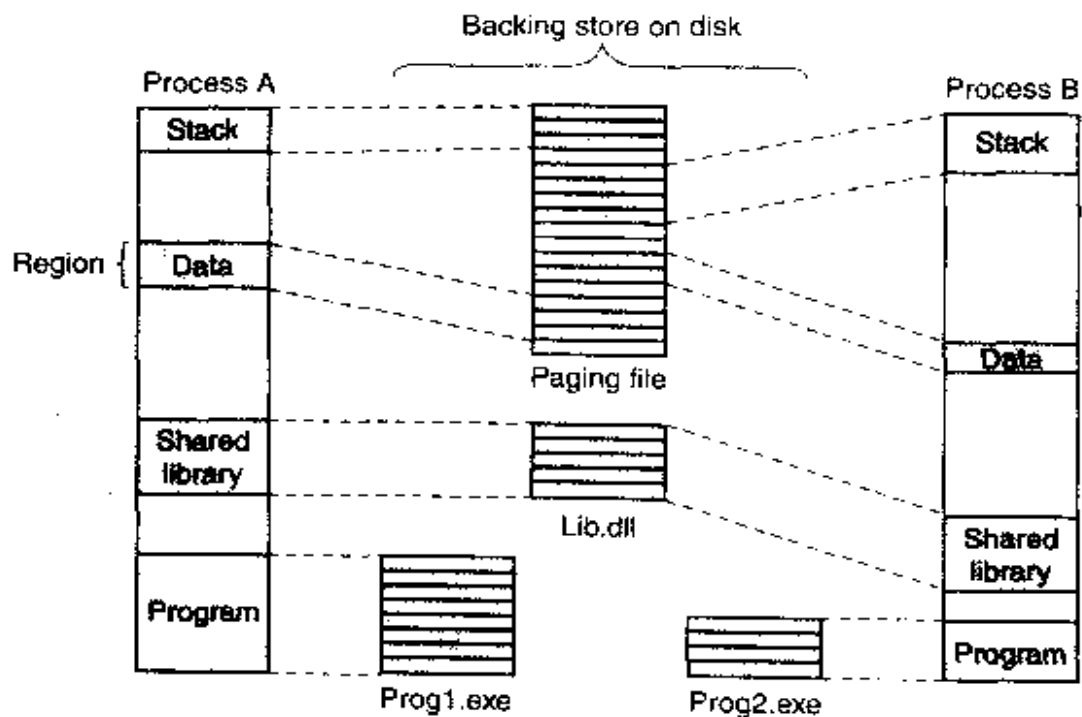
Trade-offs like this (system complexity versus better performance or more features) tend to get resolved in favor of the latter because the value of better performance or more features is clear but the downside of complexity (a bigger maintenance headache and more crashes per year) is hard to quantify. Free and reserved pages never have shadow pages on disk and references to them always cause page faults.

The shadow pages on the disk are arranged into one or more paging files. There may be up to 16 paging files, possibly spread over 16 separate disks, for higher I/O bandwidth. Each one has an initial size and a maximum size it can grow to later if needed. These files can be created at the maximum size at system installation time in order to reduce the chances that they are highly fragmented, but new ones can be created using the control panel later on. The operating system keeps track of which virtual page maps onto which part of which paging file. For (execute only) program text, the executable binary file (i.e., *.exe* or *.dll* file) contains the shadow pages; for data pages, the paging files are used.

Windows 2000, like many versions of UNIX, allows files to be mapped directly onto regions of the virtual address spaces (i.e., runs of consecutive pages). Once a file has been mapped onto the address space, it can be read or written using ordinary memory references. Memory-mapped files are implemented in the same way as other committed pages, only the shadow pages are in the user's file instead of in the paging file. As a result, while a file is mapped in, the version in memory may not be identical to the disk version (due to recent writes to the vir-

tual address space). However, when the file is unmapped or is explicitly flushed, the disk version is brought up to date.

Windows 2000 explicitly allows two or more processes to map onto the same part of the same file at the same time, possibly at different virtual addresses, as shown in Fig. 11-24. By reading and writing memory words, the processes can now communicate with each other and pass data back and forth at very high bandwidth, since no copying is required. Different processes may have different access permissions. Since all the processes using a mapped file share the same pages, changes made by one of them are immediately visible to all the others, even if the disk file has not yet been updated. Care is also taken that if another process opens the file for normal reading, it sees the current pages in RAM, not stale pages from the disk.



**Figure 11-24.** Mapped regions with their shadow pages on disk. The *lib.dll* file is mapped into two address spaces at the same time.

It is worth noting that there is a problem if two programs share a DLL file and one of them changes the file's static data. If no special action is taken, the other one will see the changed data, which is probably not what is desired. The problem is solved by mapping all pages in as read only by secretly noting that some are really writable. When a write happens to a page that is mapped read only but is really writable, a private copy of the page is made and mapped in. Now it can be written safely without affecting other users or the original copy on disk. This technique is called **copy-on-write**.

Also it is worth noting that if program text is mapped into two address spaces at different addresses, a certain problem arises with addressing. What happens if the first instruction is `JMP 300`? If process one maps the program in at address

65,536, the code can easily be patched to read JMP 65836. But what happens if a second process now maps it in at 131,072? The JMP 65836 will go to address 65,836 instead of 131,372 and the program will fail. The solution is to use only relative offsets, not absolute virtual addresses in code that is to be shared. Fortunately, most machines have instructions using relative offsets as well as instructions using absolute addresses. Compilers can use the relative offset instructions, but they have to know in advance whether to use them or the absolute ones. The relative ones are not used all the time because the resulting code is usually less efficient. Usually, a compiler flag tells them which to use. The technique of making it possible to place a piece of code at any virtual address without relocation is called **position independent code**.

Years ago, when 16-bit (or 20-bit) virtual address spaces were standard, but machines had megabytes of physical memory, all kinds of tricks were thought of to allow programs to use more physical memory than fit in the address space. Often these tricks went under the name of **bank switching**, in which a program could substitute some block of memory above the 16-bit or 20-bit limit for a block of its own memory. When 32-bit machines were introduced, people thought they would have enough address space forever. They were wrong. The problem is back. Large programs often need more than the 2 GB or 3 GB of user address space Windows 2000 allocates to them, so bank switching is back, now called **address windowing extensions**. This facility allows programs to map into shuffle chunks of memory in and out of the user portion of the address space (and especially above the dreaded 4-GB boundary). Since it is only used on servers with more than 2 GB of physical memory, we will defer the discussion until the next edition of this book (by which time even entry-level desktop machines will be feeling the 32-bit pinch).

### 11.5.2 Memory Management System Calls

The Win32 API contains a number of functions that allow a process to manage its virtual memory explicitly. The most important of these functions are listed in Fig. 11-25. All of them operate on a region consisting either of a single page or a sequence of two or more pages that are consecutive in the virtual address space.

The first four API functions are used to allocate, free, protect and query regions of virtual address space. Allocated regions always begin on 64-KB boundaries to minimize porting problems to future architectures with pages larger than current ones (up to 64 KB). The actual amount of address space allocated can be less than 64 KB, but must be a multiple of the page size. The next two give a process the ability to hardwire pages in memory so they will not be paged out and to undo this property. A real-time program might need this ability, for example. A limit is enforced by the operating system to prevent processes from getting too greedy. Actually, the pages can be removed from memory, but only if the entire

Win32 API function	Description
VirtualAlloc	Reserve or commit a region
VirtualFree	Release or decommit a region
VirtualProtect	Change the read/write/execute protection on a region
VirtualQuery	Inquire about the status of a region
VirtualLock	Make a region memory resident (i.e., disable paging for it)
VirtualUnlock	Make a region pageable in the usual way
CreateFileMapping	Create a file mapping object and (optionally) assign it a name
MapViewOfFile	Map (part of) a file into the address space
UnmapViewOfFile	Remove a mapped file from the address space
OpenFileMapping	Open a previously created file mapping object

**Figure 11-25.** The principal Win32 API functions for managing virtual memory in Windows 2000.

process is swapped out. When it is brought back, all the locked pages are reloaded before any thread can start running again. Although not shown in Fig. 11-25, Windows 2000 also has API functions to allow a process to access the virtual memory of a different process over which it has been given control (i.e., for which it has a handle).

The last four API functions listed are for managing memory-mapped files. To map a file, a file mapping object (see Fig. 11-10) must first be created, with `CreateFileMapping`. This function returns a handle to the file mapping object and optionally enters a name for it into the file system so another process can use it. The next two functions map and unmap files, respectively. The last one can be used by a process to map in a file currently also mapped in by a different process. In this way, two or more processes can share regions of their address spaces. This technique allows them to write in limited regions of each other's memory.

### 11.5.3 Implementation of Memory Management

Windows 2000 supports a single linear 4-GB demand-paged address space per process. Segmentation is not supported in any form. Theoretically, page sizes can be any power of two up to 64 KB. On the Pentium they are fixed at 4 KB; on the Itanium they can be 8 KB or 16 KB. In addition, the operating system itself can use 4-MB pages to reduce page table space consumed.

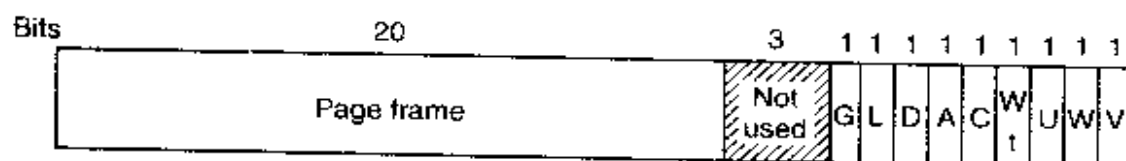
Unlike the scheduler, which selects individual threads to run and does not care much about processes, the memory manager deals entirely with processes and does not care much about threads. After all, processes, not threads, own the address space and that is what the memory manager deals with. When a region of virtual address space is allocated, as four of them have been for process *A* in

Fig. 11-24, the memory manager creates a **VAD (Virtual Address Descriptor)** for it, listing the range of addresses mapped, the backing store file and offset where it is mapped, and the protection code. When the first page is touched, the directory of page tables is created and a pointer to it inserted in the VAD. In fact, an address space is completely defined by the list of its VADs. This scheme supports sparse address spaces because unused areas between the mapped regions use no resources.

## Page Fault Handling

Windows 2000 does not use any form of prepaging. When a process starts, none of its pages are in memory. All of them are brought in dynamically as page faults occurs. On each page fault, a trap to the kernel (in the sense of Fig. 11-7) occurs. The kernel builds a machine-independent descriptor telling what happened and passes this to the memory manager part of the executive. The memory manager then checks it for validity. If the faulted page falls within a committed or reserved region, it looks up the address in the list of VADs, finds (or creates) the page table, and looks up the relevant entry.

The page table entries are different for different architectures. For the Pentium, the entry for a mapped page is shown in Fig. 11-26. Unmapped pages also have entries, but their format is somewhat different. For example, for an unmapped page that must be zeroed before it may be used, that fact is noted in the page table.



G: Page is global to all processes  
 L: Large (4-MB) page  
 D: Page is dirty  
 A: Page has been accessed  
 C: Caching enabled/disabled

Wt: Write through (no caching)  
 U: Page is accessible in user mode  
 W: Writing to the page permitted  
 V: Valid page table entry

**Figure 11-26.** A page table entry for a mapped page on the Pentium.

The most important bits in the page table entry for purposes of the paging algorithm are the *A* and *D* bits. They are fed by the hardware and keep track of whether the page has been referenced or written on, respectively, since the last time they were cleared.

Page faults come in five categories:

1. The page referenced is not committed.
2. A protection violation occurred.
3. A shared page has been written.



4. The stack needs to grow.
5. The page referenced is committed but not currently mapped in.

The first and second cases are fatal errors from which there is no recovery for the faulting process. The third case has the same symptoms as the second one (an attempt to write to a read-only page), but the treatment is different. The solution is to copy the page to a new physical page frame and map that one in read/write. This is how copy-on-write works. (If a shared page is marked writable in all processes using it, it is not copy-on-write and no fault occurs when writing to it.) The fourth case requires allocating a new page frame and mapping it in. However, the security rules require that the page contain only 0s, to prevent the process from snooping on the previous owner of the page. Thus a page of 0s must be found, or if one is not available, another page frame must be allocated and zeroed on the spot. Finally, the fifth case is a normal page fault. The page is located and mapped in.

The actual mechanics of getting and mapping pages is fairly standard, so we will not discuss this issue. The only noteworthy feature is that Windows 2000 does not read in isolated pages from the disk. Instead, it reads in runs of consecutive pages, usually about 1–8 pages, in an attempt to minimize the number of disk transfers. The run size is larger for code pages than for data pages.

### The Page Replacement Algorithm

Page replacement works like this. The system makes a serious attempt to maintain a substantial number of free pages in memory so that when a page fault occurs, a free page can be claimed on the spot, without the need to first write some other page to disk. As a consequence of this strategy, most page faults can be satisfied with at most one disk operation (reading in the page), rather than sometimes two (writing back a dirty page and then reading in the needed page).

Of course, the pages on the free list have to come from somewhere, so the real page replacement algorithm is how pages get taken away from processes and put on the free list (actually, there are four free lists, but for the moment it is simplest to think of there being just one; we will come to the details later). Let us now take a look at how Windows 2000 frees pages. To start with, the entire paging system makes heavy use of the working set concept. Each process (*not* each thread) has a working set. This set consists of the mapped-in pages that are in memory and can be thus referenced without a page fault. The size and composition of the working set fluctuates as the process' threads run, of course.

Each process' working set is described by two parameters: the minimum size and the maximum size. These are not hard bounds, so a process may have fewer pages in memory than its minimum or (under certain circumstances) more than its maximum. Every process starts with the same minimum and maximum, but these

bounds can change over time. The default initial minimum is in the range 20–50 and the default initial maximum is in the range 45–345, depending on the total amount of RAM. The system administrator can change these defaults, however.

If a page fault occurs and the working set is smaller than the minimum, the page is added. On the other hand, if a page fault occurs and the working set is larger than the maximum, a page is evicted from the working set (but not from memory) to make room for the new page. This algorithm means that Windows 2000 uses a local algorithm, to prevent one process from hurting others by hogging memory. However, the system does try to tune itself to some extent. For example, if it observes that one process is paging like crazy (and the others are not), the system may increase the size of its maximum working set, so that over time, the algorithm is a mix of local and global. There is an absolute limit on the working set size, however: even if there is only one process running, it may not take the last 512 pages, to leave some slack for new processes.

So far, so good, but the story is not over yet. Once a second, a dedicated kernel daemon thread, the **balance set manager**, checks to see if there are enough free pages. If there are not enough, it starts the **working set manager** thread to examine the working sets and recover more pages. The working set manager first determines the order to examine the processes in. Large processes that have been idle for a long time are considered before small active processes and the foreground process is considered last.

The working set manager then starts inspecting processes in the chosen order. If a process' working set is currently less than its minimum or it has incurred more than a certain number of page faults since the last inspection, it is passed over. Otherwise, one or more pages are removed. The target number of pages to remove is a complicated function of the total RAM size, how tight memory is, how the current working set size compares to the process' minimum and maximum, and other parameters. All the pages are examined in turn.

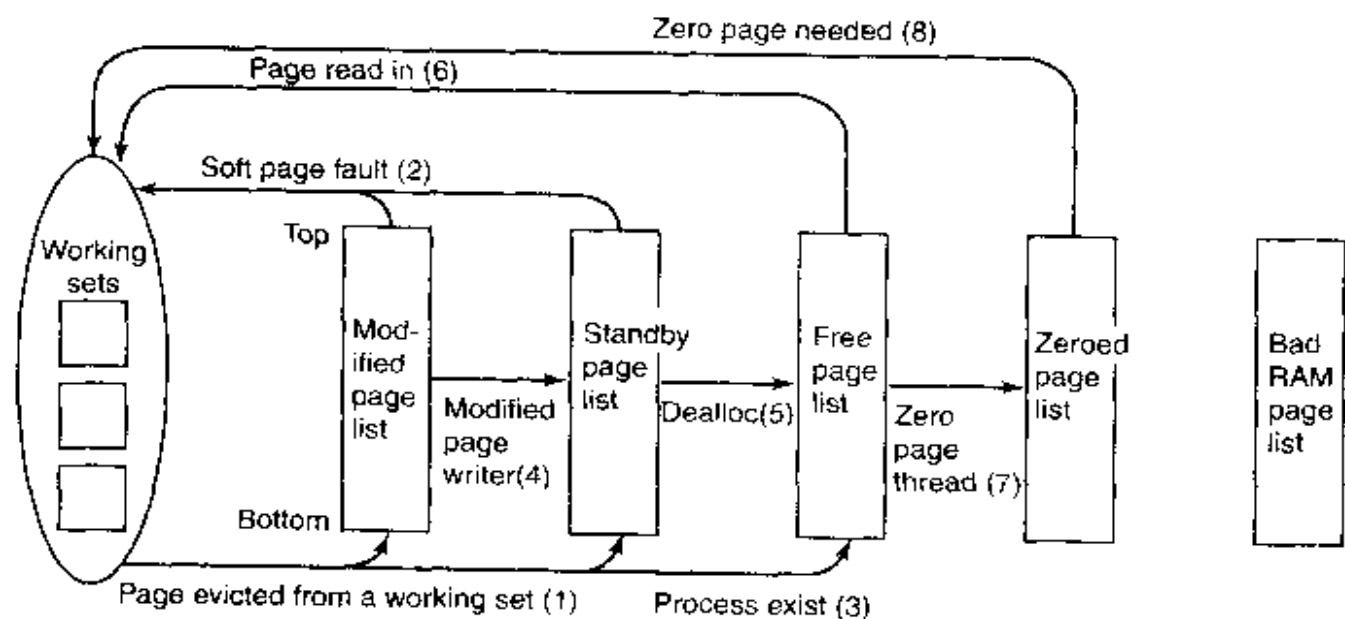
On a uniprocessor, if a page's reference bit is clear, a counter associated with the page is incremented. If the reference bit is set, the counter is set to zero. After the scan, the pages with the highest counters are removed from the working set. The thread continues examining processes until it has recovered enough pages, then it stops. If a complete pass through all processes still has not recovered enough pages, it makes another pass, trimming more aggressively, even reducing working sets below their minimum if necessary.

On a multiprocessor, looking at the reference bit does not work because although the current CPU may not have touched the page recently, some other one may have. Examining another CPU's reference bits is too expensive to do. Consequently, the reference bit is not examined and the oldest pages are removed.

It should be noted that for page replacement purposes, the operating system itself is regarded as a process. It owns pages and also has a working set. This working set can be trimmed. However, parts of the code and the nonpaged pool are locked in memory and cannot be paged out under any circumstances.

## Physical Memory Management

Above we mentioned that there were actually four free lists. Now it is time to see what all of them are for. Every page in memory is either in one or more working sets or on exactly one of these four lists, which are illustrated in Fig. 11-27. The standby (clean) and modified (dirty) lists hold pages that have recently been evicted from a working set, are still in memory, and are still associated with the process that was using them. The difference between them is that clean pages have a valid copy on disk and can thus be abandoned at will, whereas dirty pages do not have an up-to-date copy on disk. The free list consists of clean pages that are no longer associated with any process. The pages on the zeroed page list are not associated with any process and are also filled with zeros. A fifth list holds any physically defective RAM pages that may exist to make sure that they are not used for anything.



**Figure 11-27.** The various page lists and the transitions between them.

Pages are moved between the working sets and the various lists by the working set manager and other kernel daemon threads. Let us examine the transitions. When the working set manager removes a page from a working set, the page goes on the bottom of the standby or modified list, depending on its state of cleanliness. This transition is shown as (1). Pages on both lists are still valid pages, so if a page fault occurs and one of these pages is needed, it is removed from the list and faulted back into the working set without any disk I/O (2). When a process exits, its nonshared pages cannot be faulted back to it, so they go on the free list (3). These pages are no longer associated with any process.

Other transitions are caused by other daemon threads. Every 4 seconds the **swapper thread** runs and looks for processes all of whose threads have been idle

for a certain number of seconds. If it finds any such processes, their kernel stacks are unpinned and their pages are moved to the standby or modified lists, also shown as (1).

Two other daemon threads, the **mapped page writer** and the **modified page writer**, wake up periodically to see if there are enough clean pages. If there are not, they take pages from the top of the modified list, write them back to disk, and then move them to the standby list (4). The former handles writes to mapped files and the latter handles writes to the paging files. The result of these writes is to transform dirty pages into clean pages.

The reason for having two threads is that a mapped file might have to grow as a result of the write, and growing it requires access to on-disk data structures to allocate a free disk block. If there is no room in memory to bring them in when a page has to be written, a deadlock could result. The other thread can solve the problem by writing out pages to a paging file, which never grows. Nobody ever said Windows 2000 was simple.

The other transitions in Fig. 11-27 are as follows. If a process unmaps a page, the page is no longer associated with a process and can go on the free list (5), except for the case that it is shared. When a page fault requires a page frame to hold the page about to be read in, the page frame is taken from the free list (6), if possible. It does not matter that the page may still contain confidential information because it is about to be overwritten in its entirety. The situation is different when a stack grows.

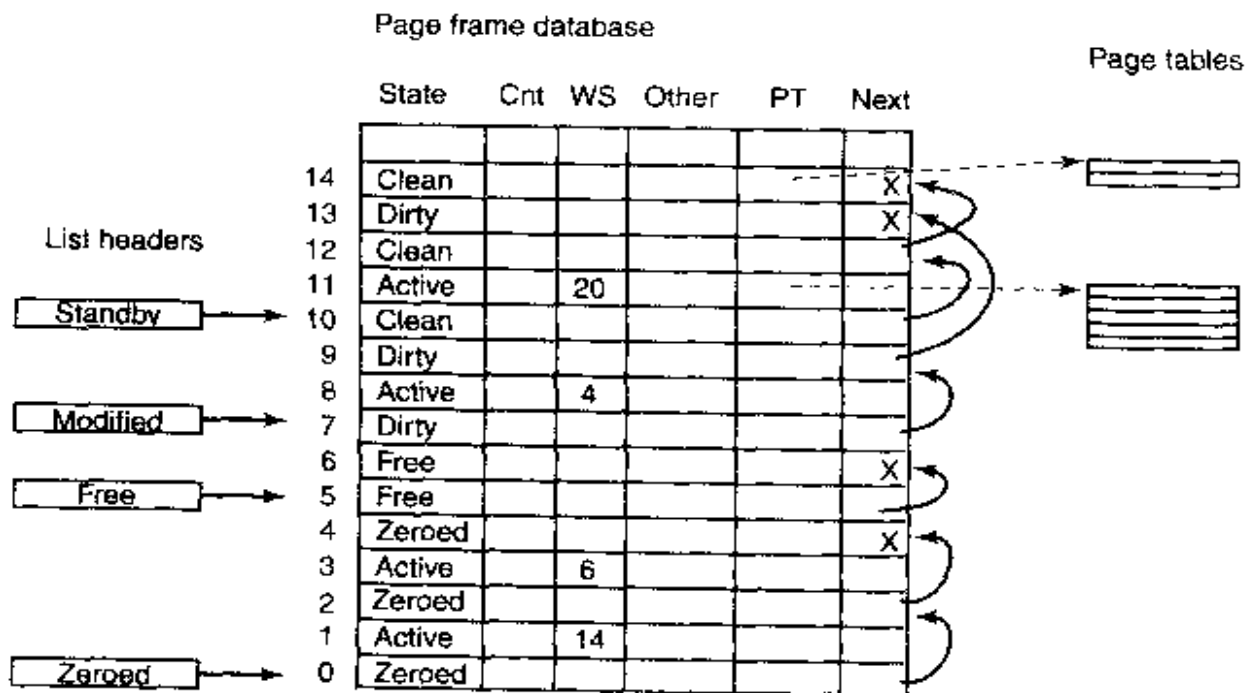
In that case, an empty page frame is needed and the security rules require the page to contain all zeros. For this reason, another kernel daemon thread, the **zero page thread**, runs at the lowest priority (see Fig. 11-19), erasing pages that are on the free list and putting them on the zeroed page list (7). Whenever the CPU is idle and there are free pages, they might as well be zeroed since a zeroed page is potentially more useful than a free page.

The existence of all these lists leads to some subtle policy choices. For example, suppose that a page has to be brought in from disk and the free list is empty. The system is now forced to choose between taking a clean page from the standby list (which might otherwise have been faulted back in later) or an empty page from the zeroed page list (throwing away the work done in zeroing it). Which is better? If the CPU is idle a lot and the zero page thread gets to run often, taking a zeroed page is better because there is no shortage of them. However, if the CPU is always busy and the disk is mostly idle, it is better to take a page from the standby list to avoid the CPU cost of having to zero another page later if a stack grows.

Another puzzle. How aggressively should the daemons move pages from the modified list to the standby list? Having clean pages around is better than having dirty pages around (since they can be reused instantly), but an aggressive cleaning policy means more disk I/O and there is some chance that a newly-cleaned page may be faulted back into a working set and dirtied again.

In general, Windows 2000 resolves these kinds of conflicts through complex heuristics, guesswork, historical precedent, rules of thumb, and administrator-controlled parameter settings. Furthermore, the code is so complex that the designers are loathe to touch parts of it for fear of breaking something somewhere else in the system that nobody really understands any more.

To keep track of all the pages and all the lists, Windows maintains a page frame database with as many entries as there are RAM pages, as shown in Fig. 11-28. This table is indexed by physical page frame number. The entries are fixed length, but different formats are used for different kinds of entries (e.g., valid versus invalid). Valid entries maintain the page's state and a count of how many page tables point to the page, so the system can tell when the page is no longer in use. Pages that are in a working set tell which one. There is also a pointer to the page table pointing to the page, if any (shared pages are handled specially), a link to the next page on the list (if any), and various other fields and flags, such as read in progress, write in progress, etc.



**Figure 11-28.** Some of the major fields in the page frame database for a valid page.

All in all, memory management is a highly complex subsystem with many data structures, algorithms, and heuristics. It attempts to be largely self tuning, but there are also many knobs that administrators can tweak to affect system performance. A number of these knobs and the associated counters can be viewed using tools in the various tool kits mentioned earlier. Probably the most important thing to remember here is that memory management in real systems is a lot more than just one simple paging algorithm like clock or aging.

## 11.6 INPUT/OUTPUT IN WINDOWS 2000

The goal of the Windows 2000 I/O system is to provide a framework for efficiently handling a very wide variety of I/O devices. Current input devices include various kinds of keyboards, mice, touch pads, joysticks, scanners, still cameras, television cameras, bar code readers, microphones, and laboratory rats. Current output devices include monitors, printers, plotters, beamers, CD-recorders, and sound cards. Storage devices include floppy disks, IDE and SCSI hard disks, CD-ROMs, DVDs, Zip drives, and tape drives. Finally, other devices include clocks, networks, telephones, and camcorders. No doubt many new I/O devices will be invented in the years to come, so Window 2000 has been designed with a general framework to which new devices can easily be attached. In the following sections we will examine some of the issues relating to I/O.

### 11.6.1 Fundamental Concepts

The I/O manager is on intimate terms with the plug-and-play manager. The basic idea behind plug and play is that of an enumerable bus. Many buses, including PC Card, PCI, USB, IEEE 1394, and SCSI, have been designed so that the plug-and-play manager can send a request to each slot and ask the device there to identify itself. Having discovered what is out there, the plug-and-play manager allocates hardware resources, such as interrupt levels, locates the appropriate drivers, and loads them into memory. As each one is loaded, a **driver object** is created for it. For some buses, such as SCSI, enumeration happens only at boot time, but for other buses, such as USB and IEEE 1394, it can happen at any moment, requiring close contact between the plug-and-play manager, the bus driver (which actually does the enumeration), and the I/O manager.

The I/O manager is also closely connected with the power manager. The power manager can put the computer into any of six states, roughly described as:

1. Fully operational.
2. Sleep-1: CPU power reduced, RAM and cache on; instant wake-up.
3. Sleep-2: CPU and RAM on; CPU cache off; continue from current PC.
4. Sleep-3: CPU and cache off; RAM on; restart from fixed address.
5. Hibernate: CPU, cache, and RAM off; restart from saved disk file.
6. Off: Everything off; full reboot required.

I/O devices can also be in various power states. Turning them on and off is handled by the power manager and I/O manager together. Note that states 2 through 6 are only used when the CPU has been idle for a shorter or longer time interval.

Somewhat surprisingly, all the file systems are technically I/O drivers. Requests for data blocks from user processes are initially sent to the cache manager. If the cache manager cannot satisfy the request from the cache, it has the I/O manager call the proper file system driver to go get the block it needs from disk.

An interesting feature of Windows 2000 is its support for **dynamic disks**. These disks may span multiple partitions and even multiple disks and may be reconfigured on the fly, without even having to reboot. In this way, logical volumes are no longer constrained to a single partition or even a single disk so that a single file system may span multiple drives in a transparent way.

Another interesting aspect of Windows 2000 is its support for asynchronous I/O. It is possible for a thread to start an I/O operation and then continue executing in parallel with the I/O. This feature is especially important on servers. There are various ways the thread can find out that the I/O has completed. One is to specify an event object at the time the call is made and then wait on it eventually. Another is to specify a queue to which a completion event will be posted by the system when the I/O is done. A third is to provide a callback procedure that the system calls when the I/O has completed.

### 11.6.2 Input/Output API Calls

Windows 2000 has over 100 separate APIs for a wide variety of I/O devices, including mice, sound cards, telephones, tape drives, etc. Probably the most important is the graphics system, for which there are thousands of Win32 API calls. In Sec. 5.7.3 we began our discussion of the Window graphical system. Here we will continue, mentioning a few of the Win32 API categories, each of which has many calls. A brief summary of the categories are given in Fig. 11-29. As we mentioned in Chap. 5, multiple 1500-page books have been written on the graphics portion of the Win32 API.

Win32 calls exist to create, destroy, and manage windows. Windows have a vast number of styles and options that can be specified, including titles, borders, colors, sizes, and scroll bars. Windows can be fixed or movable, of constant size or resizable. Their properties can be queried and messages can be sent to them.

Many windows contain menus, so there are Win32 calls for creating and deleting menus and menu bars. Dynamic menus can be popped up and removed. Menu items can be highlighted, dimmed out, or cascaded.

Dialog boxes are popped up to inform the user of some event or ask a question. They may contain buttons, sliders, or text fields to be filled in. Sounds can also be associated with dialog boxes, for example for warning messages.

There are hundreds of drawing and painting functions available, ranging from setting a single pixel to doing complex region clipping operations. Many calls are provided for drawing lines and closed geometric figures of various kinds, with detailed control over textures, colors, widths, and many other attributes.

API group	Description
Window management	Create, destroy, and manage windows.
Menus	Create, destroy, and append to menus and menu bars
Dialog boxes	Pop up a dialog box and collect information
Painting and drawing	Display points, lines, and geometric figures
Text	Display text in some font, size, and color
Bitmaps and icons	Placement of bitmaps and icons on the screen
Colors and palettes	Manage the set of colors available
The clipboard	Pass information from one application to another
Input	Get information from the mouse and keyboard

Figure 11-29. Some categories of Win32 API calls.

Another group of calls relates to displaying text. Actually, the text display call, `TextOut`, is straightforward. It is the management of the color, point sizes, typefaces, character widths, glyphs, kerning, and other typesetting details where the complexity comes in. Fortunately, the rasterization of text (conversion to bitmaps) is generally done automatically.

Bitmaps are small rectangular blocks of pixels that can be placed on the screen using the `BitBlt` Win32 call. They are used for icons and occasionally text. Various calls are provided for creating, destroying, and managing icon objects.

Many displays use a color mode with only 256 or 65,536 of the  $2^{24}$  possible colors in order to represent each pixel with only 1 or 2 bytes, respectively. In these cases a color palette is needed to determine which 256 or 65,536 colors are available. The calls in this group create, destroy, and manage palettes, select the nearest available color to a given color, and try to make colors on the screen match colors on color printers.

Many Windows 2000 programs allow the user to select some data (e.g., a block of text, part of a drawing, a set of cells in a spreadsheet), put it on the clipboard, and allow it to be pasted into another application. The clipboard is generally used for this transmission. Many clipboard formats are defined, including text, bitmaps, objects, and metafiles. The latter are sets of Win32 calls that when executed draw something, allowing arbitrary drawings to be cut and pasted. This group of calls puts things on the clipboard, takes things off the clipboard, and generally manages it.

Finally, we come to input. There are no Win32 calls for GUI applications for reading input from the keyboard because GUI applications are event driven. The main program consists of a big loop getting input messages. When the user types something interesting, a message is sent to the program telling it what just came in. On the other hand, there are calls relating to the mouse, such as reading its (x, y) position and the state of its buttons. Some of the input calls are actually output



calls, though, such as selecting a mouse cursor icon and moving it around the screen (basically, this is output to the screen). For nonGUI applications, it is possible to read from the keyboard.

### 11.6.3 Implementation of I/O

We could go on more-or-less indefinitely about the Win32 graphics calls, but now it is time to look at how the I/O manager implements graphics and other I/O functions. The main function of the I/O manager is to create a framework in which different I/O devices can operate. The basic structure of the framework is a set of device-independent procedures for certain aspects of I/O plus a set of loaded device drivers for communicating with the devices.

### 11.6.4 Device Drivers

To make sure that device drivers work well with the rest of Windows 2000, Microsoft has defined a **Windows Driver Model** that device drivers are expected to conform with. Furthermore, it also has provided a tool kit that is designed to help driver writers produce conformant drivers. In this section we will briefly examine this model. Conformant drivers must meet all of the following requirements as well as some others:

1. Handle incoming I/O requests, which arrive in a standard format.
2. Be as object based as the rest of Windows 2000.
3. Allow plug-and-play devices to be dynamically added or removed.
4. Permit power management, where applicable.
5. Be configurable in terms of resource usage.
6. Be reentrant for use on multiprocessors.
7. Be portable across Windows 98 and Windows 2000

I/O Requests are passed to drivers in the form of a standardized packet called an **IRP (I/O Request Packet)**. Conformant drivers must be able to handle them. Drivers must be object based in the sense of supporting a specific list of methods that the rest of the system can call. They must also correctly deal with other Windows 2000 objects when given an object handle to deal with.

Conformant drivers must fully support plug and play, which means that if a device managed by the driver is suddenly added or removed from the system, the driver must be prepared to accept this information and act on it, even in the case that the device currently being accessed is suddenly removed. Power management must also be supported for devices for which this is relevant. For example, if the system decides it is now time to go into a low-power hibernation mode, all

devices that are capable of doing this must do so to save energy. They must also wake up when told to do so.

Drivers must be configurable, which means not having any built-in assumptions about which interrupt lines or I/O ports certain devices use. For example, the printer port on the IBM PC and its successors has been at address 0x378 for more than 20 years and it is unlikely to change now. But a printer driver that has this address hard coded into it is not conformant.

Being multiprocessor safe is also a requirement as Windows 2000 was designed for use on multiprocessors. Concretely, this requirement means while a driver is actively running and processing one request on behalf of one CPU, a second request may come in on behalf of a different CPU. The second CPU may begin executing the driver code simultaneously with the first one. The driver must function correctly even when being executed concurrently by two or more CPUs, which implies that all sensitive data structures may only be accessed from inside critical regions. Just assuming that there will not be any other calls until the current one is finished is not permitted.

Finally, conformant drivers must work not only on Windows 2000 but also on Windows 98. It may be necessary to recompile the driver on each system however, and use of C preprocessor commands to isolate platform dependencies is permitted.

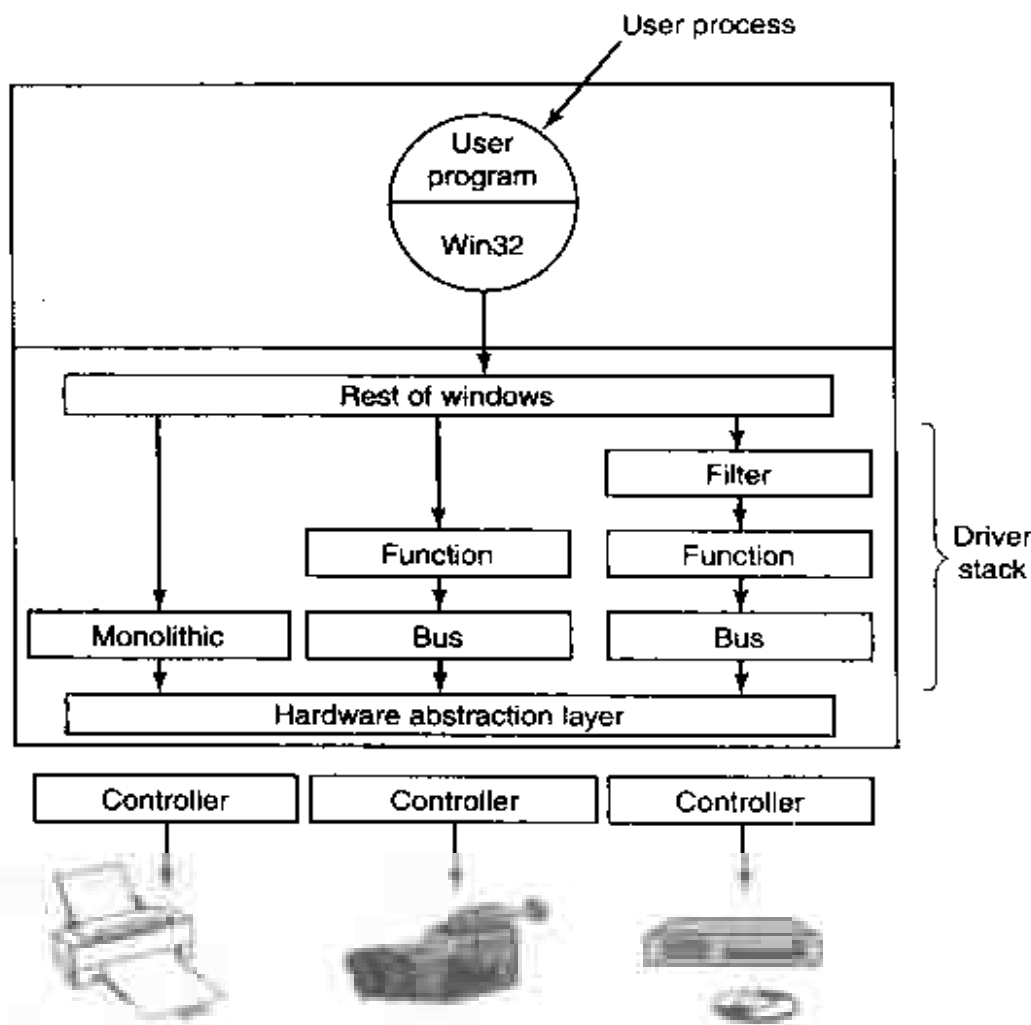
In UNIX, drivers are located by using their major device numbers. Windows 2000 uses a different scheme. At boot time, or when a new hot pluggable plug-and-play device is attached to the computer, Windows 2000 automatically detects it and calls the plug-and-play manager. The manager queries the device to find out what the manufacturer and model number are. Equipped with this knowledge, it looks on the hard disk in a certain directory to see if it has the driver. If it does not, it displays a dialog box asking the user to insert a floppy disk or CD-ROM with the driver. Once the driver is located, it is loaded into memory.

Each driver must supply a set of procedures that can be called to get its services. The first one, called *DriverEntry*, initializes the driver. It is called just after the driver is loaded. It may create tables and data structures, but must not touch the device yet. It also fills in some of the fields of the driver object created by the I/O manager when the driver was loaded. The fields in the driver object include pointers to all the other procedures that drivers must supply. In addition, for each device controlled by the driver (e.g., each IDE disk controlled by the IDE disk driver), a **device object** is created and initialized to point to the driver object. These driver objects are entered into a special directory, \???. Given a device object, the driver object can be located easily, and hence its methods can be called.

A second required procedure is *AddDevice*, which is called once (by the plug-and-play manager) for each device to be added. Once this has been accomplished, the driver is called with the first IRP, which sets up the interrupt vector and actually initializes the hardware. Other procedures that drivers must contain

are the interrupt service procedure, various timer management procedures, a fast I/O path, DMA control, a way to cancel currently executing requests, and many more. All in all, Windows 2000 drivers are so complex that multiple books have been written about them (Cant, 1999; Oney, 1999; and Viscarola and Mason, 1999).

A driver in Windows 2000 may do all the work by itself, as the printer driver does in Fig. 11-30 (just as an example). On the other hand, drivers may also be stacked, which means that a request may pass through a sequence of drivers, each doing part of the work. Two stacked drivers are also illustrated in Fig. 11-30.



**Figure 11-30.** Windows 2000 allows drivers to be stacked.

One common use for stacked drivers is to separate out the bus management from the functional work of actually controlling the device. Bus management on the PCI bus is quite complicated on account of many kinds of modes and bus transactions, and by separating this work from the device-specific part, driver writers are freed from learning how to control the bus. They can just use the standard bus driver in their stack. Similarly, USB and SCSI drivers have a device-specific part and a generic part, with common drivers used for the generic part.

Another use of stacking drivers is to be able to insert **filter drivers** into the stack. A filter driver performs some transformation on the data on the way up or down. For example, a filter driver could compress data on the way to the disk or encrypt data on the way to the network. Putting the filter here means that neither the application program nor the true device driver have to be aware of it and it works automatically for all data going to (or coming from) the device.

## 11.7 THE WINDOWS 2000 FILE SYSTEM

Windows 2000 supports several file systems, the most important of which are **FAT-16**, **FAT-32**, and **NTFS (NT File System)**. FAT-16 is the old MS-DOS file system. It uses 16-bit disk addresses, which limits it to disk partitions no larger than 2 GB. FAT-32 uses 32-bit disk addresses and supports disk partitions up to 2 TB. NTFS is a new file system developed specifically for Windows NT and carried over to Windows 2000. It uses 64-bit disk addresses and can (theoretically) support disk partitions up to  $2^{64}$  bytes, although other considerations limit it to smaller sizes. Windows 2000 also supports read-only file systems for CD-ROMs and DVDs. It is possible (even common) to have the same running system have access to multiple file system types available at the same time.

In this chapter we will treat the NTFS file system because it is a modern file system unencumbered by the need to be fully compatible with the MS-DOS file system, which was based on the CP/M file system designed for 8-inch floppy disks more than 20 years ago. Times have changed and 8-inch floppy disks are not quite state of the art any more. Neither are their file systems. Also, NTFS differs both in user interface and implementation in a number of ways from the UNIX file system, which makes it a good second example to study. NTFS is a large and complex system and space limitations prevent us from covering all of its features, but the material presented below should give a reasonable impression of it.

### 11.7.1 Fundamental Concepts

Individual file names in NTFS are limited to 255 characters; full paths are limited to 32,767 characters. File names are in Unicode, allowing people in countries not using the Latin alphabet (e.g., Greece, Japan, India, Russia, and Israel) to write file names in their native language. For example,  $\phi\iota\lambda\epsilon$  is a perfectly legal file name. NTFS fully supports case sensitive names (so *foo* is different from *Foo* and *FOO*). Unfortunately, the Win32 API does not fully support case-sensitivity for file names and not at all for directory names, so this advantage is lost to programs restricted to using Win32 (e.g., for Windows 98 compatibility).

An NTFS file is not just a linear sequence of bytes, as FAT-32 and UNIX files are. Instead, a file consists of multiple attributes, each of which is represented by a stream of bytes. Most files have a few short streams, such as the name of the file and its 64-bit object ID, plus one long (unnamed) stream with the data. However, a file can also have two or more (long) data streams as well. Each stream has a name consisting of the file name, a colon, and the stream name, as in *foo:stream1*. Each stream has its own size and is lockable independently of all the other streams. The idea of multiple streams in a file was borrowed from the Apple Macintosh, in which files have two streams, the data fork and the resource fork. This concept was incorporated into NTFS to allow an NTFS server be able to serve Macintosh clients.

File streams can be used for purposes other than Macintosh compatibility. For example, a photo editing program could use the unnamed stream for the main image and a named stream for a small thumbnail version. This scheme is simpler than the traditional way of putting them in the same file one after another. Another use of streams is in word processing. These programs often make two versions of a document, a temporary one for use during editing and a final one when the user is done. By making the temporary one a named stream and the final one the unnamed stream, both versions automatically share a file name, security information, timestamps, etc. with no extra work.

The maximum stream length is  $2^{64}$  bytes. To get some idea of how big a  $2^{64}$ -byte stream is, imagine that the stream were written out in binary, with each of the 0s and 1s in each byte occupying 1 mm of space. The  $2^{67}$ -mm listing would be 15 light-years long, reaching far beyond the solar system, to Alpha Centauri and back. File pointers are used to keep track of where a process is in each stream, and these are 64 bits wide to handle the maximum length stream, which is about 18.4 exabytes.

The Win32 API function calls for file and directory manipulation are roughly similar to their UNIX counterparts, except most have more parameters and the security model is different. Opening a file returns a handle, which is then used for reading and writing the file. For graphical applications, no file handles are predefined. Standard input, standard output, and standard error have to be acquired explicitly if needed; in console mode they are preopened, however. Win32 also has a number of additional calls not present in UNIX.

### 11.7.2 File System API Calls in Windows 2000

The principal Win32 API functions for file management are listed in Fig. 11-31. There are actually many more, but these give a reasonable first impression of the basic ones. Let us now examine these calls briefly. **CreateFile** can be used to create a new file and return a handle to it. This API function must also be used to open existing files as there is no **FileOpen** API function. We have not listed the

Win32 API function	UNIX	Description
CreateFile	open	Create a file or open an existing file; return a handle
DeleteFile	unlink	Destroy an existing file
CloseHandle	close	Close a file
ReadFile	read	Read data from a file
WriteFile	write	Write data to a file
SetFilePointer	lseek	Set the file pointer to a specific place in the file
GetFileAttributes	stat	Return the file properties
LockFile	fcntl	Lock a region of the file to provide mutual exclusion
UnlockFile	fcntl	Unlock a previously locked region of the file

**Figure 11-31.** The principal Win32 API functions for file I/O. The second column gives the nearest UNIX equivalent.

parameters for the API functions because they are so voluminous. As an example, CreateFile has seven parameters, which are roughly summarized as follows:

1. A pointer to the name of the file to create or open.
2. Flags telling whether the file can be read, written, or both.
3. Flags telling whether multiple processes can open the file at once.
4. A pointer to the security descriptor, telling who can access the file.
5. Flags telling what to do if the file exists/does not exist.
6. Flags dealing with attributes such as archiving, compression, etc.
7. The handle of a file whose attributes should be cloned for the new file.

The next six API functions in Fig. 11-31 are fairly similar to the corresponding UNIX system calls. The last two allow a region of a file to be locked and unlocked to permit a process to get guaranteed mutual exclusion to it.

Using these API functions, it is possible to write a procedure to copy a file, analogous to the UNIX version of Fig. 6-5. Such a code fragment (without any error checking) is shown in Fig. 11-32. It has been designed to mimic our UNIX version. In practice, one would not have to program a copy file program since CopyFile is an API function (which executes something close to this program as a library procedure).

Windows 2000 NTFS is a hierarchical file system, similar to the UNIX file system. The separator between component names is \ however, instead of /, a fossil inherited from MS-DOS. There is a concept of a current working directory and path names can be relative or absolute. Hard and symbolic links are supported,

```

/* Open files for input and output. */
inhandle = CreateFile("data", GENERIC_READ, 0, NULL, OPEN_EXISTING, 0, NULL);
outhandle = CreateFile("newf", GENERIC_WRITE, 0, NULL, CREATE_ALWAYS,
    FILE_ATTRIBUTE_NORMAL, NULL);

/* Copy the file. */
do {
    s = ReadFile(inhandle, buffer, BUF_SIZE, &count, NULL);
    if (s && count > 0) WriteFile(outhandle, buffer, count, &ocnt, NULL);
} while (s > 0 && count > 0);

/* Close the files. */
CloseHandle(inhandle);
CloseHandle(outhandle);

```

**Figure 11-32.** A program fragment for copying a file using the Windows 2000 API functions.

the former implemented by having multiple directory entries, as in UNIX, and the latter implemented using reparse points (discussed later in this chapter). In addition, compression, encryption, and fault tolerance are also supported. These features and their implementations will be discussed later in this chapter.

The major directory management API functions are given in Fig. 11-33, again along with their nearest UNIX equivalents. The functions should be self explanatory.

Win32 API function	UNIX	Description
CreateDirectory	mkdir	Create a new directory
RemoveDirectory	rmdir	Remove an empty directory
FindFirstFile	opendir	Initialize to start reading the entries in a directory
FindNextFile	readdir	Read the next directory entry
MoveFile	rename	Move a file from one directory to another
SetCurrentDirectory	chdir	Change the current working directory

**Figure 11-33.** The principal Win32 API functions for directory management. The second column gives the nearest UNIX equivalent, when one exists.

### 11.7.3 Implementation of the Windows 2000 File System

NTFS is a highly complex and sophisticated file system. It was designed *from scratch*, rather than being an attempt to improve the old MS-DOS file system. Below we will examine a number of its features, starting with its structure, then moving on to file name lookup, file compression, and file encryption.

## File System Structure

Each NTFS volume (e.g., disk partition) contains files, directories, bitmaps, and other data structures. Each volume is organized as a linear sequence of blocks (clusters in Microsoft's terminology), with the block size being fixed for each volume and ranging from 512 bytes to 64 KB, depending on the volume size. Most NTFS disks use 4-KB blocks as a compromise between large blocks (for efficient transfers) and small blocks (for low internal fragmentation). Blocks are referred to by their offset from the start of the volume using 64-bit numbers.

The main data structure in each volume is the **MFT (Master File Table)**, which is a linear sequence of fixed-size 1-KB records. Each MFT record describes one file or directory. It contains the file's attributes, such as its name and timestamps, and the list of disk addresses where its blocks are located. If a file is extremely large, it is sometimes necessary to use two or more MFT records to contain the list of all the blocks, in which case the first MFT record, called the **base record**, points to the other MFT records. This overflow scheme dates back to CP/M, where each directory entry was called an extent. A bitmap keeps track of which MFT entries are free.

The MFT is itself a file and as such can be placed anywhere within the volume, thus eliminating the problem with defective sectors in the first track. Furthermore, the file can grow as needed, up to a maximum size of  $2^{48}$  records.

The MFT is shown in Fig. 11-34. Each MFT record consists of a sequence of (attribute header, value) pairs. Each attribute begins with a header telling which attribute this is and how long the value is because some attribute values are variable length, such as the file name and the data. If the attribute value is short enough to fit in the MFT record, it is placed there. If it is too long, it is placed elsewhere on the disk and a pointer to it is placed in the MFT record.

The first 16 MFT records are reserved for NTFS metadata files, as shown in Fig. 11-34. Each of the records describes a normal file that has attributes and data blocks, just like any other file. Each of these files has a name that begins with a dollar sign to indicate that it is a metadata file. The first record describes the MFT file itself. In particular, it tells where the blocks of the MFT file are located so the system can find the MFT file. Clearly, Windows 2000 needs a way to find the first block of the MFT file in order to find the rest of the file system information. The way it finds the first block of the MFT file is to look in the boot block, where its address is installed at system installation time.

Record 1 is a duplicate of the early part of the MFT file. This information is so precious that having a second copy can be critical in the event one of the first blocks of the MFT ever goes bad. Record 2 is the log file. When structural changes are made to the file system, such as adding a new directory or removing an existing one, the action is logged here before it is performed, in order to increase the chance of correct recovery in the event of a failure during the operation. Changes to file attributes are also logged here. In fact, the only changes not



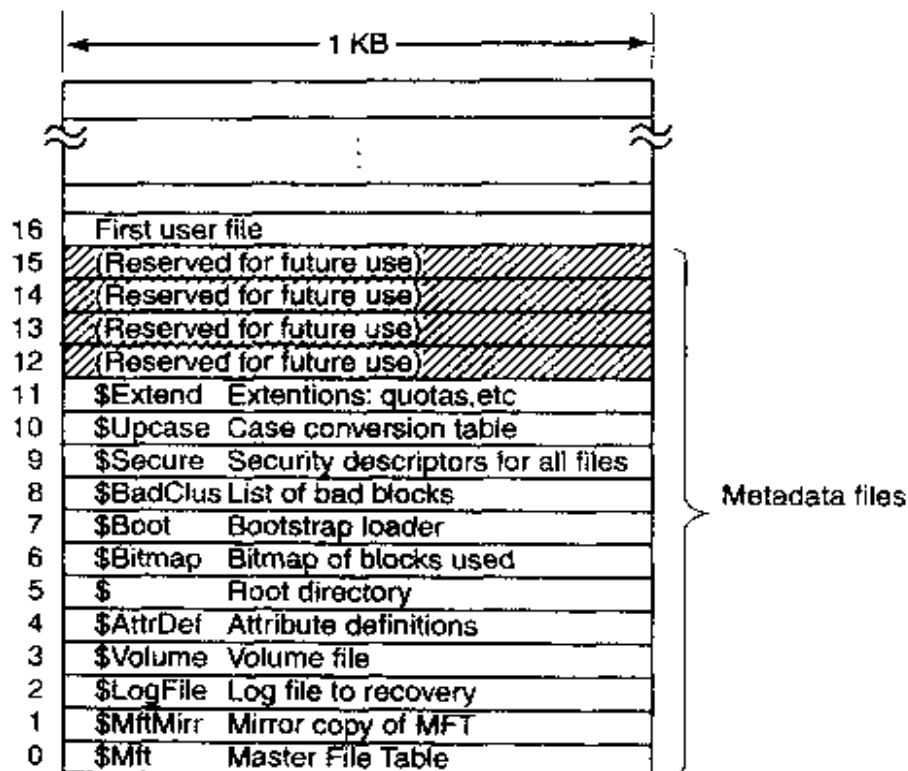


Figure 11-34. The NTFS master file table.

logged here are changes to user data. Record 3 contains information about the volume, such as its size, label, and version.

As mentioned above, each MFT record contains a sequence of (attribute header, value) pairs. The *\$AttrDef* file is where the attributes are defined. Information about this file is in MFT record 4. Next comes the root directory, which itself is a file and can grow to arbitrary length. It is described by MFT record 5.

Free space on the volume is kept track of with a bitmap. The bitmap is itself a file and its attributes and disk addresses are given in MFT record 6. The next MFT record points to the bootstrap loader file. Record 8 is used to link all the bad blocks together to make sure they never occur in a file. Record 9 contains the security information. Record 10 is used for case mapping. For the Latin letters A-Z case mapping is obvious (at least for people who speak Latin). Case mapping for other languages, such as Greek, Armenian, or Georgian (the country, not the state), is less obvious to Latin speakers, so this file tells how to do it. Finally, record 11 is a directory containing miscellaneous files for things like disk quotas, object identifiers, reparse points, and so on. The last 4 MFT records are reserved for future use.

Each MFT record consists of a record header followed by a sequence of (attribute header, value) pairs. The record header contains a magic number used for validity checking, a sequence number updated each time the record is reused for a new file, a count of references to the file, the actual number of bytes in the record

used, the identifier (index, sequence number) of the base record (used only for extension records), and some other miscellaneous fields. Following the record header comes the header of the first attribute, then the first attribute value, the second attribute header, the second attribute value, and so on.

NTFS defines 13 attributes that can appear in MFT records. These are listed in Fig. 11-35. Each MFT record consists of a sequence of attribute headers, each of which identifies the attribute it is heading and gives the length and location of the value field along with a variety of flags and other information. Usually, attribute values follow their attribute headers directly, but if a value is too long to fit in the MFT record, it may be put in a separate disk block. Such an attribute is said to be a **nonresident attribute**. The data attribute is an obvious candidate. Some attributes, such as the name, may be repeated, but all attributes must appear in a fixed order in the MFT record. The headers for resident attributes are 24 bytes long; those for nonresident attributes are longer because they contain information about where to find the attribute on disk.

Attribute	Description
Standard information	Flag bits, timestamps, etc.
File name	File name in Unicode; may be repeated for MS-DOS name
Security descriptor	Obsolete. Security information is now in \$Extend\$Secure
Attribute list	Location of additional MFT records, if needed
Object ID	64-bit file identifier unique to this volume
Reparse point	Used for mounting and symbolic links
Volume name	Name of this volume (used only in \$Volume)
Volume information	Volume version (used only in \$Volume)
Index root	Used for directories
Index allocation	Used for very large directories
Bitmap	Used for very large directories
Logged utility stream	Controls logging to \$LogFile
Data	Stream data; may be repeated

**Figure 11-35.** The attributes used in MFT records.

The standard information field contains the file owner, security information, the timestamps needed by POSIX, the hard link count, the read-only and archive bits, etc. It is a fixed-length field and is always present. The file name is variable length in Unicode. In order to make files with nonMS-DOS names accessible to old 16-bit programs, files can also have an 8 + 3 MS-DOS name. If the actual file name conforms to the MS-DOS 8 + 3 naming rule, a secondary MS-DOS name is not used.

In NT 4.0, security information could be put in an attribute, but in Windows 2000 it all goes into a single file so that multiple files can share the same security

descriptions. The attribute list is needed in case the attributes do not fit in the MFT record. This attribute then tells where to find the extension records. Each entry in the list contains a 48-bit index into the MFT telling where the extension record is and a 16-bit sequence number to allow verification that the extension record and base records match up.

The object ID attribute gives the file a unique name. This is sometimes needed internally. The reparse point tells the procedure parsing the file name to do something special. This mechanism is used for mounting and symbolic links. The two volume attributes are only used for volume identification. The next three attributes deal with how directories are implemented. Small ones are just lists of files but large ones are implemented using B+ trees. The logged utility stream attribute is used by the encrypting file system.

Finally, we come to the attribute that everyone has been waiting for: the data. The stream name, if present, goes in this attribute header. Following the header is either a list of disk addresses telling which blocks the file contains, or for files of only a few hundred bytes (and there are many of these), the file itself. Putting the actual file data in the MFT record is called an **immediate file** (Mullender and Tanenbaum, 1987).

Of course, most of the time the data does not fit in the MFT record, so this attribute is usually nonresident. Let us now take a look at how NTFS keeps track of the location of nonresident attributes, in particular data.

The model for keeping track of disk blocks is that they are assigned in runs of consecutive blocks, where possible, for efficiency reasons. For example, if the first logical block of a file is placed in block 20 on the disk, then the system will try hard to place the second logical block in block 21, the third logical block in 22, and so on. One way to achieve these runs is to allocate disk storage several blocks at a time, if possible.

The blocks in a file are described by a sequence of records, each one describing a sequence of logically contiguous blocks. For a file with no holes in it, there will be only one such record. Files that are written in order from beginning to end all belong in this category. For a file with one hole in it (e.g., only blocks 0–49 and blocks 60–79 are defined), there will be two records. Such a file could be produced by writing the first 50 blocks, then seeking forward to logical block 60 and writing another 20 blocks. When a hole is read back, all the missing bytes are zeros.

Each record begins with a header giving the offset of the first block within the file. Next comes the offset of the first block not covered by the record. In the example above, the first record would have a header of (0, 50) and would provide the disk addresses for these 50 blocks. The second one would have a header of (60, 80) and would provide the disk addresses for these 20 blocks.

Each record header is followed by one or more pairs, each giving a disk address and run length. The disk address is the offset of the disk block from the start of its partition; the run length is the number of blocks in the run. As many

pairs as needed can be in the run record. Use of this scheme for a three-run, nine-block file is illustrated in Fig. 11-36.

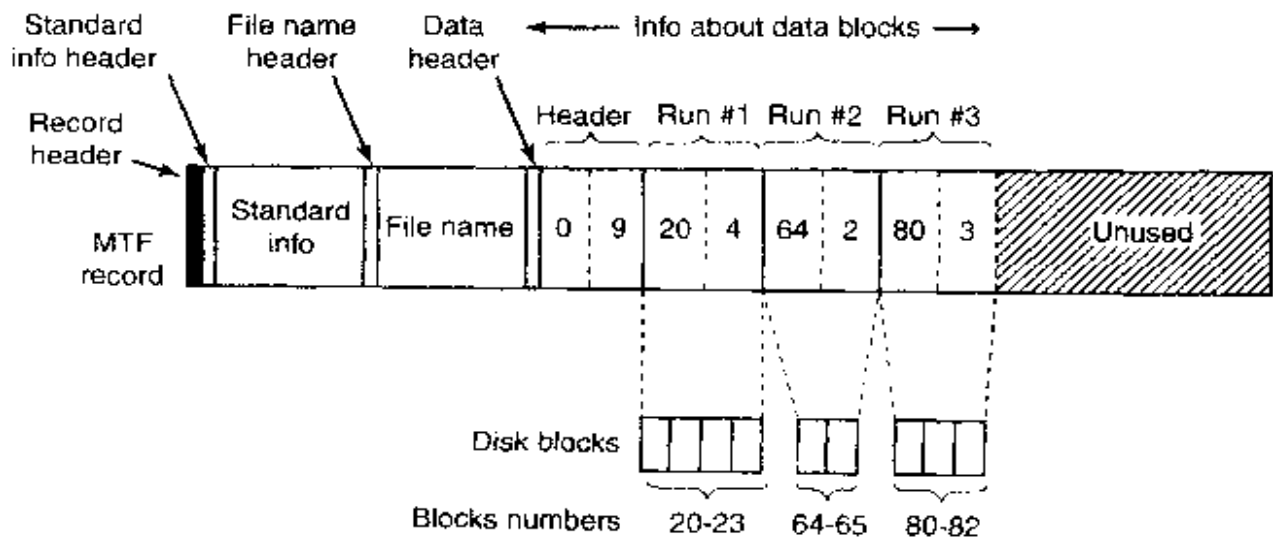


Figure 11-36. An MFT record for a three-run, nine-block file.

In this figure we have an MFT record for a short file (short here means that all the information about the file blocks fits in one MFT record). It consists of the three runs of consecutive blocks on the disk. The first run is blocks 20-23, the second is blocks 64-65, and the third is blocks 80-82. Each of these runs is recorded in the MFT record as a (disk address, block count) pair. How many runs there are depends on how good a job the disk block allocator did in finding runs of consecutive blocks when the file was created. For a  $n$ -block file, the number of runs can be anything from 1 up to and including  $n$ .

Several comments are worth making here. First, there is no upper limit to the size of files that can be represented this way. In the absence of address compression, each pair requires two 64-bit numbers in the pair for a total of 16 bytes. However, a pair could represent 1 million or more consecutive disk blocks. In fact, a 20-MB file consisting of 20 separate runs of 1 million 1-KB blocks each fits easily in one MFT record, whereas a 60-KB file scattered into 60 isolated blocks does not.

Second, while the straightforward way of representing each pair takes  $2 \times 8$  bytes, a compression method is available to reduce the size of the pairs below 16. Many disk addresses have multiple high-order zero-bytes. These can be omitted. The data header tells how many are omitted, that is, how many bytes are actually used per address. Other kinds of compression are also used. In practice, the pairs are often only 4 bytes.

Our first example was easy: all the file information fit in one MFT record. What happens if the file is so large or highly fragmented that the block information does not fit in one MFT record? The answer is simple: use two or more MFT records. In Fig. 11-37 we see a file whose base record is in MFT record 102. It

has too many runs for one MFT record, so it computes how many extension records it needs, say, two, and puts their indices in the base record. The rest of the record is used for the first  $k$  data runs.

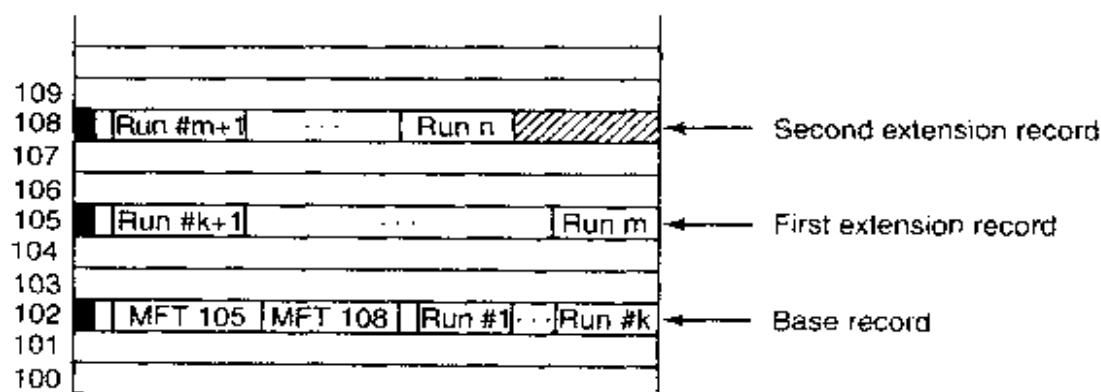


Figure 11-37. A file that requires three MFT records to store all its runs.

Note that Fig. 11-37 contains some redundancy. In theory, it should not be necessary to specify the end of a sequence of runs because this information can be calculated from the run pairs. The reason for “overspecifying” this information is to make seeking more efficient: to find the block at a given file offset, it is only necessary to examine the record headers, not the run pairs.

When all the space in record 102 has been used up, storage of the runs continues with MFT record 105. As many runs are packed in this record as fit. When this record is also full, the rest of the runs go in MFT record 108. In this way many MFT records can be used to handle large fragmented files.

A problem arises if so many MFT records are needed that there is no room in the base MFT to list all their indices. There is also a solution to this problem: the list of extension MFT records is made nonresident (i.e., stored on disk instead of in the base MFT record). Then it can grow as large as needed.

An MFT entry for a small directory is shown in Fig. 11-38. The record contains a number of directory entries, each of which describes one file or directory. Each entry has a fixed-length structure followed by a variable-length file name. The fixed part contains the index of the MFT entry for the file, the length of the file name, and a variety of other fields and flags. Looking for an entry in a directory consists of examining all the file names in turn.

Large directories use a different format. Instead of listing the files linearly, a B+ tree is used to make alphabetical lookup possible and to make it easy to insert new names in the directory in the proper place.

## File Name Lookup

We now have enough information to see how file name lookup occurs. When a user program wants to open a file, it typically makes a call like

```
CreateFile("C:\maria\web.htm", ...)
```

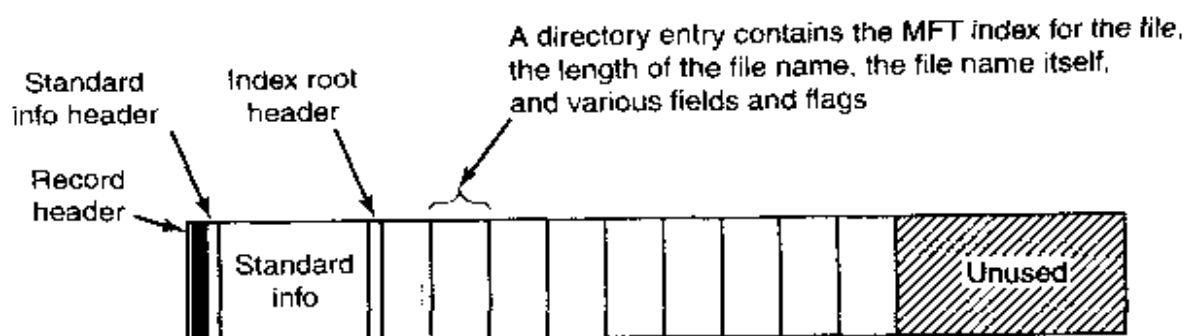


Figure 11-38. The MFT record for a small directory.

This call goes to the user-level shared library, *kernel32.dll*, where `\\?` is prepended to the file name giving

`\\?\\C:\\maria\\web.htm`

It is this name that is passed as a parameter to the system call `NtFileCreate`.

Then the operating system starts the search at the root of the object manager's name space (see Fig. 11-12). It then looks in the directory `\\?` to find `C:`, which it will find. This file is a symbolic link to another part of the object manager's name space, the directory `\\Device`. The link typically ends at an object whose name is something like `\\Device\\HarddiskVolume1`. This object corresponds to the first partition of the first hard disk. From this object it is possible to determine which MFT to use, namely the one on this partition. These steps are shown in Fig. 11-39.

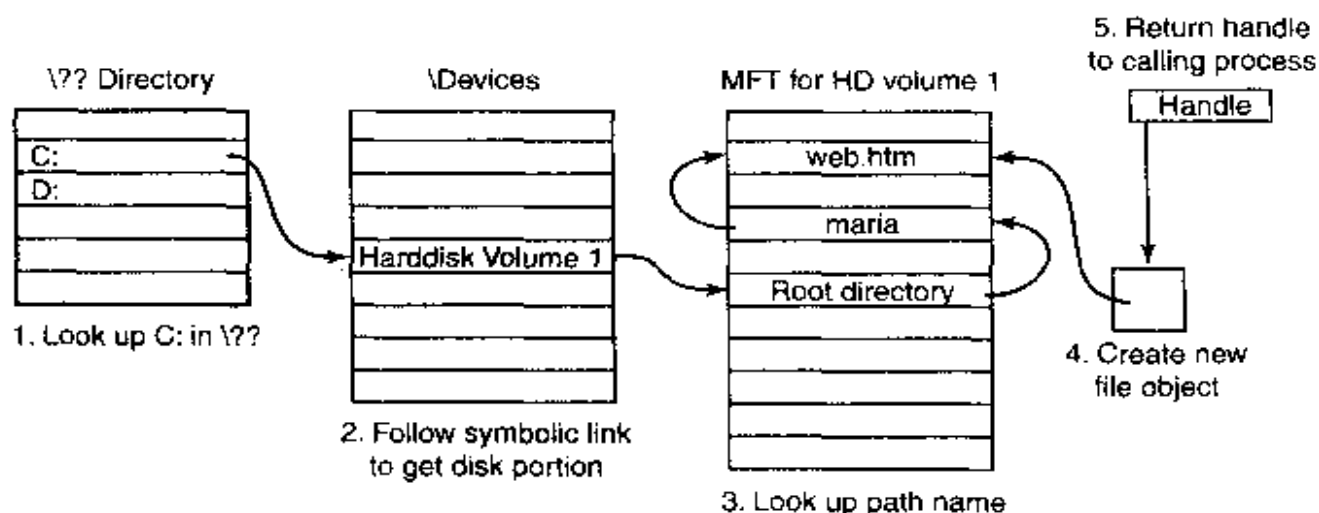


Figure 11-39. Steps in looking up the file `C:\\maria\\web.htm`.

The parsing of the file name continues now at the root directory, whose blocks can be found from entry 5 in the MFT (see Fig. 11-34). The string "maria" is now looked up in the root directory, which returns the index into the MFT for the directory *maria*. This directory is then searched for the string "web.htm". If

successful, the result is a new object created by the object manager. The object, which is unnamed, contains the index of the MFT record for the file. A handle to this object is returned to the calling process. On subsequent `ReadFile` calls, the handle is provided, which allows the object manager to find the index and then the contents of the MFT record for the file. If a thread in a second process opens the file again, it gets a handle to a new file object.

In addition to regular files and directories, NTFS supports hard links in the UNIX sense, and also symbolic links using a mechanism called **reparse points**. It is possible to tag a file or directory as a reparse point and associate a block of data with it. When the file or directory is encountered during a file name parse, exception processing is triggered and the block of data is interpreted. It can do various things, including redirecting the search to a different part of the directory hierarchy or even to a different partition. This mechanism is used to support both symbolic links and mounted file systems.

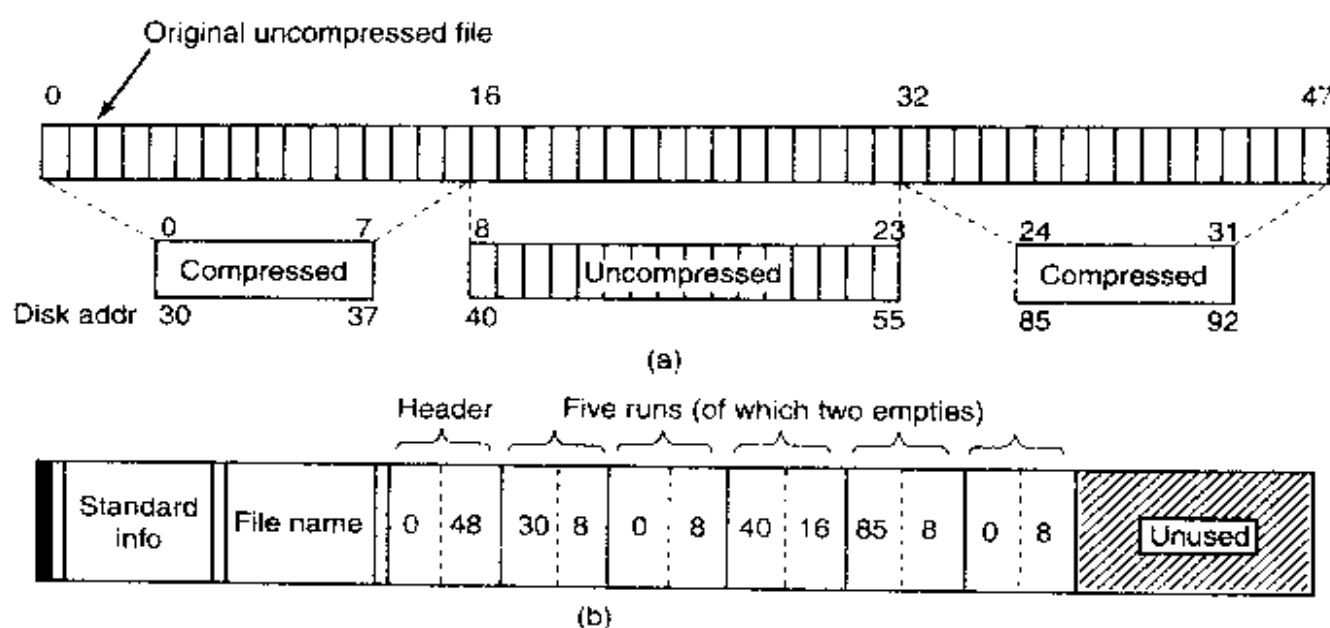
## File Compression

NTFS supports transparent file compression. A file can be created in compressed mode, which means that NTFS automatically tries to compress the blocks as they are written to disk and automatically uncompresses them when they are read back. Processes that read or write compressed files are completely unaware of the fact that compression and decompression are going on.

Compression works as follows. When NTFS writes a file marked for compression to disk, it examines the first 16 (logical) blocks in the file, irrespective of how many runs they occupy. It then runs a compression algorithm on them. If the resulting data can be stored in 15 or fewer blocks, the compressed data are written to the disk, preferably in one run, if possible. If the compressed data still take 16 blocks, the 16 blocks are written in uncompressed form. Then blocks 16-31 are examined to see if they can be compressed to 15 blocks or less, and so on.

Figure 11-40(a) shows a file in which the first 16 blocks have successfully compressed to eight blocks, the second 16 blocks failed to compress, and the third 16 blocks have also compressed by 50%. The three parts have been written as three runs and stored in the MFT record. The “missing” blocks are stored in the MFT entry with disk address 0 as shown in Fig. 11-40(b). Here the header (0, 48) is followed by five pairs, two for the first (compressed) run, one for the uncompressed run, and two for the final (compressed) run.

When the file is read back, NTFS has to know which runs are compressed and which are not. It sees that based on the disk addresses. A disk address of 0 indicates that it is the final part of 16 compressed blocks. Disk block 0 may not be used for storing data, to avoid ambiguity. Since it contains the boot sector, using it for data is impossible anyway.



**Figure 11-40.** (a) An example of a 48-block file being compressed to 32 blocks.  
 (b) The MFT record for the file after compression.

Random access to compressed files is possible, but tricky. Suppose that a process does a seek to block 35 in Fig. 11-40. How does NTFS locate block 35 in a compressed file? The answer is that it has to read and decompress the entire run first. Then it knows where block 35 is and can pass it to any process that reads it. The choice of 16 blocks for the compression unit was a compromise. Making it shorter would have made the compression less effective. Making it longer would have made random access more expensive.

## File Encryption

Computers are used nowadays to store all kinds of sensitive data, including plans for corporate takeovers, tax information, and love letters (love email?), which the owners do not especially want revealed to anyone. Information loss can happen when a laptop computer is lost or stolen, a desktop system is rebooted using an MS-DOS floppy disk to bypass Windows 2000 security, or a hard disk is physically removed from one computer and installed on another one with an insecure operating system. Even the simple act of going to the bathroom and leaving the computer unattended and logged in can be a huge security breach.

Windows 2000 addresses these problem by having an option to encrypt files, so even in the event the computer is stolen or rebooted using MS-DOS, the files will be unreadable. The normal way to use Windows 2000 encryption is to mark certain directories as encrypted, which causes all the files in them to be encrypted, and new files moved to them or created in them to be encrypted as well. The actual encryption and decryption is not done by NTFS itself, but by a driver called EFS (**E**ncrypting **F**ile **S**ystem), which is positioned between NTFS and the user



process. In this way, application programs are unaware of encryption and NTFS itself is only partially involved in it.

To understand how the encrypting file system works, it is necessary to understand how modern cryptography works. For this purpose, a brief review was given in Sec. 9.2. Readers not familiar with the basics of cryptography should read that section before continuing.

Now let us see how Windows 2000 encrypts files. When the user asks a file to be encrypted, a random 128-bit file key is generated and used to encrypt the file block by block using a symmetric algorithm parametrized by this key. Each new file encrypted gets a different 128-bit random file key, so no two files use the same encryption key, which increases security in case one key is compromised. The current encryption algorithm is a variant of **DES (Data Encryption Standard)**, but the EFS architecture supports the addition of new algorithms in the future. Encrypting each block independently of all the others is necessary to make random access still possible.

The file key has to be stored somewhere so the file can be decrypted later. If it were just stored on the disk in plaintext, then someone who stole or found the computer could easily decrypt the file, defeating the purpose of encrypting the files. For this reason, the file keys must all be encrypted before they are stored on the disk. Public-key cryptography is used for this purpose.

After the file is encrypted, the location of the user's public key is looked up using information in the registry. There is no danger of storing the public key's location in the registry because if a thief steals the computer and finds the public key, there is no way to deduce the private key from it. The 128-bit random file key is now encrypted with the public key and the result stored on disk along with the file, as shown in Fig. 11-41.

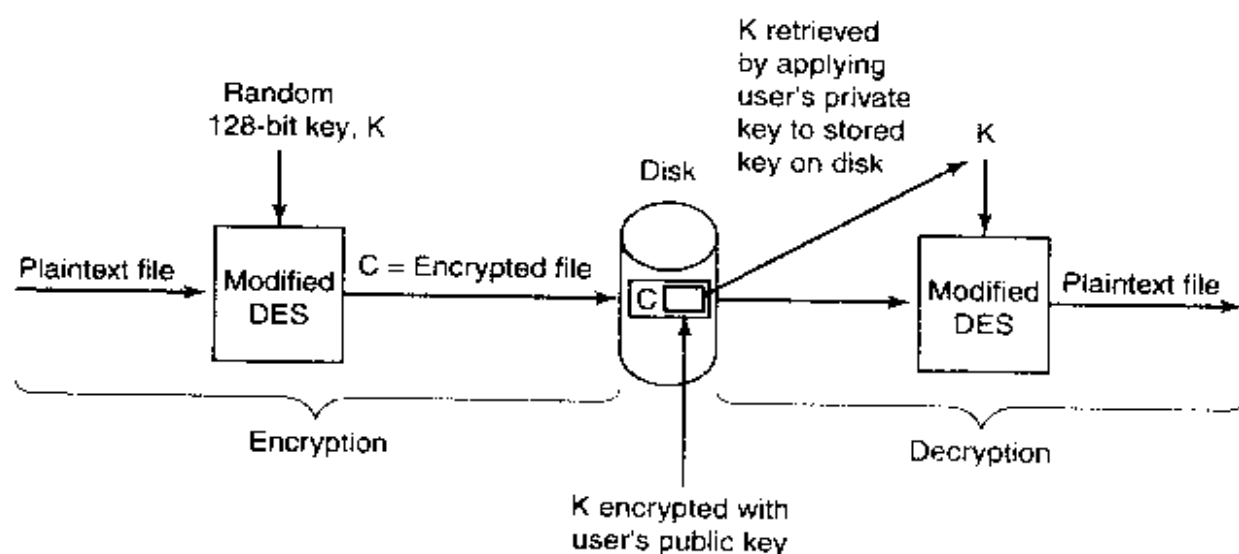


Figure 11-41. Operating of the encrypting file system.

To decrypt a file, the encrypted 128-bit random file key is fetched from disk. However, to decrypt it and retrieve the file key, the user must present the private

key. Ideally, this should be stored on a smart card, external to the computer, and only inserted in a reader when a file has to be decrypted. Although Windows 2000 supports smart cards, it does not store private keys on them.

Instead, the first time a user encrypts a file using EFS, Windows 2000 generates a (private key, public key) pair and stores the private key on disk encrypted using a symmetric encryption algorithm. The key used for the symmetric algorithm is derived either from the user's login password or from a key stored on the smart card, if smart card login is enabled. In this way, EFS can decrypt the private key at login time and keep it within its own virtual address space during normal operation so it can decrypt the 128-bit file keys as needed without further disk accesses. When the computer is shut down, the private key is erased from EFS' virtual address space so anyone stealing the computer will not have access to the private key.

A complication occurs when multiple users need access to the same encrypted file. Currently the shared use of encrypted files by multiple users is not supported. However, the EFS architecture could support sharing in the future by encrypting each file's key multiple times, once with the public key of each authorized user. All of these encrypted versions of the file key could be attached to the file.

The potential need to share encrypted files is one reason why this two-key system is used. If all files were encrypted by their owner's key, there would be no way to share any files. By using a different key to encrypt each file, this problem can be solved.

Having a random file key per file but encrypting it with the owner's symmetric key does not work because having the symmetric encryption key just lying around in plain view would ruin the security—generating the decryption key from the encryption key is too easy. Thus (slow) public-key cryptography is needed to encrypt the file keys. Because the encryption key is public anyway, having it lying around is not dangerous.

The other reason the two-key system is used is performance. Using public-key cryptography to encrypt each file would be too slow. It is much more efficient to use symmetric-key cryptography to encrypt the data and public-key cryptography to encrypt the symmetric file key.

## 11.8 SECURITY IN WINDOWS 2000

Having just looked at encryption in the file system, this is a good time to examine security in general. NT was designed to meet the U.S. Department of Defense's C2 security requirements (DoD 5200.28-STD), the Orange Book, which we studied in Chap. 9. This standard requires operating systems to have certain properties in order to be classified as secure enough for certain kinds of

military work. Although Windows 2000 was not specifically designed for C2 compliance, it inherits many security properties from NT, including the following:

1. Secure login with antispoofting measures.
2. Discretionary access controls.
3. Privileged access controls.
4. Address space protection per process.
5. New pages must be zeroed before being mapped in.
6. Security auditing.

Let us review these items briefly (none of which are met by Windows 98, incidentally).

Secure login means that the system administrator can require all users to have a password in order to log in. Spoofing is when a malicious user writes a program that displays the login prompt or screen and then walks away from the computer in the hope that an innocent user will sit down and enter a name and password. The name and password are then written to disk and the user is told that login has failed. Windows 2000 prevents this attack by instructing users to hit CTRL-ALT-DEL to log in. This key sequence is always captured by the keyboard driver, which then invokes a system program that puts up the genuine login screen. This procedure works because there is no way for user processes to disable CTRL-ALT-DEL processing in the keyboard driver.

Discretionary access controls allow the owner of a file or other object to say who can use it and in what way. Privileged access controls allow the system administrator (superuser) to override them when needed. Address space protection simply means that each process has its own protected virtual address space not accessible by any unauthorized process. The next item means that when a stack grows, the pages mapped in are initialized to zero so processes cannot find any old information put there by the previous owner (hence the zeroed page list in Fig. 11-27, which provides a supply of zeroed pages for this purpose). Finally, security auditing allows the administrator to produce a log of certain security-related events.

In the next section we will describe the basic concepts behind Windows 2000 security. After that we will look at the security system calls. Finally, we will conclude by seeing how security is implemented.

### 11.8.1 Fundamental Concepts

Every Windows 2000 user (and group) is identified by a **SID (Security ID)**. SIDs are binary numbers with a short header followed by a long random component. Each SID is intended to be unique worldwide. When a user starts up a

process, the process and its threads run under the user's SID. Most of the security system is designed to make sure that each object can be accessed only by threads with authorized SIDs.

Each process has an **access token** that specifies its SID and other properties. It is normally assigned at login time by *winlogon* and is shown in Fig. 11-42, although processes should call *GetTokenInformation* to acquire this information since it may change in the future. The header contains some administrative information. The expiration time field could tell when the token ceases to be valid, but it is currently not used. The *Groups* fields specify the groups to which the process belongs; this is needed for POSIX conformance. The default **DACL (Discretionary ACL)** is the access control list assigned to objects created by the process if no other ACL is specified. The user SID tells who owns the process. The restricted SIDs are to allow untrustworthy processes to take part in jobs with trustworthy processes but with less power to do damage.

Finally, the privileges listed, if any, give the process special powers, such as the right to shut the machine down or access files to which access would otherwise be denied. In effect, the privileges split up the power of the superuser into several rights that can be assigned to processes individually. In this way, a user can be given some superuser power, but not all of it. In summary, the access token tells who owns the process and which defaults and powers are associated with it.

Header	Expiration time	Groups	Default DACL	User SID	Group SID	Restricted SIDs	Privileges
--------	-----------------	--------	--------------	----------	-----------	-----------------	------------

Figure 11-42. Structure of an access token.

When a user logs in, *winlogon* gives the initial process an access token. Subsequent processes normally inherit this token on down the line. A process' access token initially applies to all the threads in the process. However, a thread can acquire a different access token during execution, in which case the thread's access token overrides the process' access token. In particular, a client thread can pass its access token to a server thread to allow the server to access the client's protected files and other objects. This mechanism is called **impersonation**.

Another basic concept is the **security descriptor**. Every object has a security descriptor associated with it that tells who can perform which operations on it. A security descriptor consists of a header followed by a DACL with one or more **ACEs (Access Control Elements)**. The two main kinds of elements are Allow and Deny. An allow element specifies a SID and a bitmap that specifies which operations processes with that SID may perform on the object. A deny element works the same way, except a match means the caller may not perform the operation. For example, Ida has a file whose security descriptor specifies that everyone

has read access, Elvis has no access. Cathy has read/write access, and Ida herself has full access. This simple example is illustrated in Fig. 11-43. The SID Everyone refers to the set of all users, but it is overridden by any explicit ACEs that follow.

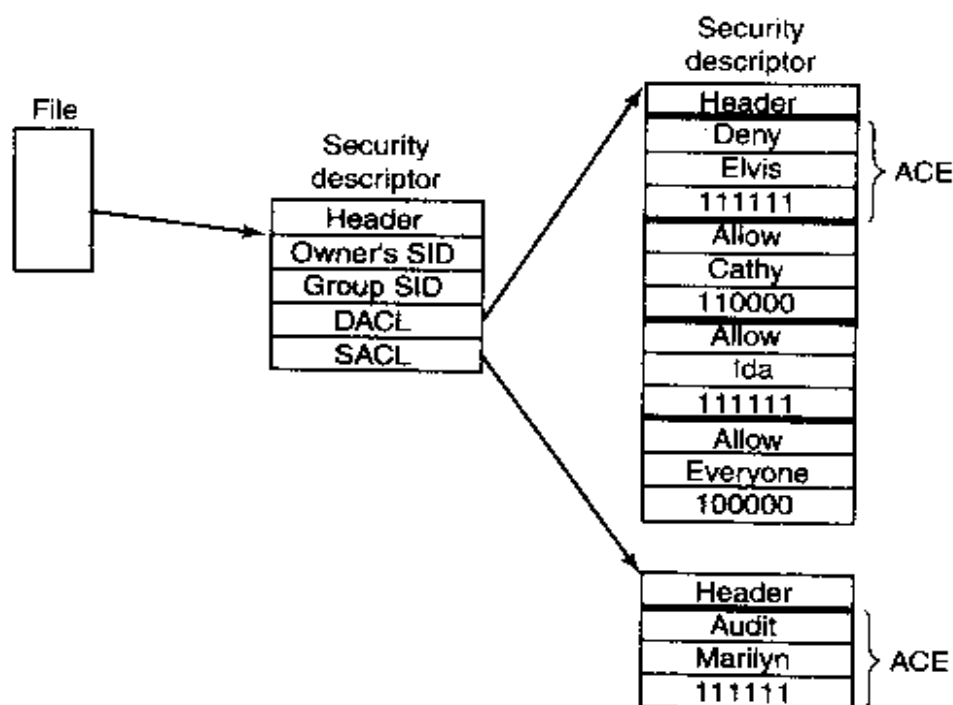


Figure 11-43. An example security descriptor for a file.

In addition to the DACL, a security descriptor also has a **SACL** (**S**ystem **A**ccess **C**ontrol **L**ist), which is like a DACL except that it specifies not who may use the object, but which operations on the object are recorded in the system-wide security event log. In Fig. 11-43, every operation that Marilyn performs on the file will be logged. Windows 2000 provides additional auditing features to log sensitive accesses.

## 11.8.2 Security API Calls

Most of the Windows 2000 access control mechanism is based on security descriptors. The usual pattern is that when a process creates an object, it provides a security descriptor as one of the parameters to the `CreateProcess`, `CreateFile`, or other object creation call. This security descriptor then becomes the security descriptor attached to the object, as we saw in Fig. 11-43. If no security descriptor is provided in the object creation call, the default security in the caller's access token (see Fig. 11-42) is used instead.

Many of the Win32 API security calls relate to the management of security descriptors, so we will focus on those here. The most important calls are listed in Fig. 11-44. To create a security descriptor, storage for it is first allocated and then initialized using `InitializeSecurityDescriptor`. This call fills in the header. If the

owner SID is not known, it can be looked up by name using `LookupAccountSid`. It can then be inserted into the security descriptor. The same holds for the group SID, if any. Normally, these will be the caller's own SID and one of the caller's groups, but the system administrator can fill in any SIDs.

Win32 API function	Description
<code>InitializeSecurityDescriptor</code>	Prepare a new security descriptor for use
<code>LookupAccountSid</code>	Look up the SID for a given user name
<code>SetSecurityDescriptorOwner</code>	Enter the owner SID in the security descriptor
<code>SetSecurityDescriptorGroup</code>	Enter a group SID in the security descriptor
<code>InitializeAcl</code>	Initialize a DACL or SACL
<code>AddAccessAllowedAce</code>	Add a new ACE to a DACL or SACL allowing access
<code>AddAccessDeniedAce</code>	Add a new ACE to a DACL or SACL denying access
<code>DeleteAce</code>	Remove an ACE from a DACL or SACL
<code>SetSecurityDescriptorDacl</code>	Attach a DACL to a security descriptor

**Figure 11-44.** The principal Win32 API functions for security.

At this point the security descriptor's DACL (or SACL) can be initialized with `InitializeAcl`. ACL entries can be added using `AddAccessAllowedAce`, and `AddAccessDeniedAce`. These calls can be repeated multiple times to add as many ACE entries as are needed. `DeleteAce` can be used to remove an entry, more like on an existing ACL than on one being constructed for the first time. When the ACL is ready, `SetSecurityDescriptorDacl` can be used to attach it to the security descriptor. Finally, when the object is created, the newly minted security descriptor can be passed as a parameter to have it attached to the object.

### 11.8.3 Implementation of Security

Security in a standalone Windows 2000 system is implemented by a number of components, most of which we have already seen (networking is a whole other story and beyond the scope of this book). Logging in is handled by *winlogon* and authentication is handled by *lsass* and *msgina.dll* as discussed in Sec. 11.4.5. The result of a successful login is a new shell with its associated access token. This process uses the SECURITY and SAM keys in the registry. The former sets the general security policy and the latter contains the security information for the individual users, as discussed in Sec. 11.2.3.

Once a user is logged in, security operations happen when an object is opened for access. Every `OpenXXX` call requires the name of the object being opened and the set of rights needed. During processing of the open, the security manager (see Fig. 11-7) checks to see if the caller has all the rights required. It performs this check by looking at the caller's access token and the DACL associated with the

object. It goes down the list of entries in the ACL in order. As soon as it finds an entry that matches the caller's SID or one of the caller's groups, the access found there is taken as definitive. If all the rights the caller needs are available, the open succeeds; otherwise it fails.

DACLs can have Deny entries as well as Allow entries, as we have seen. For this reason, it is usual to put entries denying access ahead of entries granting access in the ACL, so that a user who is specifically denied access cannot get in via a back door by being a member of a group that has legitimate access.

After an object has been opened, a handle to it is returned to the caller. On subsequent calls, the only check that is made is whether the operation now being tried was in the set of operations requested at open time, to prevent a caller from opening a file for reading and then trying to write on it. Any log entries required by the SACL are made.

## 11.9 CACHING IN WINDOWS 2000

The Windows 2000 cache manager does caching for performance reasons, conceptually similar to caches in other operating systems. However, its design has some unusual properties that are worth looking at briefly.

The cache manager's job is to keep file system blocks that have been used recently in memory to reduce access time on any subsequent reference. Windows 2000 has a single integrated cache that works for all the file systems in use, including NTFS, FAT-32, FAT-16, and even CD-ROM file systems. This means that the file systems do not need to maintain their own caches.

As a consequence of the design goal to have a single integrated cache despite the presence of multiple file systems, the cache manager is located in an unusual position in the system, as we saw in Fig. 11-7. It is not part of the file system because there are multiple independent file systems that may have nothing in common. Instead it operates at a higher level than the file systems, which are technically drivers under control of the I/O manager.

The Windows 2000 cache is organized by virtual block, not physical block. To see what this means, remember that traditional file caches keep track of blocks by two-part addresses of the form (partition, block), where the first member denotes a device and partition and the second member is a block number within that partition. The Windows 2000 cache manager does not do that. Instead, it uses (file, offset) to refer to a block. (Technically, streams are cached rather than files, but we will ignore that detail below.)

The reason for this unorthodox arrangement is that when a request comes in to the cache manager, it is specified as (file, offset) because that is all the calling process knows. If cache blocks had been labeled by a (partition, block) tag, the cache manager would have no way of knowing which (file, offset) block cor-

responds to which (partition, block) block since it is the file systems that maintain those mappings.

Let us now examine how the cache manager works. When a file is referenced, the cache manager maps a 256-KB chunk of kernel virtual address space onto the file. If the file is larger than 256 KB, only a portion of the file is mapped. The total amount of virtual address space the cache manager can use is determined at boot time and depends on the amount of RAM present. If the cache manager runs out of 256-KB chunks of virtual address space, it must unmap an old file before mapping in a new one.

Once a file is mapped, the cache manager can satisfy requests for its blocks by just copying from kernel virtual address space to the user buffer. If the block copied is not in physical memory, a page fault will occur and the memory manager will satisfy the fault in the usual way. The cache manager is not even aware of whether the block was in the cache or not. The copy always succeeds.

The operation of the cache manager is shown in Fig. 11-45 in the case of an NTFS file system on a SCSI disk and a FAT-32 file system on an IDE disk. When a process does a read on a file, the request is directed to the cache manager. If the block needed is in the cache, it is copied to the user immediately. If it is not in the cache, the cache manager gets a page fault when trying to copy it. When the page fault has been handled, the block is copied to the calling process.

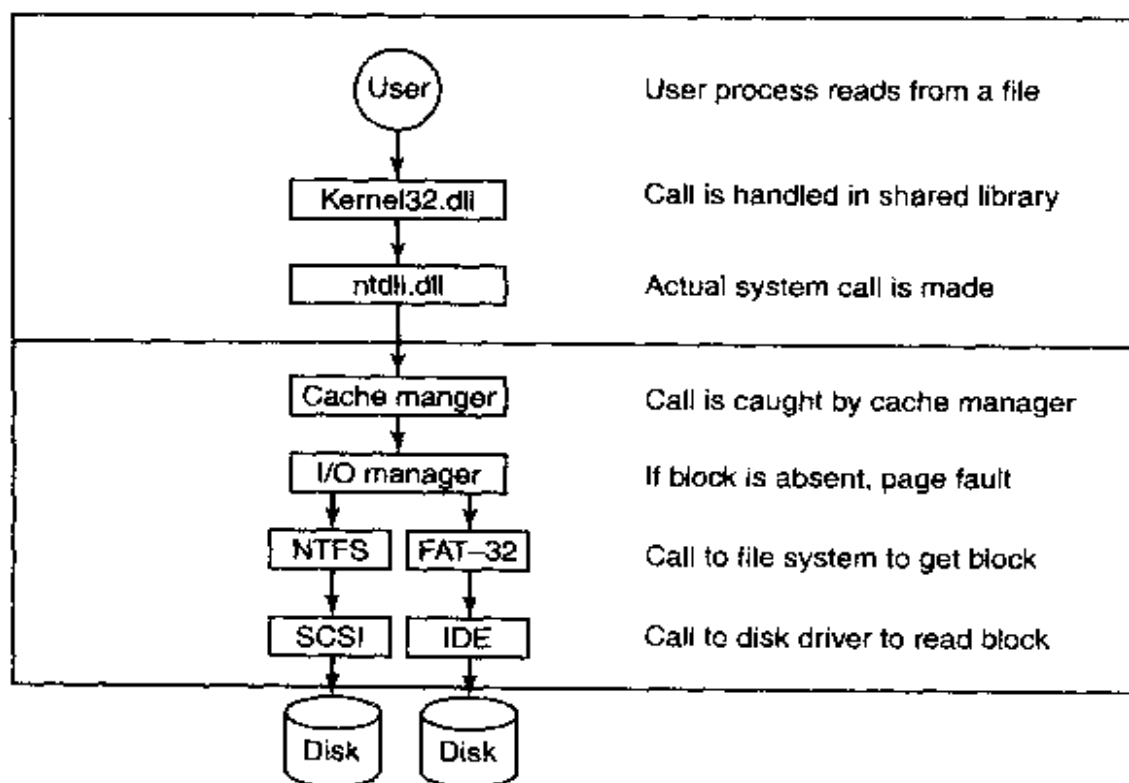


Figure 11-45. The path through the cache to the hardware.

As a consequence of this design, the cache manager does not know how many of its mapped pages are in physical memory or even how large its cache is. Only



the memory manager knows for sure. This approach allows the memory manager to dynamically trade off the size of the cache against memory for user pages. If there is little file activity but there are many processes active, the memory manager can use most of physical memory for process pages. On the other hand, if there is a lot of file activity and few processes, more physical memory can be devoted to the cache.

Another property the cache manager has is that coherence is maintained between memory-mapped files and files that are open for reading and writing. Consider, for example, a situation in which one process opens some file for reading and writing and a second process maps that file onto its address space. What happens if the second process writes on the file directly and immediately thereafter the first one reads in a block that was just changed? Does it get stale data?

The answer is no. In both cases—open files and mapped files—the cache manager maps a 256-KB piece of its virtual address space onto the file. The file is only mapped once, no matter how many processes have it open or mapped. In the case of a mapped file, both the cache manager and the user process share pages in memory. When a read request has to be satisfied, the cache manager just copies the page from memory to the user buffer, which always reflects the exact current state of the file because the cache manager is using the same pages as the process that has the file mapped in.

## 11.10 SUMMARY

Windows 2000 is structured in the HAL, the kernel, the executive, and a thin system services layer that catches incoming system calls. In addition, there are a variety of device drivers, including the file systems and GDI. The HAL hides certain differences in hardware from the upper layers. The kernel attempts to hide the remaining differences from the executive to make it almost completely machine independent.

The executive is based on memory objects. User processes create them and get back handles to manipulate them later. Executive components can also create objects. The object manager maintains a name space into which objects can be inserted for subsequent lookup.

Windows 2000 supports processes, jobs, threads, and fibers. Processes have virtual address spaces and are containers for resources. Threads are the unit of execution and are scheduled by the operating system. Fibers are lightweight threads that are scheduled entirely in user space. Jobs are collections of processes and are used for assigning resource quotas. Scheduling is done using a priority algorithm in which the highest priority ready thread runs next.

Windows 2000 supports a demand-paged virtual memory. The paging algorithm is based on the working set concept. The system maintains several lists of

free pages, so that when a page fault occurs, there is generally a free page available. The free page lists are fed by trimming the working sets using complex formulas that try to throw out pages that have not been used in a long time.

I/O is done by device drivers, which follow the Windows Device Model. Each driver starts out by initializing a driver object that contains the addresses of the procedures that the system can call to add devices or perform I/O. Drivers can be stacked to act as filters.

The NTFS file system is based on a master file table, which has one record per file or directory. Each file has multiple attributes, which can either be in the MFT record or nonresident, on the disk. NTFS supports compression and encryption, among other features.

Security is based on access control lists. Each process has an access control token that tells who it is and what special privileges it has, if any. Each object has a security descriptor associated with it. The security descriptor points to a discretionary access control list that contains access control entries that can allow or deny access to individuals or groups.

Finally, Windows 2000 maintains a single system-wide cache for all the file systems. It is a virtual cache rather than a physical one. Requests for disk blocks go first to the cache. If they cannot be satisfied, then the appropriate file system is called to fetch the needed blocks.

## PROBLEMS

1. Give one advantage and one disadvantage of the registry versus having individual *.ini* files?
2. A mouse can have 1, 2, or 3 buttons. All three types are in use. Does the HAL hide this difference from the rest of the operating system? Why or why not?
3. The HAL keeps track of time starting in the year 1601. Give an example of an application where this feature is useful.
4. The POSIX subsystem needs to implement UNIX-style signals. If a user hits the key for the Quit signal, is this scheduled as a DPC or an APC?
5. Many components of the executive (Fig. 11-7) call other components of the executive. Give three examples of one component calling another one, but use (six) different components in all.
6. Win32 does not have signals. If they were to be introduced, they could be per process, per thread, both, or neither. Make a proposal and explain why it is a good idea.
7. An alternative to using DLLs is to statically link each program with precisely those library procedures it actually calls, no more and no less. If this scheme were to be introduced, would it make more sense on client machines or on server machines?

8. The file *ntdll.dll* exports 1179 function calls while *ntoskrnl.exe* exports 1209 function calls. Is this a bug? What might be the cause of this discrepancy?
9. Objects managed by the object manager are variable-sized, with different objects having different sizes. Can an object start at an arbitrary byte in the nonpaged pool?  
*Hint:* You do not need any information about Windows 2000 other than what was given in the text.
10. Is there any limit on the number of different operations that can be defined on an executive object? If so, where does this limit come from? If not, why not?
11. The Win32 API call *WaitForMultipleObjects* allows a thread to block on a set of synchronization objects whose handles are passed as parameters. As soon as any one of them is signaled, the calling thread is released. Is it possible to have the set of synchronization objects include two semaphores, one mutex, and one critical section? Why or why not? *Hint:* This is not a trick question but it does require some careful thought.
12. Name three reasons why a process might be terminated.
13. Consider the situation of Fig. 11-19 in which the system is about to schedule a thread. Assuming that each thread is compute bound, how long does it take before some thread at priority 3 gets to run on Windows 2000 Professional?
14. Suppose that the quantum is set to 20 msec and the current thread, at priority 24, has just started a quantum. Suddenly an I/O operation completes and a priority 28 thread is made ready. About how long does it have to wait to get serviced?
15. In Windows 2000, the current priority is always greater than or equal to the base priority. Are there any circumstances in which it would make sense to have the current priority be lower than the base priority? If so, give an example. If not, why not?
16. Some MS-DOS programs were written in assembly language using instructions such as *IN port* and *OUT port*. Can such programs be run under Windows 2000? If not, can you think of a way they could be supported?
17. Name two ways to give better response time to important processes.
18. Position-independent code was briefly discussed in the text. It is a technique for allowing two processes to share the same procedure at different virtual addresses. Can this problem be solved by just setting the page tables of the two processes appropriately? Explain your answer.
19. Shared libraries in Windows 2000 are contained in *.dll* files that multiple processes may map in at the same time. If two processes need to map in the same shared library at different virtual addresses, a problem will occur. How could this problem be solved on the Pentium using a property of its memory architecture? If the solution requires changes to Windows 2000 to implement, state what changes are needed.
20. If a region of virtual address space is reserved but not committed, do you think a VAD is created for it? Defend your answer.
21. Which of the transitions shown in Fig. 11-27 are policy decisions, as opposed to required moves forced by system events (e.g., a process exiting and freeing its pages)?

22. Suppose that a page is shared and in two working sets at once. If it is evicted from one of the working sets, where does it go in Fig. 11-27?
23. When a process unmaps a clean page, it makes the transition (5) in Fig. 11-27. Why is there no transition to the modified list when a dirty page is unmapped?
24. Assuming that both RAM and video RAM each can read or write a 32-bit word in 10 nsec, how long does it take to paint the background on an XGA screen in the best case?
25. A file has the following mapping. Give the MFT run entries.

Offset	0	1	2	3	4	5	6	7	8	9	10
Disk address	50	51	52	22	24	25	26	53	54	-	60

26. Consider the MFT record of Fig. 11-36. Suppose that the file grew and a 10th block was assigned to the end of the file. The number of this block is 66. What would the MFT record look like now?
27. In Fig. 11-40(b), the first two runs are each of length 8 blocks. Is it just an accident they are equal or does this have to do with the way compression works? Explain your answer.
28. Suppose that you wanted to build Windows 2000 Lite. Which of the fields of Fig. 11-42 could be removed without weakening the security of the system?
29. The *regedit* command can be used to export part or all of the registry to a text file under all current versions of Windows. Save the registry several times during a work session and see what changes. If you have access to a Windows computer on which you can install software or hardware, find out what changes when a program or device are added or removed.
30. Write a UNIX program that simulates writing an NTFS file with multiple streams. It should accept a list of one or more files as arguments and write an output file that contains one stream with the attributes of all arguments and additional streams with the contents of each of the arguments. Now write a second program for reporting on the attributes and streams and extracting all the components.
31. Write a program that generates an MS-DOS 8 + 3 filename given any file name. Use the Windows 2000 algorithm.