# *Index*

## Symbols

\ (backslash)
   disallowed in filenames, 23
   for root directory, 16, 57

## A

absolute pathnames, file objects, 177
access
   Direct Memory Access (DMA), 256
   file object, 177
   mount points and, 28
   mounted logical volumes, 24
   violations, memory and, 69
   virtual address space, 196-201
AcquireFileForNtCreateSection( ), 547-549
AcquireForCcFlush( ), 551
AcquireForLazyWrite( ), 289, 354
   example of, 553-554
AcquireForModWrite( ), 549-551
AcquireForReadAhead( ), 289, 352
administration
   event logging, 86-93
   status reports to, 66
aliasing, 185
allocating, 39–41
   access violations and, 69
   to Cache Manager, 248-249
   device objects, 140-141
   dispatcher object storage, 98
   ERSOURCE structures storage, 110
   for event log entries, 91-92
   event object storage, 101
   Executive spin lock storage, 96
   file objects, 176
   file streams, 263, 346
      changing file stream size, 487
   FSD design and, 364
   IRPs, 144-146, 150, 636-639
   lookaside lists and, 45
   MEM_ types for, 208
   multiple stack locations, 155
   termination handlers and, 84
   virtual address spaces, 206-209
   VPB structure, 369
   zones, 42
ANSI strings, converting, 47
APCs (asynchronous procedure calls), 10, 107
   APC_LEVEL, 10
   postprocessing IRPs, 150, 151, 154
APIs (application programming interfaces), 5
   as event log viewer, 87, 90
   MPR module, 61
arbitrary threads, 131-133
associated IRPs, 147, 647-650
asynchronous
   file object operations, 177
   I/O, 124, 458, 464-476
   IRP handling, 149-150
   paging I/O requests, 166

## S

# About the Author

Rajeev Nagar has been working on operating systems (specifically storage management systems) for the past six years. He has designed and implemented kernel software for the Windows NT, AIX, HPUX, and SunOS platforms. His file system development work has included local, disk-based file systems, networked file systems, and distributed file systems. His undergraduate degree is in computer engineering, and he has a master's degree in computer science. Rajeev has implemented an OSF distributed file system client on the Windows NT platform, as well as other filter drivers for storage management products.

# Colophon

A vulture is featured on the cover of *Windows NT File System Internals.* Vultures are divided into two famlies—New World vultures, a family that includes the majestic but near-extinct California condor, and Old World vultures. Both families are closely related to eagles and hawks, but, unlike their relatives, vultures are carrion eaters, not hunters. A vulture will rarely kill for food. Instead, they sit by and wait for another animal to die before starting to dine. Vultures often live in open country where herds of large mammals, such as cattle, can be found. They fly in slow circles, searching the ground for dead, sick, or injured animals. They also watch for running packs of jackals or hyenas, who often lead them to food. When food has been spotted, the vulture swoops down to the ground, and other circling vultures follow.

Both Old World and New World vultures have heads and necks that are almost bare, covered only by a thin layer of down. Many vultures have a thick ruff of feathers around their neck. These adaptations allow the vulture to place its head deep inside carcasses without soiling its plumage. The digestive enzymes of the vulture allow it to survive on decaying meat that would be toxic to other animals.

Although the modern view of vultures is often one of disgust and comtempt, some ancient cultures revered them as embodiments of immortality.

Edie Freedman designed the cover of this book, using a 19th-century engraving from the Dover Pictorial Archive. The cover layout was produced with Quark XPress 3-3 using the ITC Garamond font. Whenever possible, our books use Rep-Kover™, a durable and flexible lay-flat binding. If the page count exceeds Rep-Kover's limit, perfect binding is used.

The inside layout was designed by Nancy Priest and implemented in FrameMaker 5.0 by Mike Sierra. The text and heading fonts are ITC Garamond Light and Garamond Book. The illustrations that appear in the book were created in Macromedia Freehand 7.0 by Robert Romano. This colophon was written by Clairemarie Fisher O'Leary.