

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Sthultz
Identity Theft and Financial Fraud Research and Operation Center

The Philosophy Of Digital Data Hiding

Digital data hiding is actually a cluster-concept that spans many contexts. In modern times, non-physical data hiding is usually associated with digital forms such as cryptography, steganography, and watermarking. Although related in the sense that they all are means to achieve secure or proprietary communications, there are differences among their three activities at a number of levels – some of which are quite subtle. To illustrate, the cryptographer's interest is primarily with obscuring the content of a message, but not the communication of the message. The steganographer, on the other hand is concerned with hiding the very communication of the message, while the digital watermarker attempts to add sufficient metadata to a message to establish ownership, provenance, source, etc.

Cryptography and steganography share the feature that the object of interest is embedded, hidden or obscured, whereas the object of interest in watermarking is the host or carrier which is being protected by the object that is embedded, hidden or obscured. Further, watermarking and steganography may be used with or without cryptography; and imperceptible watermarking shares functionality with steganography, whereas perceptible watermarking does not. Overviews exist for cryptography [13], steganography [9], [14], and watermarking [3].

However, there is also a physical aspect of digital data hiding. In this case, digital storage locations are used to hide or conceal data. Obviously, these storage locations must be somewhat obscure or detection would be trivial. Examples include covert channeling (e.g. within TCP/IP packet headers or Loki's use of the ICMP options field) or obscure or infrequently used addressable space in memory media. One of the earliest examples of this arose when microcomputers were first introduced. There was a difference between the number of tracks that a floppy drive controller could access (usually 81 or 82) and the number of tracks that were recognized by the operating system (DOS recognized only 80). The upper 2 tracks were used to hide data by vendors and hackers alike. For awhile they were even used for product licensing.

A modern analogue might involve taking advantage of the physical characteristics of a storage medium to hide data. The steganographic tool, Camouflage, was such a program (camouflage.unfiction.com). Camouflage embeds messages in the file slack (see below). This simple approach to data hiding has the advantage that the characteristics of the host or carrier message is unaffected, and are transparent to the host operating system and file managers. The disadvantage is that the hidden message is easily recovered with a hex-editor. Although this approach has not met with great success (and Camouflage is no longer supported) it is a great segue into the art of data hiding by taking advantage of the physical characteristics of computer systems.

The ability to hide data in computers is a byproduct of the system and peripheral architectures. If all storage were bit-addressable at the operating system level, there would be no place to hide data. For efficiency considerations, addressability has to be at more abstract levels (typically words in primary, and blocks in secondary). Such abstractions create digital warrens where data may go unnoticed or, in some cases, be inaccessible.

Digital Storage And File Systems

Digital storage encompasses a wide range of media including diskettes, hard drives, zip disks, USB flash drives, compact flash cards, CD ROMs, DVDs, etc. Since the structures of most of these media can be related to hard drives, a discussion of hard drive architecture will serve to illustrate the various data hiding mechanisms.

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Stultz
 Identity Theft and Financial Fraud Research and Operation Center

Hard drive technology existed long before the advent of personal computers (although the early personal computers started out with only cassette tape or diskette storage!). As PCs evolved the emphasis on backward compatibility and the need for user convenience (i.e. making data structures totally transparent to the user) greatly influenced the development of secondary and tertiary storage technology – both physically and logically. An unintended consequence has been the creation of many places where data can be intentionally hidden or unintentionally left behind.

A functioning hard drive actually consists of a geometric structure and a set of nested data structures: hard drive, partition, file system, file, record, and field. Hidden data can be found at each of these levels. We will assume that the reader is familiar with basic hard disk geometry (cylinders, tracks, blocks, clusters, and sectors).

The following describes the various hiding mechanisms (as illustrated in Figure 1), starting at the level of the hard drive itself and then working down through the set of nested data structures.

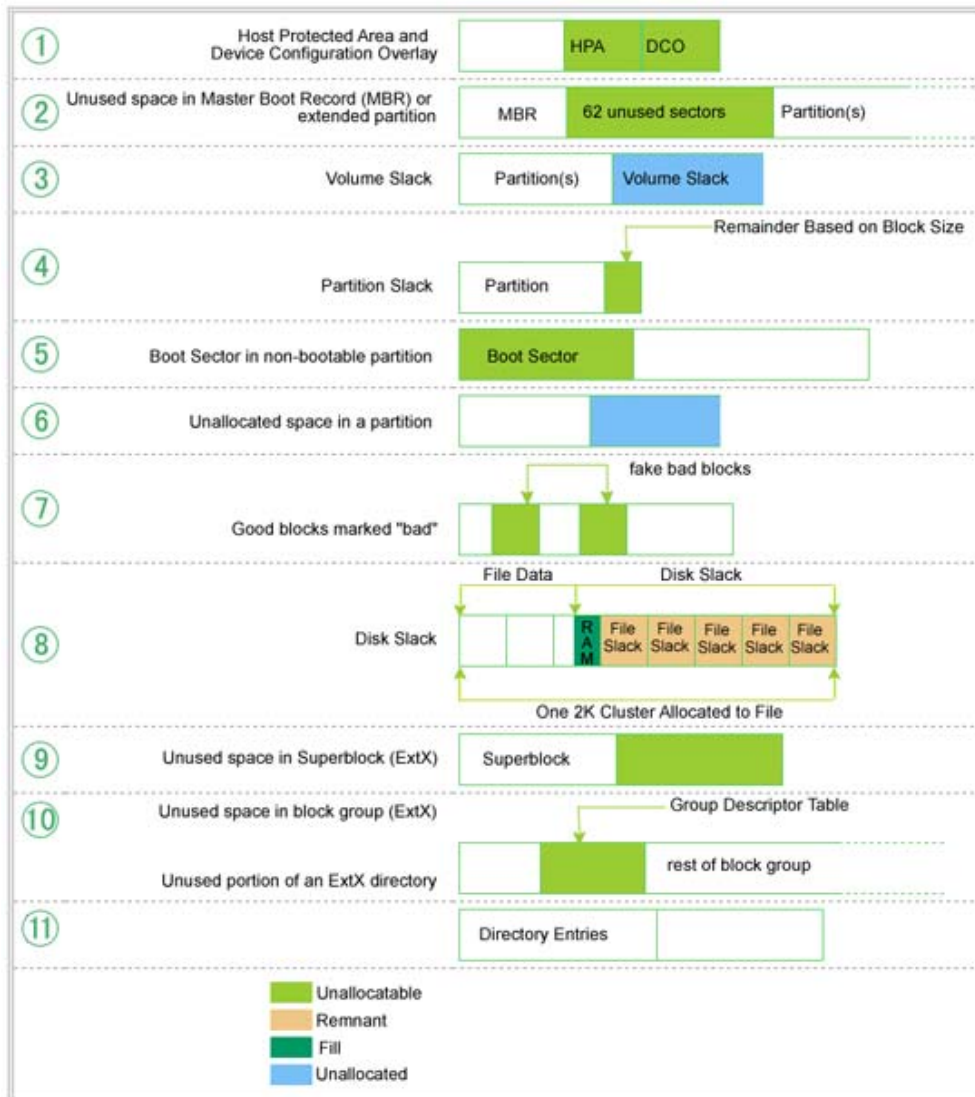


Figure 1: Digital Disk Warrens

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Stultz

Identity Theft and Financial Fraud Research and Operation Center

Some hard drives can have a reserved area of the disk called the Host Protected Area (HPA) (see Figure 1, item 1). Device Configuration Overlay allows modification of the apparent features provided by a hard drive, e.g., the number of available clusters. This was designed to be an area where computer vendors could store data that is protected from normal user activities. It is not affected by operating system utilities (format, delete, etc.) and cannot be accessed without the use of a special program that reconfigures the controller to access all physical blocks. It is not difficult, however, to write a program to access these areas, write data to them, and subsequently return the area to a HPA. This is an example of a hiding method that takes advantage of what is more or less a “physical” feature of the drive architecture.

At the next layer, common operating systems typically require that a hard drive be partitioned into virtual file systems before it can be used. This is true even if the entire hard drive is to be mapped onto a single partition. A partition is a set of consecutive blocks on a hard disk that appear to the operating system as a separate logical volume (drive for Windows vs. directory or mount point for Unix). Note that there are several different types of partition formats. This discussion is confined to the partition format that has evolved from the original PC (often referred to as DOS partitions) and will not apply to Apple, xBSD, Sun Solaris, GPT partitioning, or multiple disk volumes. A recommended source for additional detail appears as a reference [6]. Even so, the various data hiding techniques at the data and file system levels are partition format independent.

Every hard drive using a DOS partition has space reserved at the beginning of the drive for a Master Boot Record (MBR) (see Figure 1, item 2). This will often contain the boot code necessary to begin the initial program load of an operating system and will always contain a partition table (provided we are dealing with partitioned media; e.g., floppy disks are not typically partitioned media while fixed disks are) defining the size and location of up to four partitions. Since the MBR only requires a single sector and partitions must start on a cylinder boundary, this results in 62 sectors of empty MBR space where data can be hidden.

As disk sizes grew beyond the limitations of existing operating systems, there arose a need for more than four partitions on a single hard disk. Extended partitions (as opposed to primary partitions) were then designed that could contain multiple logical partitions. Each of these extended partitions contains a structure similar to the MBR, leaving another 62 sectors within each extended partition with the potential of harboring more hidden data. Extended partitions may contain at most one file system and one extended partition, so this design permits us to nest the extended partitions to satisfy our volume requirements. Of course, each iteration creates yet another convenient hiding place for data.

If the partitions on a hard drive do not use up all of the available space, the remaining area cannot be accessed by the operating system by conventional means (e.g., through Windows Explorer). This wasted space is called volume slack (see Figure 1, item 3). It is possible to create two or more partitions, put some data into them, and then delete one of the partitions. Since deleting the partition does not actually delete the data, that data is now hidden.

Once partitions have been defined, we're ready to move up to the next layer and create an organizational structure for each partition. Before an operating system can store and access data within a partition, a file system must be defined. Modern operating systems support one or more native file systems. A file system allocates data in blocks (or clusters) where a block consists of one or more consecutive sectors. This allocation scheme allows a smaller number of references to handle larger amounts of data, but limits us to accessing data within the file system as block-sized chunks rather than sector sized chunks. Overall, this tends to make storage and access far more efficient

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Sthultz
Identity Theft and Financial Fraud Research and Operation Center

than referencing each individual sector. However, if the total number of sectors in a partition is not a multiple of the block size, there will be some sectors at the end of the partition that cannot be accessed by the operating system using any typical means. This is referred to as partition slack (see Figure 1, item 4) and is another place where data can be hidden.

Every partition contains a boot sector, even if that partition is not bootable. The boot sectors in non-bootable partitions (see Figure 1, item 5) is available to hide data.

Any space in a partition not currently allocated (i.e., unallocated space) , to a particular file (see Figure 1, item 6) cannot be accessed by the operating system. Until that space has been allocated to a file, it could contain hidden data.

It is possible to manipulate the file system metadata that identifies bad blocks (e.g. the File Allocation Table in a FAT file system or \$BadClus in NTFS) so that usable blocks are marked as bad and therefore will no longer be accessed by the operating system. Such metadata manipulation (see Figure 1, item 7) will produce blocks that can store hidden data.

Disk slack (see Figure 1, item 8) is a byproduct of a strategy to accelerate file management. Modern operating systems write data in complete “blocks” where a block could be a sector (the minimal addressable unit of a disk) or a cluster (same concept as block in Microsoft's terms). If a file is not an exact multiple of the sector size, the operating system must pad the last sector and, in some cases (with older operating systems), this padding is data from memory (hence the historical term “ RAM slack ”). Modern operating systems tend to pad this area with nulls. If the total amount of data written does not fill an entire block, the remainder of the block from the sector boundary of the last sector within the block actually used by the file to the actual end of the block will remain unused and will likely contain data from a previously deleted file (file slack). It may also be effectively used to hide ephemeral data.

All of the above applies to nearly every file system in common use today, including FAT/FAT32, NTFS and Linux Ext-based file systems. There are some potential data hiding places in Linux file systems that require a more detailed description.

Ext2 and Ext3 (ExtX) file systems are divided into sections called block groups. Block groups are used to store file names, metadata, and file content. Information about block group size and configuration is stored in a superblock at the beginning of the file system, copies of which are scattered throughout the partition. The block following the superblock (if present) or the first block in every group (if not present) contains a group descriptor table with group descriptors describing the layout of each block group.

An ExtX superblock has 1,024 bytes allocated to it and the last 788 bytes are unused. There also might be some reserved area behind the superblock, depending upon the block size. We call this Superblock slack (see Figure 1, item 9).

There is a reserved area behind the ExtX group descriptor since the group descriptor is only 32 bytes long and the block bitmap that follows it must start on a block boundary. This means there is a minimum of 992 bytes (1,006 if you count the padding at the end of the group descriptor) where data could be hidden and more if the block size is larger than 1,024. We refer to this as ExtX group descriptor slack (see Figure 1, item 10).

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Stultz
Identity Theft and Financial Fraud Research and Operation Center

ExtX directories are like any other file and are allocated in blocks. The space between the last directory entry and the end of the block is unused and can be used to hide data. This is directory slack (see Figure 1, item 11).

Figure 2 is a graphical illustration of the relative volatility of the various data hiding areas discussed above (cf. also, [10]). The degree of persistence of the hidden data is dependant upon the characteristics of the particular area where it is hiding and the type of disk activity that has occurred since the data was written there. Figure 2 illustrates the relative persistence through normal disk activity and possible re-partitioning and/or re-formatting. The numbers in the figure correspond to the numbered items in Figure 1.

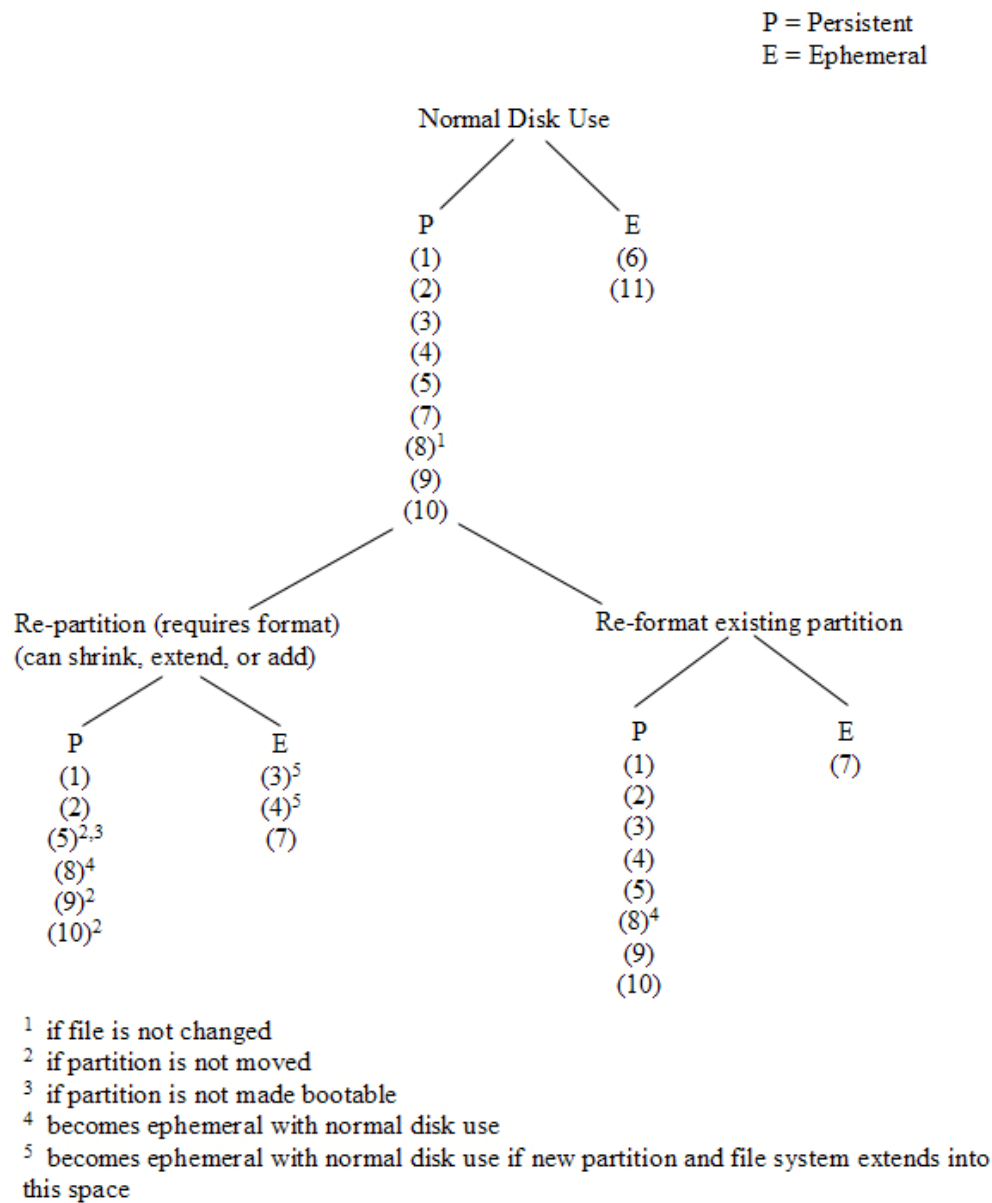


Figure 2: Relative Volatility of Data Hiding Areas

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Sthultz
Identity Theft and Financial Fraud Research and Operation Center

NTFS file systems also offer some unique opportunities for data hiding. The NTFS file systems used today contain innovations that provide efficient file access (for instance, B-Tree organization of files within directories) and readily accessible metadata files to manage disk organization (Microsoft's version of resource forks called Alternate Data Streams), and some other small file storage oddities as well.

When seeking to hide data, there are various strategies that might be employed. As mentioned in a more general sense, metadata manipulation may be used to conceal covert data in bad clusters (\$BadClus). In fact, this same concept can be extended easily on an NTFS file system by working directly with the \$Bitmap file.

The \$Bitmap file contains a complete map marking the allocation status of every addressable cluster in the partition. Should a consistency check be run, it would become obvious should someone modify this table to hide data, but otherwise this provides a wonderful avenue for hiding data in a way that allows that data to persist for the life of the file system. Depending upon the purpose, these are far better approaches than using file slack which persists only for the life of the file.

NTFS provides some other nooks and crannies. For instance, Alternate Data Streams are actually additional \$FILE entries associated with a parent file record within the Master File Table. We could not find Microsoft documentation regarding the number of alternate data streams that may be associated with a file or folder, but empirical testing conducted by Jason Fossen (www.fossen.net) suggests that the maximum is 4106 regardless of the size of the ADSs themselves. These streams are not identified by typical file management tools (e.g., Windows Explorer) so are hidden at that level. However, several utilities are available that report and manipulate alternate data streams [1]. Alternate data streams persist for the life of the attached file or folder as long as that file or folder remains in an NTFS file structure.

Small files also offer some interesting possibilities, especially when considered in conjunction with alternate data streams. NTFS has a rather unusual capability in that if a file is so small that the entire content of the file can fit within the boundaries of the Master File Table entry for the file, NTFS will store the file there. This is, of course, a convenient place to hide data. But it becomes even more interesting when a file is deleted. When a file is deleted, the clusters that were in use are released to the file system for reallocation. This includes the MFT entry itself. What if one were to create several thousand files, thus consuming several thousand MFT entries; once these files were created, a final file could be created that fits entirely within the MFT. All of these files are now deleted. In essence, we have created a piece of hidden disk space that resides within an allocated file (\$MFT and \$MFTMirror) that will persist until enough files are created to reuse that particular MFT entry. For further detail, see [2].

Forensic Implications

The implications of this type of intentional data hiding can be serious in the context of forensic analysis ([4],[11],[12]). Typically, forensic analysis of systems reveals that bad actors do not often take any extraordinary means to hide data beyond, perhaps, either wiping the media or encrypting the data. In fact, in most cases the data is not even encrypted or wiped. What impact would a deep knowledge of the on disk structures coupled with a desire to hide data have on an analysis? To consider this, we have created several sample disk images and have embedded two pieces of data using the methods discussed in this article. After embedding the data, we used a tool that is commonly used by law enforcement for forensic analysis to determine how easily an analyst could uncover this hidden data.

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Stultz
Identity Theft and Financial Fraud Research and Operation Center

The data that was hidden was a text string, "Hidden Message" and a GIF image containing an image of the word "Hidden". Please understand, of course, that if the investigator already knows something of what he is looking for, for example, the word "Hidden", the difficulty of recovering this data is trivial. What we are seeking to measure is whether or not the tool alerts the analyst to the presence of data in unusual locations and how likely it is that the analyst would discover the materials.

The first test was run using FAT16. As a form of control, an image file was copied to the disk using the typical file system tools. Data was hidden using two simple methods: First, the FAT16, when formatted, automatically preferred a cluster size of 4096 bytes. The result of this was that the boot sector of 512 bytes is immediately followed by three empty sectors. Our "Hidden Message" text was inserted here. In order to conceal the presence of the GIF image, the first FAT was modified so that clusters 24 through 29 were marked as bad (See Figure 3). Finally, the GIF image was placed onto the disk beginning in cluster 24.

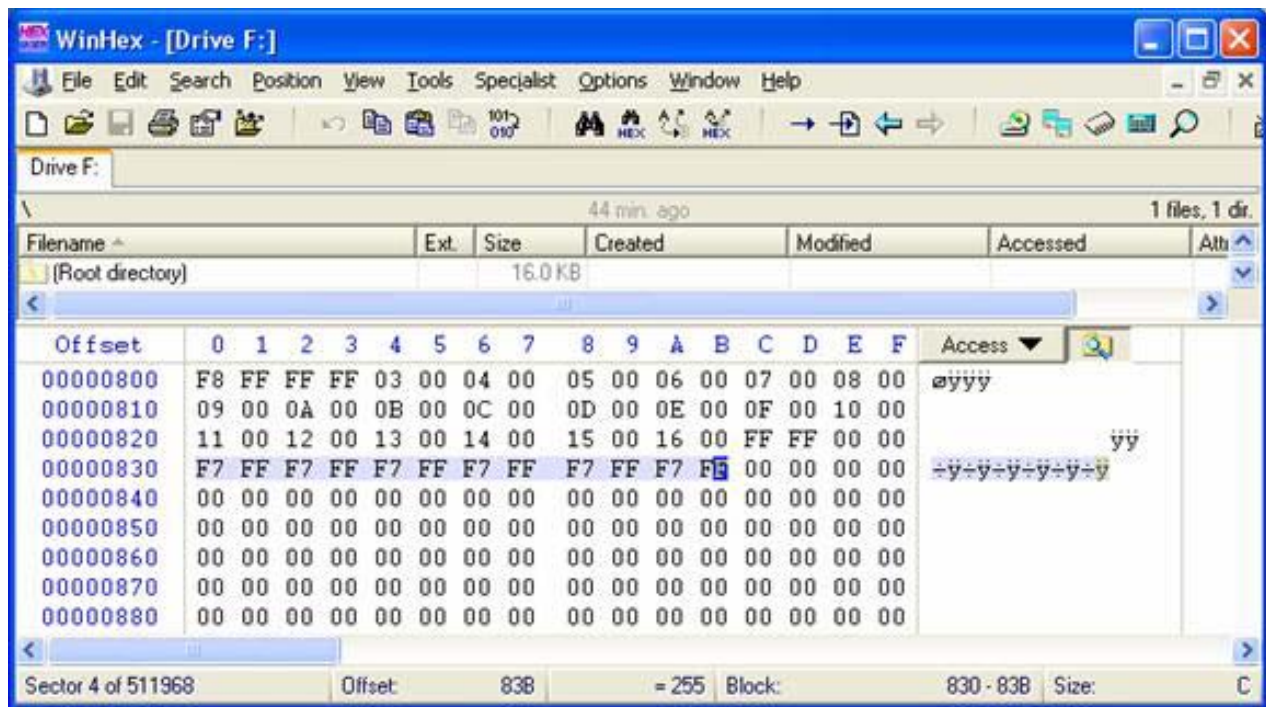


Figure 3: File Allocation Table with clusters 24-29 marked "bad."

With our data hidden, we used AccessData's Forensic Toolkit (better known as FTK) to acquire an image of the drive and perform data carving and full text indexing. See Figure 4.

Of particular interest to us were the results from the file carver. File carvers are forensic tools (or data recovery tools) that analyze the data on a drive (typically sector by sector or block by block) without any regard for the actual logical organization of the disk. Each sector or block is then checked for known file signatures and, if one is found, the data is extracted or marked appropriately. One of the reasons that a GIF image was selected is that it has a very standard and easily recognizable "magic number" or fingerprint: Almost every one of them begins with the marker, "GIF89a". While FTK did find

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Stultz
Identity Theft and Financial Fraud Research and Operation Center

the file and mark it as a “Lost File Chain” marked as bad sectors, it did not successfully identify this as an image.

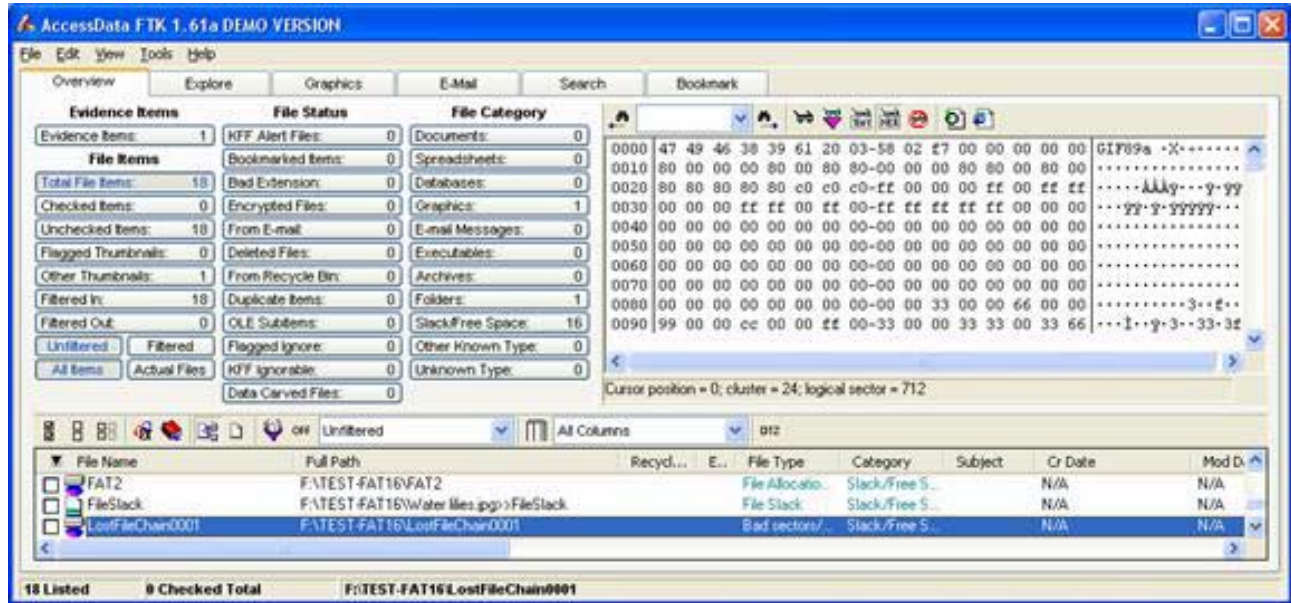


Figure 4: Forensics Toolkit results showing hidden GIF file.

Even though the “hidden message” is sitting in what should be a completely barren segment of the drive, FTK makes no special mention of this text. Unless the analyst is in the habit of examining all reserved sectors for additional data or already knows to search for a keyword of “Hidden”, it is very unlikely that this data would be found. From our experience, this is typical behavior for a forensics toolkit. In short, it is unusual for any forensic analysis tool currently on the market to draw the analyst's attention to a piece of data that is residing in an atypical location.

The next set of tests was performed using an NTFS file system with a cluster size of 4096. Data was hidden using several of the techniques discussed. First, in the reserved area between the primary partition table and the first partition, “Hidden Message” was again embedded. Since there was additional space, we also embedded our image file, but this time the file was not embedded beginning at the beginning of a sector. Many file carvers speed their operation by only examining the first few bytes of a sector or only the first sector in a cluster.

As an additional test, the \$Bitmap file was modified to mark the final clusters of the volume as “In Use” and the GIF image was copied into these now reserved clusters.

The “Hidden Message” was also embedded using the Alternate Data Stream technique discussed.

To accomplish this, a host file was created followed by 1000 associated streams. The data that we wished to conceal was then copied into what amounts to stream 1001 and the host file was then deleted. Finally, additional files were copied onto the drive to obscure the original host file MFT entries, but with the cushion of 1000 alternate data streams, the original hidden data would still be untouched. With all of these tasks completed, the drive was imaged and analyzed using FTK. FTK was able to identify and carve the hidden data in the slack space between the partition table and the first partition,

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Sthultz
Identity Theft and Financial Fraud Research and Operation Center

even though it wasn't on a sector boundary. Unfortunately, while it was successful with the image file, it still failed to "notice" the hidden text. What if the file that was hidden didn't match a known fingerprint or were encoded in some way? FTK would fail to notify the analyst. It was also successful in identifying the image file tucked away at the tail end of the disk in space that we "reserved". While this is good, there was no notification that space was marked "In Use" that was not actually allocated to any disk structure or file. Obviously, the same problem exists as does for the gap between the partition table and the partition.

What about the alternate data streams? FTK, surprisingly, made no mention of the fact that there were still approximately 800 remnants of alternate data streams floating around in the MFT and, of course, did not notify us to the presence of the "Hidden Message" unless we already knew what to search for.

Overall, this means that for us to be successful today with forensic investigations, we need really sharp analysts who are willing to go above and beyond to find data when we are dealing with clever malcontents. Our tests were performed on a small drive (256 megabytes). Obviously, hunting down data on a 500 gigabyte drive is nigh on impossible without good tools and good initial information. In the long run, as drives continue to increase in size and those who desire to conceal activities become better at their trade, the forensic community will be at a greater and greater disadvantage unless we become proactive in this "Arms Race" and move away from strict signature based analysis of media and start building in some sort of anomaly detection capability.

Perspectives

It should be pointed out that this discussion only includes methods of hiding data within the media and file system structure and does not address other data hiding methods such as:

- Altered BIOS Parameters
- Registry entries using new keys or unused keys
- Swap files
- Renamed files (e.g. as .dll)
- Binding of one executable file to another
- Steganography
- Hiding data within documents (e.g. as metadata or using a white font)
- Hiding data within html files
- Merging Microsoft Office documents
- Encrypted files
- Compressed files

Conclusion

Knowing how data can be hidden within the media and file system structure means knowing how that data can be found. Note that we are only talking about the data hiding mechanisms that have been discussed here; methods such as encryption and steganography present their own sets of analysis problems.

Many of those that attempt to temporarily hide data or use a computer with nefarious intent are aware that there are ways in which this hidden data can be found. They will therefore use one of the available "disk wiper" utilities in an attempt to eliminate the evidence. What they are not aware of is that many of these utilities are ineffective in eliminating all hidden data. See [2] for more information.

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Sthultz
Identity Theft and Financial Fraud Research and Operation Center

A forensic analysis tool can be as simple as a hex editor. While effective, this approach is very labor intensive given the size of today's hard drives. There are many commercial and open source digital forensic tools available that will automate this process. However, one must exercise caution when using these tools, since not all of them are equally effective in finding all types of hidden data. It is often advisable to use more than one of these tools to perform a digital forensic investigation. For further information, see [5], [7], [8].

Finding hidden data is further complicated by the fact that there are a number of anti-forensic tools being developed. One source is the research being done by the Metasploit Project (www.metasploit.com/projects/antiforensics). An example is their Slacker tool that automatically encrypts and hides a set of data within the slack space of multiple files that are known to be unlikely to be changed. Data thus hidden cannot be detected by currently available digital forensic analysis software.

Data hiding becomes an increasingly important topic as the sophistication of information technologists on both sides of the law increases.

REFERENCES:

- [1] Berghel, H. and Natasa Brajkovska, "Wading through Alternate Data Streams," Communications of the ACM, 47:4 pp. 21-27 (2004).
- [2] Berghel, H. and D. Hoelzer, "What Does a Disk Wiper Wipe when A Disk Wiper Does Wipe Disks," Communications of the ACM, 49:8 [in press].
- [3] Berghel, H. and L. O'Gorman, "Protecting Ownership Rights through Digital Watermarks," IEEE Computer, 29:7 pp. 101-103 (1996).
- [4] Caloyannides, Michael A., Computer Forensics and Privacy , Artech House, 2001
- [5] Carrier, Brian, "Defining Digital Forensic Examination and Analysis Tools Using Abstraction Layers," International Journal of Digital Evidence, Winter 2003, Volume 1, Issue 4
- [6] Carrier, Brian, File System Forensic Analysis , Addison-Wesley, Upper Saddle River NJ , 2005
- [7] Carvey, Harlan, Windows Forensics and Incident Recovery , Addison-Wesley, 2005
- [8] Casey, Eoghan, Handbook of Computer Crime Investigation , Academic Press, 2002
- [9] Cole, E., Hiding in Plain Sight: Steganography and the Art of Covert Communication , Wiley Publishing, Indianapolis , 2003.
- [10] Farmer, Dan and Wietse Venema, Forensic Discovery , Addison-Wesley, 2005
- [11] Kruse II, Warren G. and Jay G. Heiser, Computer Forensics, Incident Response Essentials , Addison-Wesley, 2002
- [12] Nelson, Bill et al, Guide to Computer Forensics and Investigations, Second Edition , Course Technology, 2006

Data Hiding Tactics for Windows and Unix File Systems

Hal Berghel, David Hoelzer, Michael Sthultz
Identity Theft and Financial Fraud Research and Operation Center

[13] Singh, S., The Code Book , Anchor Books, New York , 2000.

[14] Wang, H. and S. Wang, "Cyber Warfare: Steganography vs. Steganalysis," Communications of the ACM, v. 47, no. 10, pp 76-82.