# Global Information Assurance Certification Paper

## Copyright SANS Institute
## Author Retains Full Rights

## Interested in learning more?

Check out the list of upcoming events offering
"Advanced Incident Response, Threat Hunting, and Digital Forensics (Forensics
at http://www.giac.org/registration/gcfa

# Google Desktop Search as an Analysis Tool:

## An overview of how VMWare and Google Desktop Search can be leveraged to analyze the GDS cache on a suspect system.

## Chris Poldervaart

**November, 2006**

# Introduction:

As hard drive prices continue to fall and storage sizes increase, users are saving more and more data to their hard drives without having to deal with the house-cleaning rituals of years gone by. There is no longer the need to delete unused files or applications just to make room for new items.  Of course, this presents a problem to the users when they need to locate a particular document in the sea of information residing on their computers.

In comes the desktop search tool.  These applications provide an indexing feature that allows a user to search their entire document store with nearly instantaneous results.  Google Desktop is just one of the many options available to users to meet their computer indexing needs.  Copernic, MSN, Yahoo, dtSearch, X1, and 80-20 (now Verity) are a few others with similar capabilities.  Google's Desktop Search may be the most popular for two reasons.  The product is free, and there is no denying the brand name recognition that Google has achieved with regards to search capabilities.

While the desktop search capabilities were designed to make a user's life easier, I have found through testing and experience that they can also be used by a computer forensics examiner to recover information that may normally not be available otherwise.  This paper originates from a genuine need to conduct an analysis on a computer that had been taking advantage of the Google Desktop Search capabilities.  The methodology here is fairly straightforward and would likely carry over to many, if not all, of the other desktop search programs.  This paper focuses on the analysis of Google Desktop installations on Windows XP machines.

# Background:

## *How does Google Desktop work?*

Google Desktop analyzes files for content and sorts the text into a database that is referenced when a user performs a search, as well as copying the information to a cache that allows for easy viewing of the files contents without having to open the original document.  After the initial drive indexing is complete, files, web pages, email, and some chats are indexed on the fly as the user opens them.  Google Desktop also indexes network storage locations and removable media.  This is especially helpful to the forensic examiner because documents not residing on the local machine might appear in the Google Desktop cache.  Another interesting tidbit to note is that encrypted documents are indexed and cached when decrypted by user.  So, while the file may exist in an encrypted state on the hard drive, there may be a decrypted copy in the Google Desktop cache.

### What are the key components of interest during an investigation?

The components of most interest during a suspect system investigation are all contained within the user's "Google Desktop Search" folder. This folder holds the index information for the user's profile. Several of the files in this folder contain path, URL, email address, and email subject information in a Unicode format. It seems that some combination of these files located within the "Google Desktop Search" folder is required by the application to reconstruct the cache. It was my experience that best results were achieved by simply utilizing the entire folder structure and contents. It appears as if the Google Desktop cache does not exist in a plain text format, but rather stored in some form of encrypted or compressed format. There has been a recent paper written regarding the Google Desktop Search database by Sriram Krishnan that also suggests that this is the case. This explains why many of the keywords that may be located using the Google Desktop Search API are not recoverable by typical disk search techniques.

### Where are these components located?

The "Google Desktop Search" folder is typically located in the user's "Application Data" folder. For example:

C:\Documents and Settings\[USERNAME]\Local Settings\Application Data\Google\Google Desktop Search

This cache storage location is configurable by the user, and the following registry key can be queried for the exact location:

HKEY_CURRENT_USER\Software\Google\Google Desktop

The value for "data_dir" will display the cache storage location. This and other data values of interest are detailed below, in the section titled "Where are the settings that show the preferences the user has set?

Google Desktop Search as an Analysis Tool
Chris A. Poldervaart

3

### *What files are indexed with Google Desktop?*

According to the Google Desktop support site, "In all versions, Google Desktop searches the full text of PDF, TXT,HTML, Microsoft Word, Excel, and PowerPoint files, metadata for audio and video files (such as artist and album information), and file names for most other files on your internal hard drives and the external in Google Desktop 2 or networked drives in Google Desktop 2 you've told Google Desktop to Search." There is also the option disable the indexing of password protected Office documents and/or HTTPS web pages.

### *What files are not indexed with Google Desktop?*

Again, relying on the Google Desktop support site, we learn that "Google Desktop doesn't search hidden folders and some directories which contain system and temporary files or files within your "All Users" shared directory."

### *Specific limitations to consider:*

A 100,000 file index limit exists, according to Google. "In all versions, Google Desktop will only index 100,000 files per drive during the initial indexing period. If you have more than 100,000 files in a particular drive, Google Desktop will not index them all. However, Google Desktop will add files to your index during real-time indexing if you move them. If the files are Microsoft Office files, Google Desktop can also add them to your index if you open them." In many cases this may not pose any problem at all. In other cases, it may make a difference with regard to a successful recovery effort.

There is also a 10,000 word limit within each individual file. "...if you're searching for a word within the file, please note that Google Desktop searches only about the first 10,000 words. In a few cases, Google Desktop may index slightly fewer words to save space in your search index and on your hard drive." It is not exactly clear to me what methods are used to determine what words are indexed, and which are not. My personal experience in this area has led to strange anomalies. During my search effort I would not get hits on an email that I knew contained a certain keyword, but might get a hit on a different keyword in the same file. It was

Google Desktop Search as an Analysis Tool
Chris A. Poldervaart

4

inconsistent from message to message.  For instance, an email might have the following text:

"I really hate you and I want to kill you."

I would run a search on "hate" and come up with nothing.  Then later I would run a search for "kill" and get a hit on the email that clearly should have appeared with the "hate" search.  I suppose the reason for this is that Google Desktop treats the entire email store, whether pst or ost, as a single file, causing the 10,000 word limit to take effect, thus generating the inconsistencies with what is exactly cataloged in the indexing process.  I have not been able to confirm the precise reason for this anomaly.

## *Google Desktop Database Encryption*

Encryption of the Google Desktop database may also pose a problem for an analyst expecting to uncover nuggets of information from this valuable data repository. While database encryption is turned off by default, it is very possible that this setting can be set by the user.  While not a complete death certificate to the operation, it certainly presents another hurdle to leap.  The database encryption feature is designed to ensure that an unauthorized user is not able to access the GDS cache of another.

Google Desktop's encryption feature can only be utilized on NTFS volumes.  It utilizes the Windows EFS (Encrypting File System) to accomplish this capability, and thus would occur at the file system level above any imbedded encryption/compression utilized by GDS itself.  In order to decrypt an encrypted Google Desktop database, one would need the user's NT login password.  With that and any number of tools designed to decrypt EFS files, access to the database should be possible.  It would then need to be saved in the decrypted format for access by the GDS application.

Google Desktop Search as an Analysis Tool
Chris A. Poldervaart

5

## *Where are the settings that show the preferences the user has set?*

The registry settings for Google Desktop can be found in the user's NTUSER.DAT file. This file is located in the user's directory (C:\Documents and Settings\USER\NTUSER.DAT). The key you want to locate is \NTUSER.DAT\Software\Google\Google Desktop. Within this key are some settings of interest to the examiner:

| | |
|---|---|
| data_dir | Designates where the index and cache are stored on the user's system. It is stored in the following folder by default: C:\Documents and Settings\USER\Local Settings\Application Data\Google\Google Desktop Search |
| file_extentions_to_skip | Provides a list of file extensions Google Desktop will explicitly ignore during the indexing process. |
| installtime | This is the installation timestamp for Google Desktop. It is stored in Windows 64 bit Hexadecimal – Little Endian format. An easy way to decode this timestamp is to copy the hex value and decode it with "Decode" available for free at: http://www.digital-detective.co.uk/freetools/decode.asp |

## *How do we use this to our advantage?*

We take the index from the suspect computer and search against it for data that is relevant to our investigation. This is the same index the suspect used to facilitate the rapid searching of his own computer.

## *Real world examples:*

In a fairly recent investigation, the one that actually prompted the research and process outlined in this paper, there was a specific need to recover email conversations between the suspect and another person. The suspect was interested in the protection of the privacy of these email conversations, and was consistent in deleting them from his Inbox, Deleted Items folder, and his Sent Items folder. While some email was recovered in the unallocated space of the hard drive
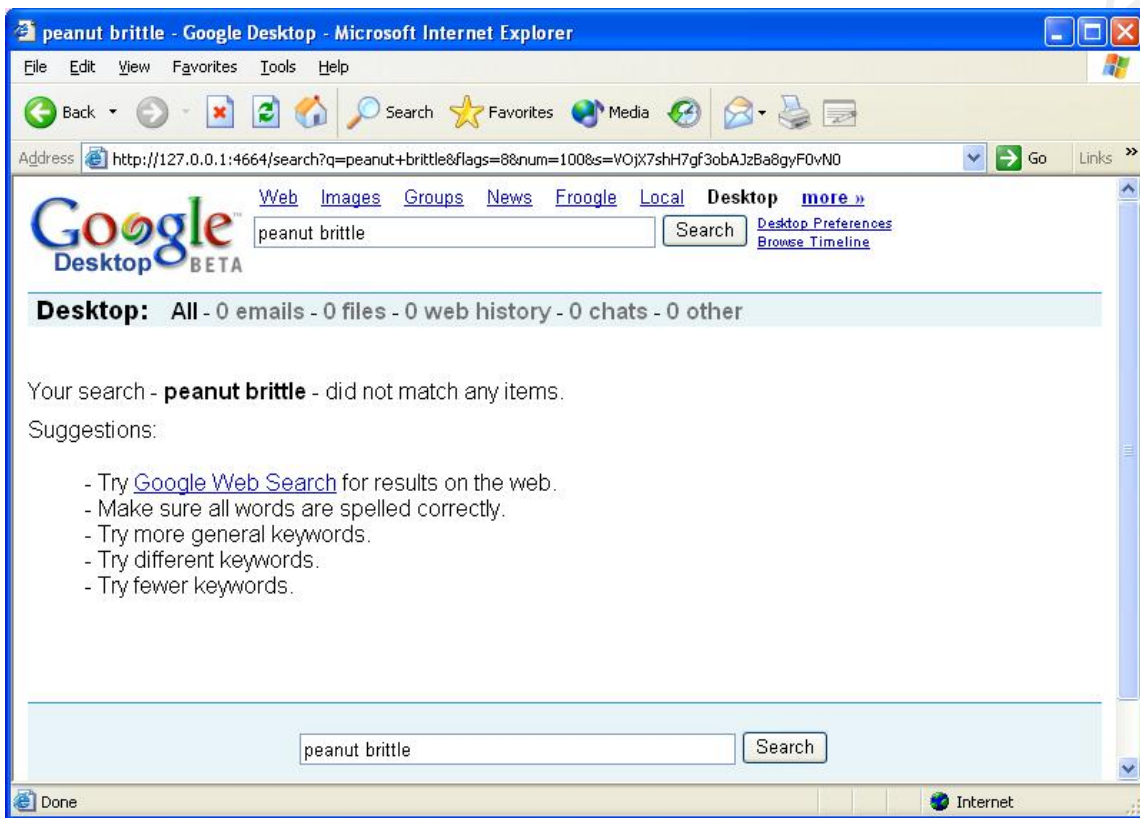
(using keyword searches in the Outlook Compressed Encryption format) we were still missing some key messages to the investigation. There was, however, some relevant search hits (although not in a usable context) located in a file named "dbeam." This file was clearly associated with the Google Desktop Search tool.

Knowing that Google Desktop keeps a keyword database and a cache of indexed email messages, an attempt was made to retrieve this information to put the relevant keywords back into context. Using the process described in this paper, numerous email messages were recovered from the suspect computer that were not available before. While the suspect was diligent in deleting suspicious messages, he neglected to realize that Google Desktop was silently caching these in the background even as he worked to destroy them. Keyword searches using the Google Desktop Search software revealed complete email conversations stored in the cache, including basic email header information, that took the investigation to a level not before realized.
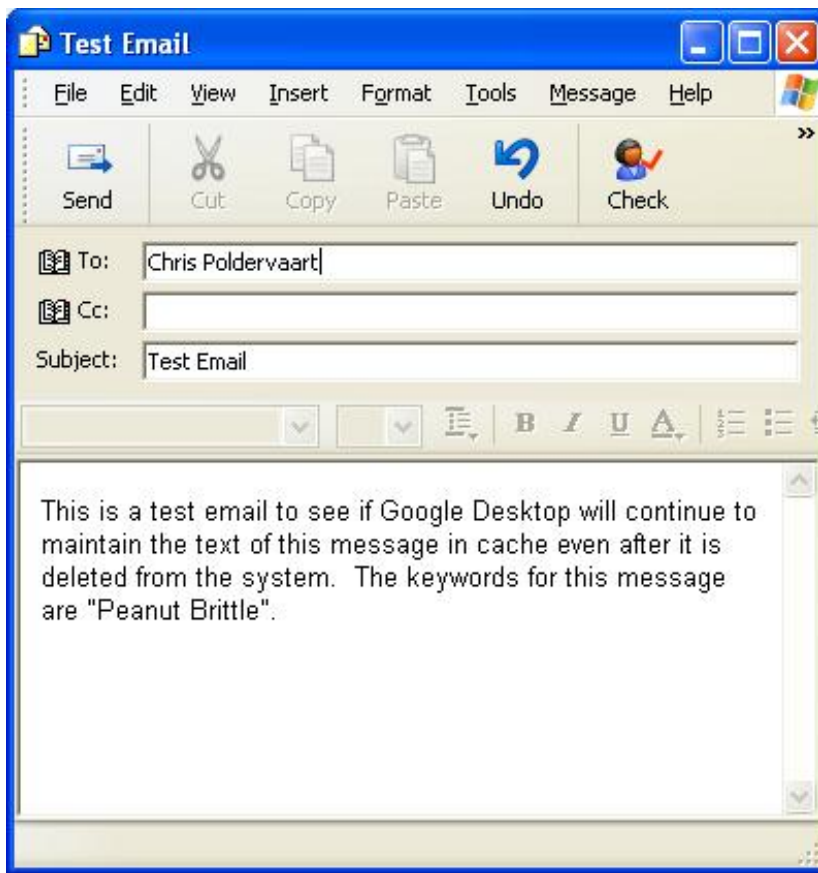
## Testing and Validation:

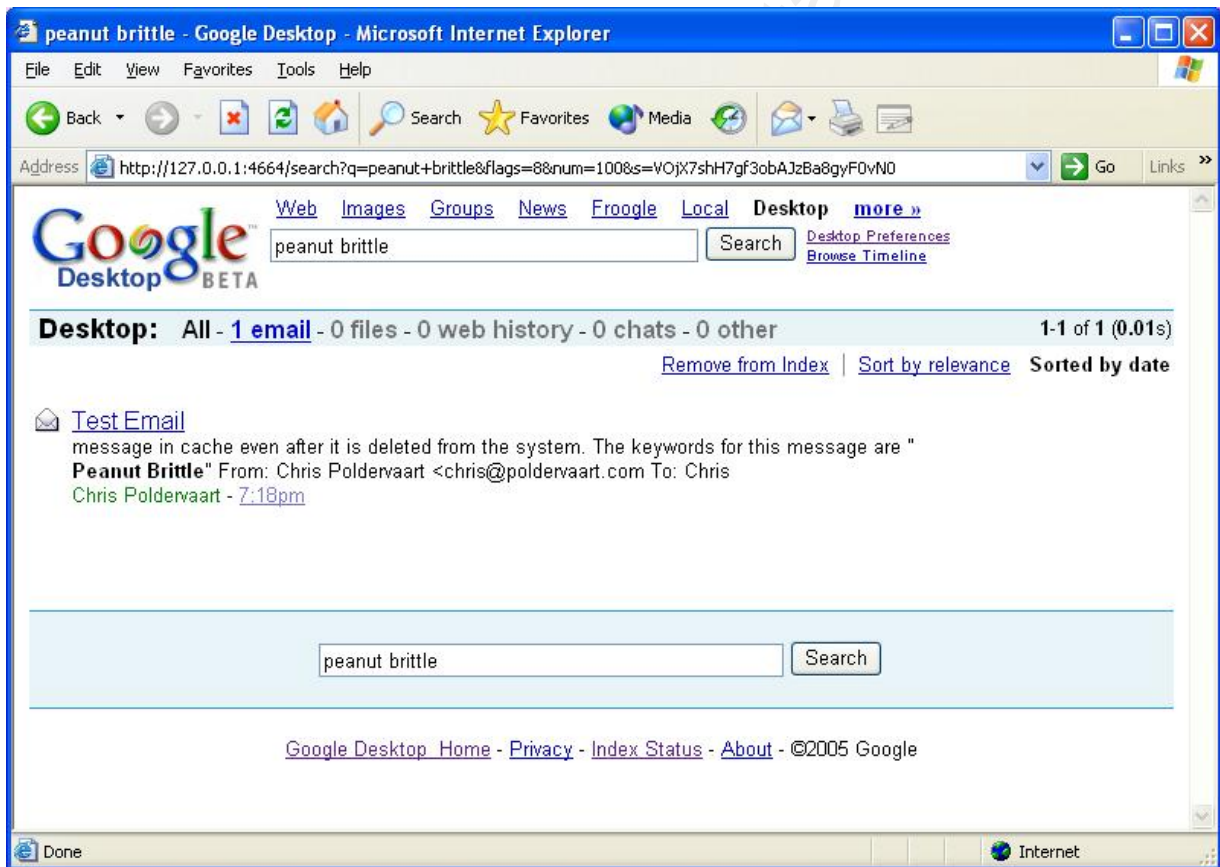I performed some initial testing to validate the real-world findings.

The first step was to utilize a fresh Windows Installation with Google Desktop Search installed immediately after. I issued a search for the phrase "peanut brittle" to verify that there were no existing results from the GDS cache:
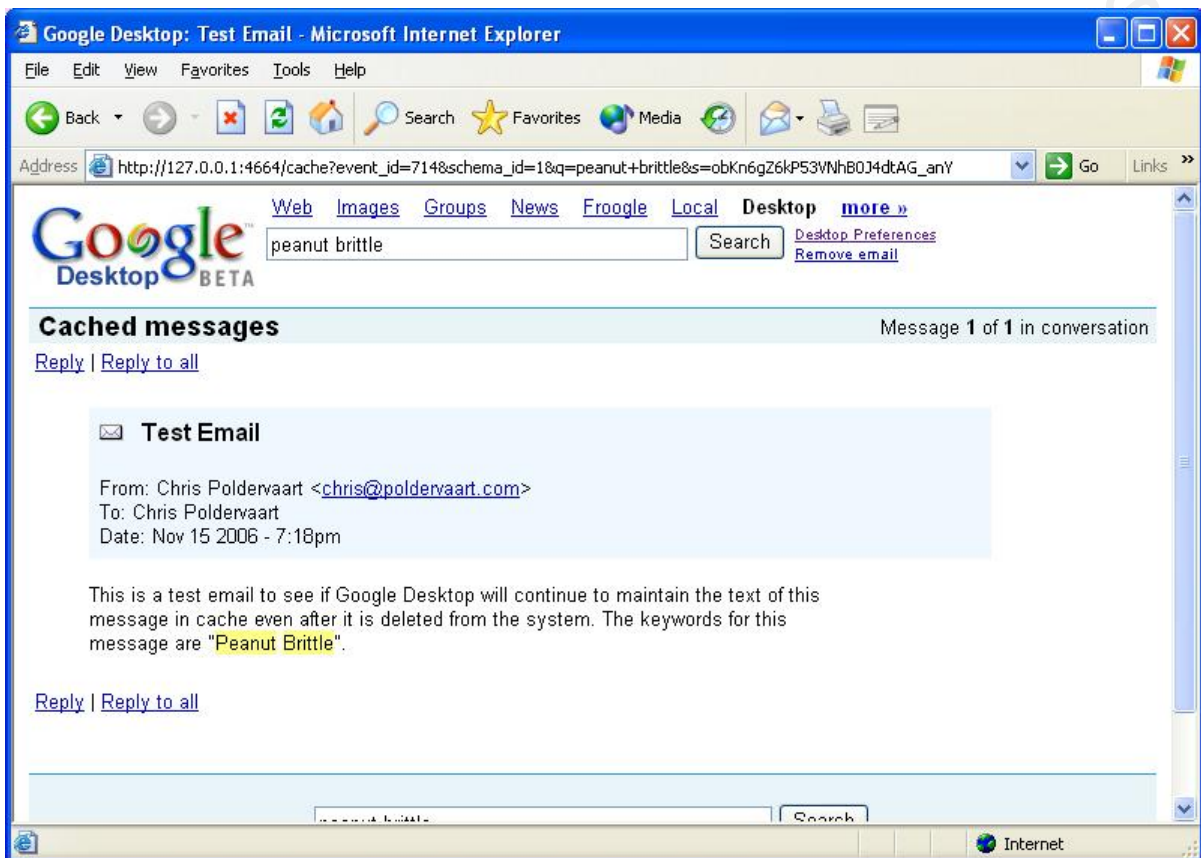
Next, I created a test email using Outlook Express with the subject "Test Email" to
"Chris Poldervaart". This email was "sent" and a copy was saved in the "sent Items"
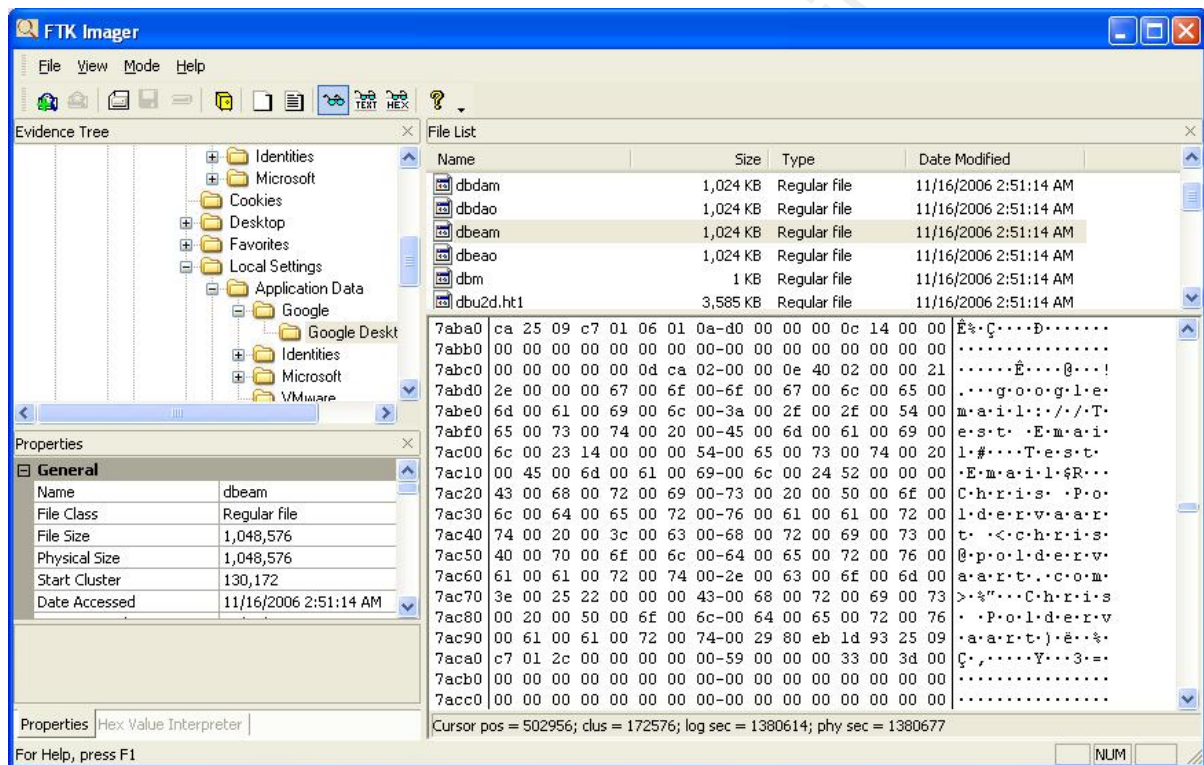folder.

Google Desktop Search as an Analysis Tool                                    9
Chris A. Poldervaart

I then re-issued the search for "peanut brittle" in Google Desktop Search to verify that the email message was in fact indexed.
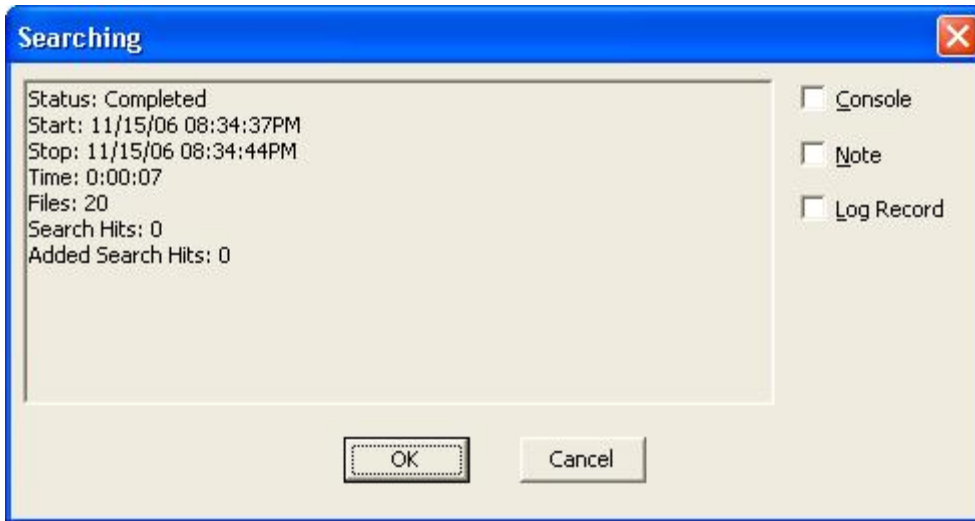
Clicking on the search hit revealed the text of the email in more detail. It is important to note that this email text is stored entirely in the GDS cache, and not pulled from the original email storage location (pst/ost). This will be clear when finding the text remains after complete deletion of the original email message.
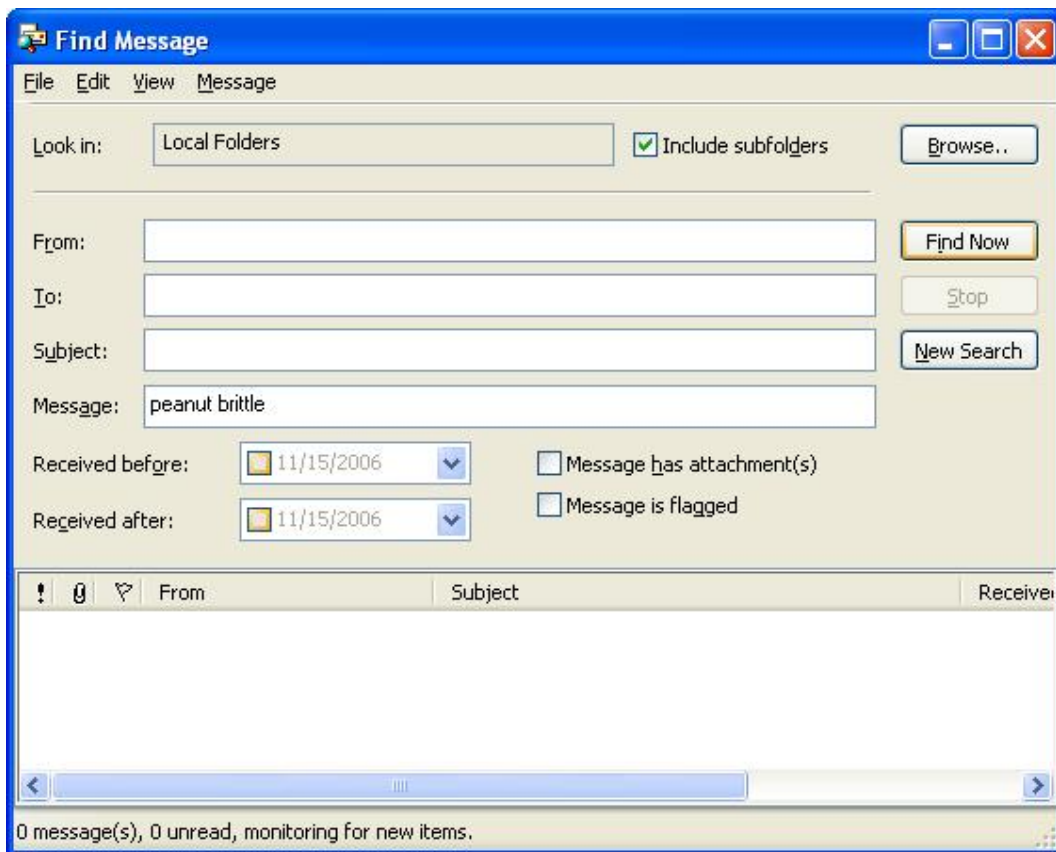
The only textual indication of the cached email could be found in the dbeam file in a Unicode format. The email subject, sender, and recipient were all indexed in this format. This can be seen by viewing the file using FTK Imager or a similar tool.
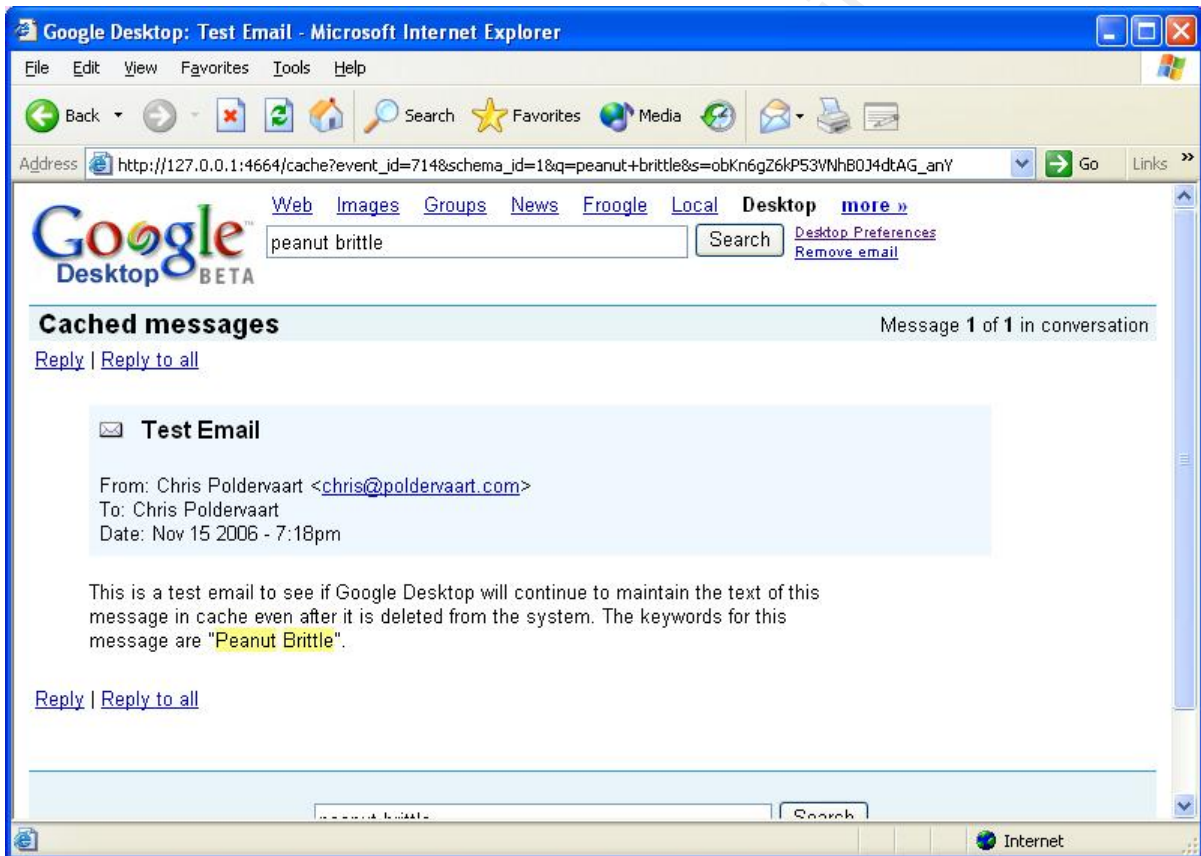


A search of all 20 files in the GDS folder for the words "peanut" or "brittle" yielded zero results, as seen in this EnCase search result. This further confirms the likelihood of encrypted or compressed cache storage.

Searching

Status: Completed
Start: 11/15/06 08:34:37PM
Stop: 11/15/06 08:34:44PM
Time: 0:00:07
Files: 20
Search Hits: 0
Added Search Hits: 0

☐ Console
☐ Note
☐ Log Record

OK    Cancel

After confirming what was available in the cache (and not visible), I deleted the message from the email storage location. I then confirmed it was not available by searching for "peanut brittle" using the default Outlook Express "Find" utility. There were no identified results.

The result of the GDS search after the message deletion indicates that the entire text of the email continues to exist in the GDS cache even after the original message was deleted. This confirms the real-world findings mentioned above.

# Process:

The process behind analyzing a suspect's Google Desktop cache is really quite simple. It is not the typical solution in the realm of computer forensic investigations, but it is important to remember that there are many techniques that can be utilized to elicit the necessary results even if it means not sticking with those tools designated as being designed specifically for "forensics." It takes stepping out of the box to see how to use the resources that are available to you when other typical methods fail to meet the needs of the analysis. The key is to validate your results and document your steps.

### *Virtual Machine Preparation:*

We need to do several things to prepare for searching the suspect's Google Desktop index and cache. First of which is to prepare a safe machine that is clean of residual data from other investigations, and also free of any pre-established Google Desktop data. I chose to make use of a VMware virtual machine for this. I created a virtual machine with the necessary and minimal software installations required to complete the task. I also disabled networking from the machine after my setup was complete to ensure that search results are not affected by Google Desktop calling home or performing any other calls to the Internet during the process. Once the virtual machine is appropriately configured and saved, the machine was cloned for each new investigation. In other words, a new and independent machine was used for each Google Desktop analysis, leaving the original machine untouched and available as a baseline. The specific process used to set up the virtual machine is detailed in the section titled "VMware Virtual Machine Setup Detail" located in the appendix to this document.

### *VM Software Installation:*

The necessary software for the process is the guest operating system, the Google Desktop software, and optionally, a forensic tool to retrieve the required files from

the suspect image from within the virtual machine.  While there are certainly other methods to get the suspect's Google Desktop files onto your virtual machine, I found a certain amount of simplicity in installing Encase Forensic Edition in the VM environment.  The VM will appropriately recognize the USB dongle, and renders a fully operational EnCase installation.  Also, in setting up the VM, I made sure to create shared folders that allowed me to access the suspect image files from within. I simply opened up my EnCase case file from within the VM, located the appropriate folders for export, and exported them directly to my virtual machine.

### *VM Google Desktop Preparation:*

Setting up the Google Desktop software in the VM was pretty straightforward.  I performed the standard installation, and even allowed the initial index to complete.  I disabled the option for Google Desktop to run on startup.  I also browsed to my Google Desktop cache storage location and deleted my profile and subsequent cache and index.  I now have a Google Desktop installation, but no data.  This is fine, since we will be importing the data from the suspect's image.  Remember, if these steps are performed during the creation of the baseline VM, then the working clones will be properly configured with minimal work left to do to get up and running searching through your suspect's cache.

### *Moving the target system's Google Desktop files to the virtual machine and use:*

As I mentioned earlier, there are several ways to move the suspect's Google Desktop Search folder structure to the Virtual Machine we have just created.  You could utilize shared folders between the Virtual Machine and the Host system.  You could use removable media, such as USB devices or DVDs (size permitting).

I have had great success using the Encase Forensic product within a VM.  I am sure that other similar forensic tools would work equally well.  I use it in the following manner:

1.  Set up Virtual Machine with Shared Folders pointing to a swappable drive bay.

Google Desktop Search as an Analysis Tool 18
Chris A. Poldervaart

2.  Insert working drive containing forensic image of suspect system into drive bay.

3.  Using forensic tool in VM, load suspect computer image into case.

4.  Export entire "Google Desktop Search" folder tree into Virtual Machine, making sure the folder structure is replaces existing structure for the VM installation.

5.  Utilize the Google Desktop search in the VM to search the suspect's cache.

## *Conclusion*

That's it.  It really is that simple.  As I mentioned before, this was very successful in an actual investigation.  The suspect was very savvy about cleaning up email and documents, but neglected to realize that Google was keeping track of everything for later retrieval.

This process can be applied to many different situations, but the point it that the investigator should find a way to take advantage of the data aggregation and preservation that a suspect user may be inadvertently be conducting for you.

Since I started this paper, Google has released the GDS API, which no doubtedly could be leveraged in ways that make this technique obsolete.  I simply haven't had time to look into it any further, and this process has served me well.
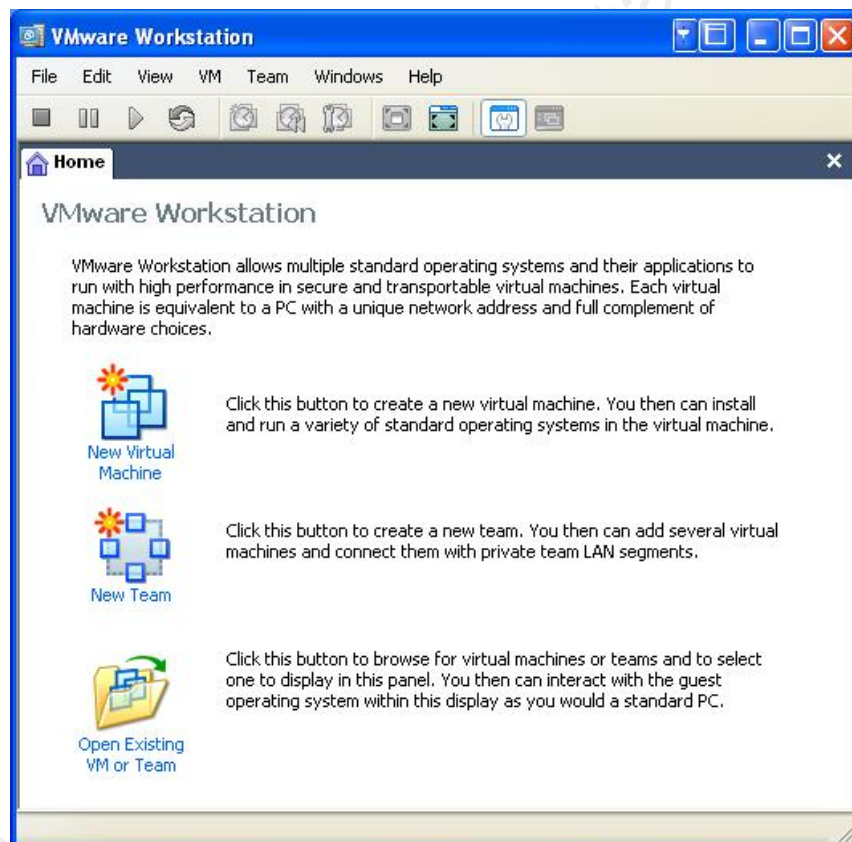
# Appendix:

## *VMware Virtual Machine Setup Detail:*

1.  To begin the process of viewing a suspect's Google Desktop cache from within your virtual machine, you must first install the VMware software on your

system. The VMware installation is fairly straight forward, and installation instructions are provided with the software.
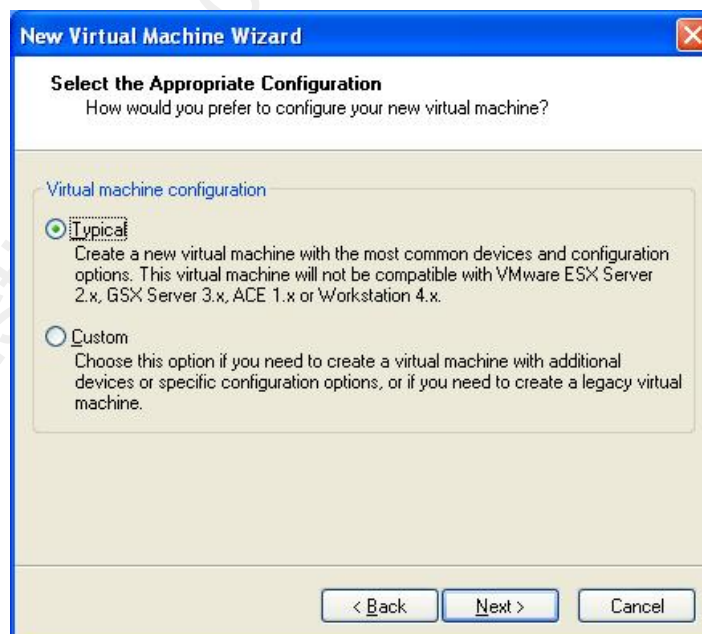
2. Once VMware is installed, you will then create a new virtual machine that will be used for your Google Desktop installation. This is accomplished by either clicking on the "New Virtual Machine" icon from the VMware home screen, or selecting File □ New □ Virtual Machine (CTRL+N).
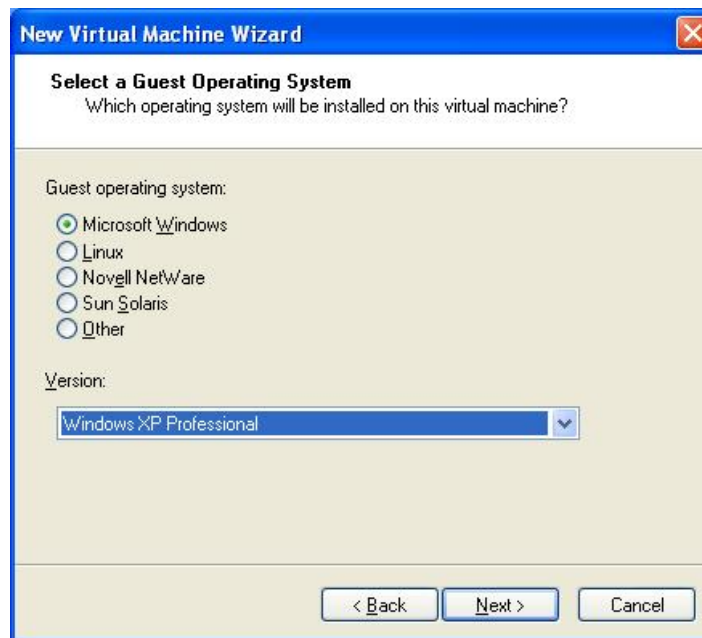


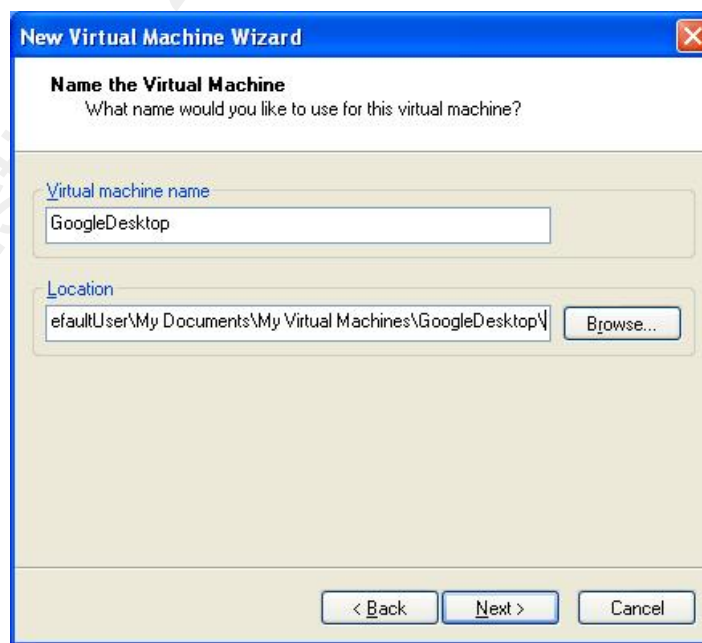3. You will then be taken through the "New Virtual Machine Wizard."

4. Click "Next" to begin setting up your Virtual Machine. For the purposes of setting up a machine to view the Google Desktop cache, you can select the "Typical" configuration option.
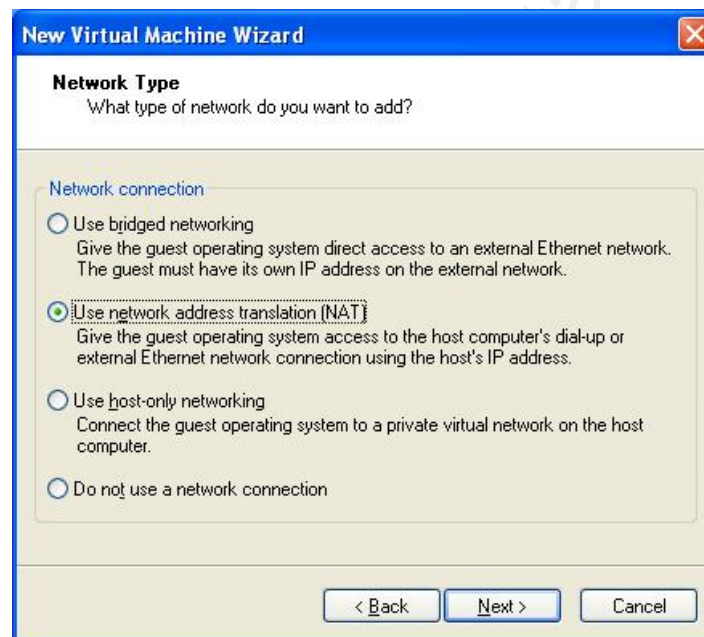


5. You will then select the Operating System that you intend to use on the Guest machine. In this case, I will be using Windows XP Professional.
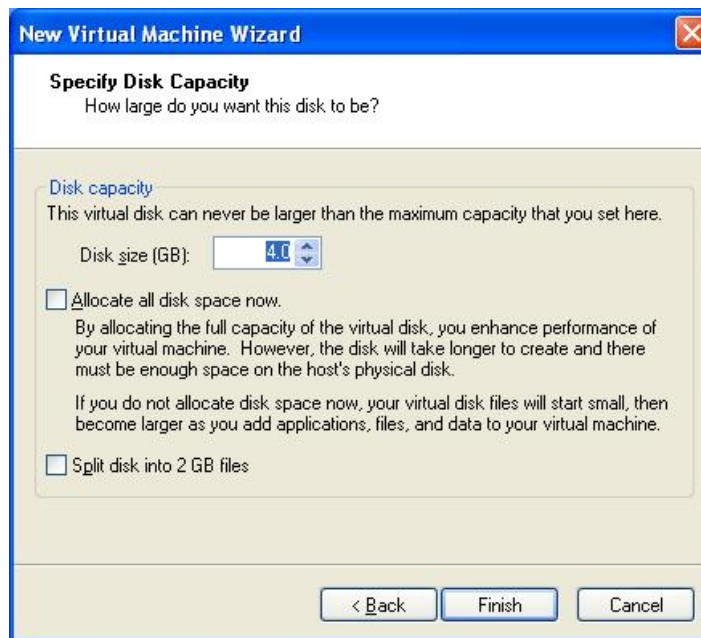
6. You then need to create a unique name for your new Virtual Machine, and
   select the location you would like the associated files to be saved. I generally
   create a new folder within the "My Virtual Machines" folder for each new
   machine I create. This seems to make organization simpler in the long run. It
   is not advisable to store more than one virtual machine in the same folder.
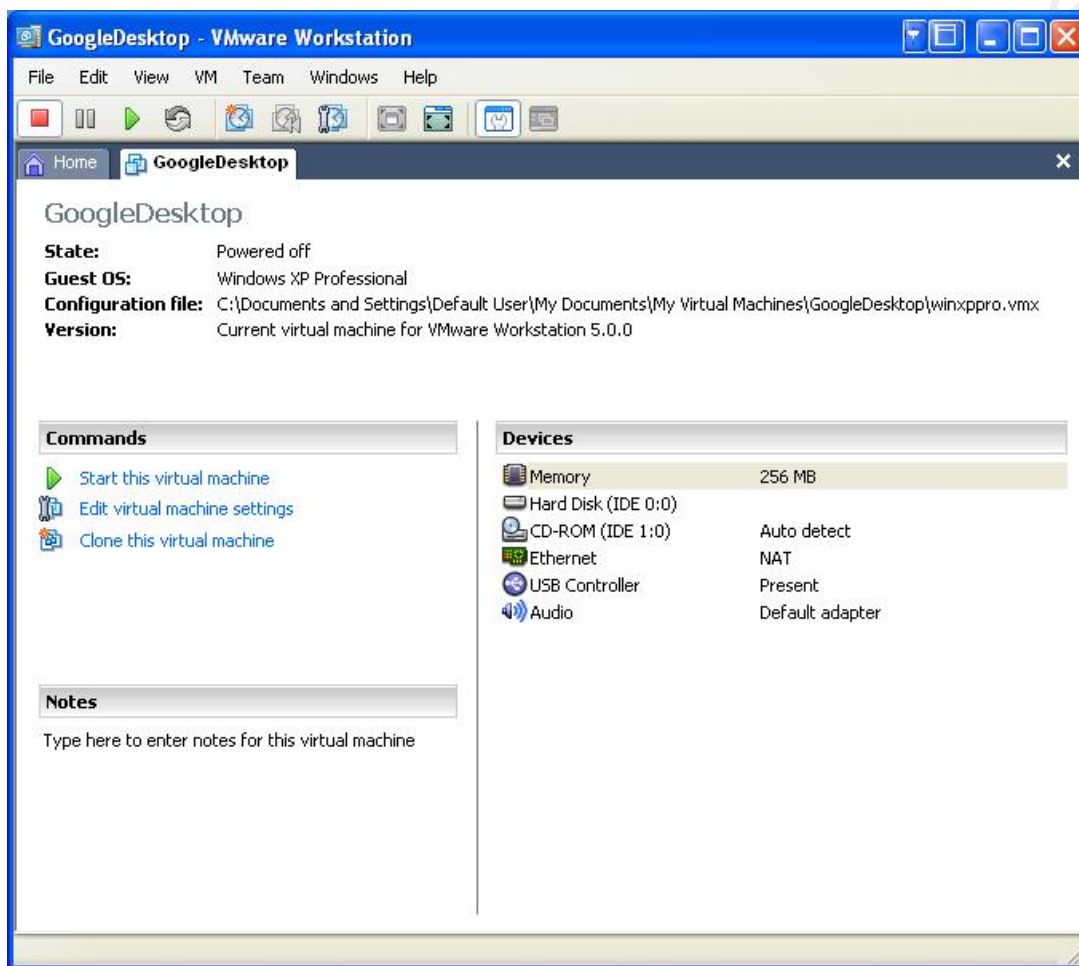
7. When I initially set up the machine, I enable NAT networking so that the virtual machine has access to the Internet. This is so that I can download the appropriate software that I am planning to work with, as well as any operating system updates I need. Since I am just using this machine to view the Google Cache data, I will not worry about OS updates or anything else. I will use the network to download the Google Desktop software and Encase, and then I will remove the network connection after that.



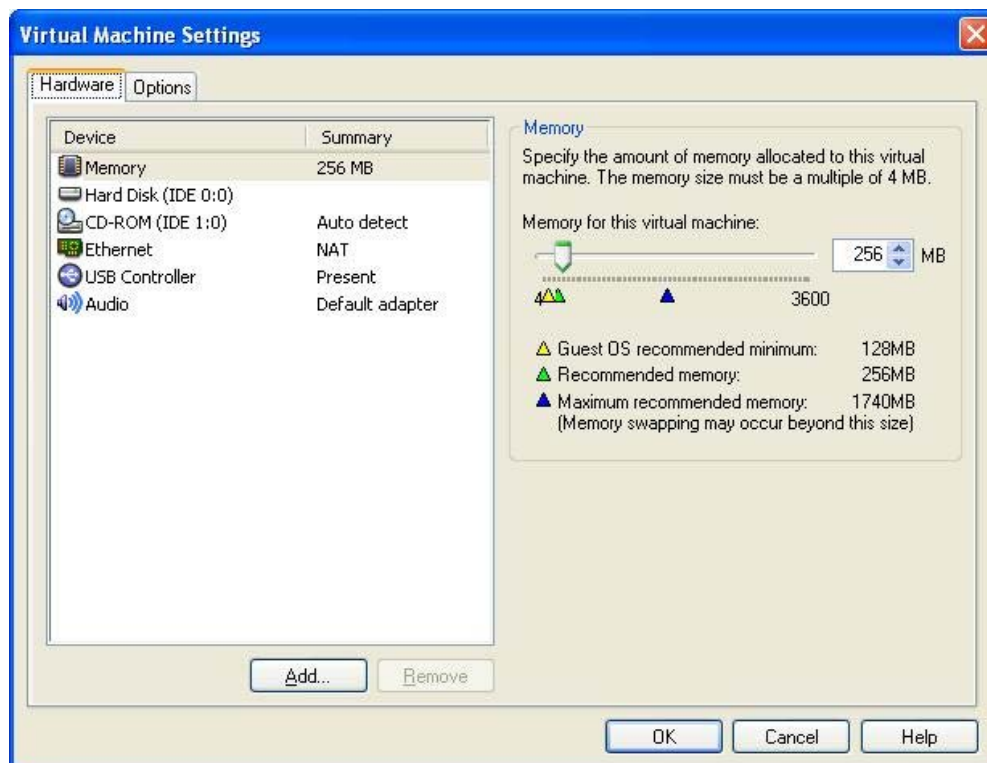8. I leave all the default selections for disk capacity, and click "Finish".
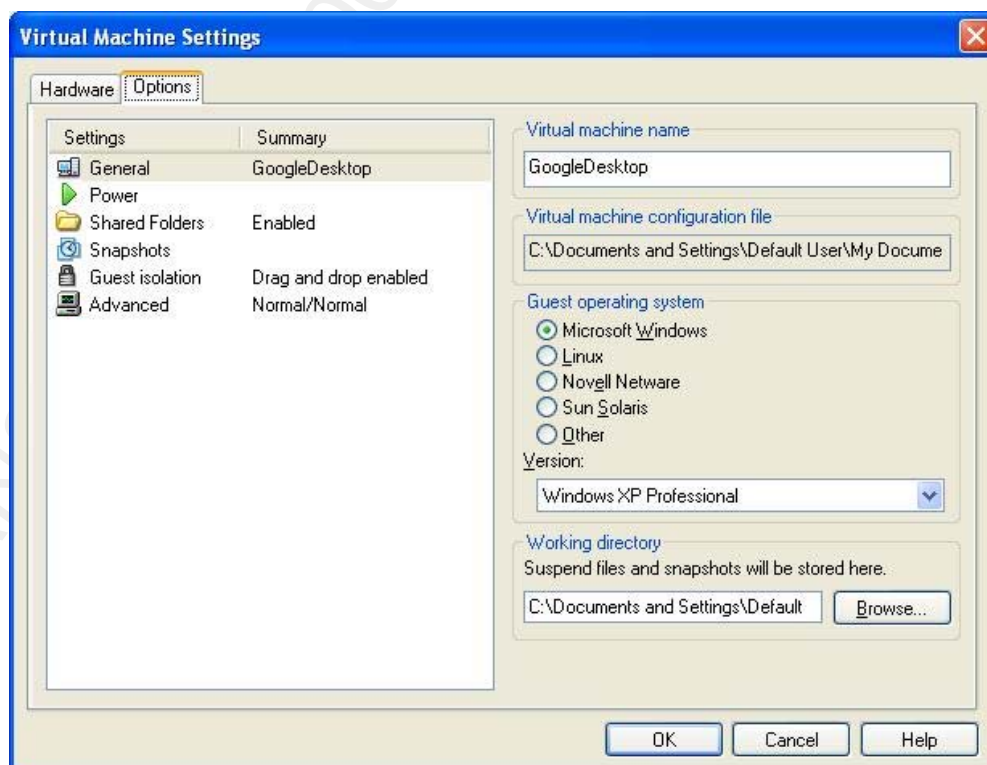
9. At this point, you should have a new Virtual Machine named "GoogleDesktop", or whatever you chose to name it. This is all that is needed to begin installing your guest operating system to the new machine. I will not cover installation of the operating system in this paper since everyone at this point should be able to handle that.

10. Once you have the guest operating system installed, you will want to create a shared folder between your host machine and guest machine so you have access to the suspect image files, as well as a location to copy relevant findings to. This is accomplished by clicking on "Edit virtual machine settings" from the "GoogleDesktop" home screen.
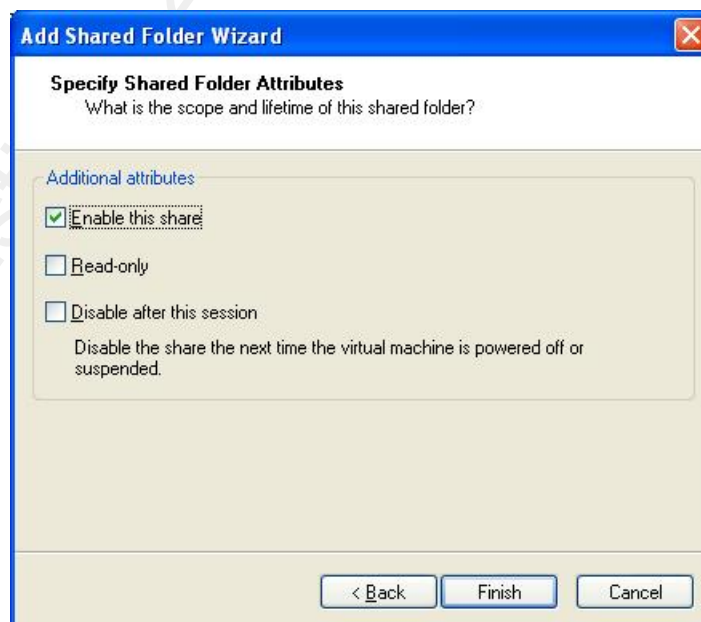
11. Click on the "Options" tab.

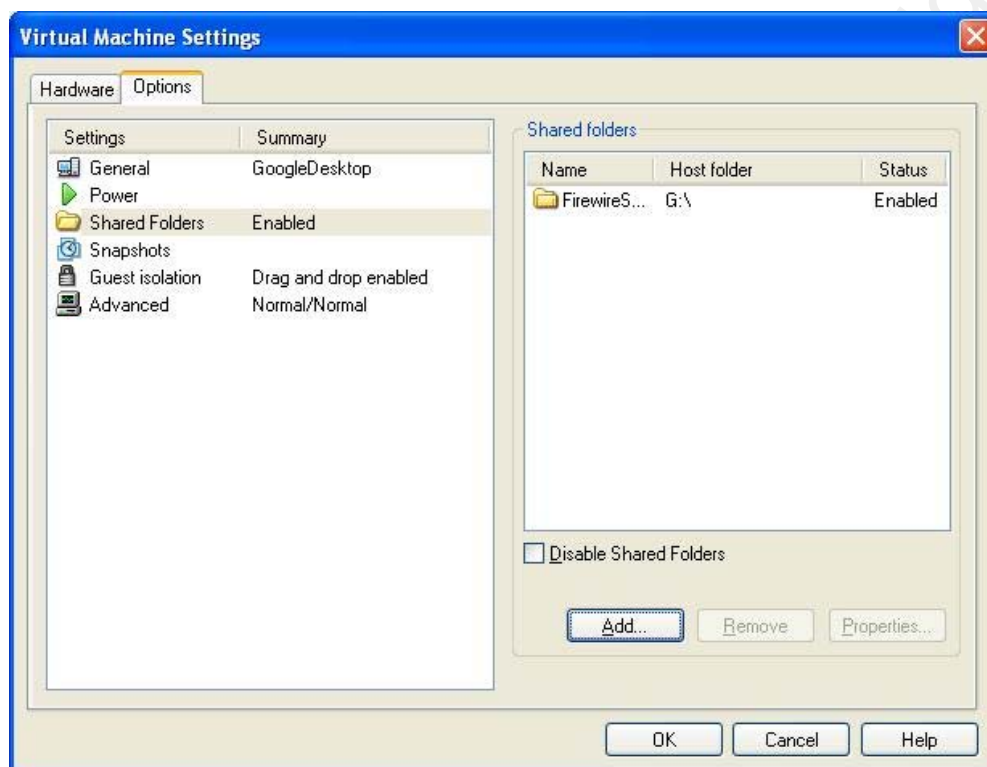12. You will then select "Shared Folders". You will enter the "Add Shared Folder Wizard."



13. I created a shared folder that points to a swappable firewire drive bay I have attached to my host system. This gives me access to whatever drive is in the drive tray, and with it being a firewire connection, the drive is hot swappable as needed.

**Add Shared Folder Wizard**

**Name the Shared Folder**
What would you like to call this shared folder?

Name
FirewireSwap

Host folder
G:\    [Browse...]

[< Back]  [Next >]  [Cancel]

14. I chose to leave the default attributes. If you wanted, you could make your access to the shared folder read only. This is a personal choice. I chose to allow write access so that I can copy relevant findings to the shared drive so all of my case data will remain in the same location. I am not concerned with my image files, since they are individually set to a read only state anyway.

**Add Shared Folder Wizard**

**Specify Shared Folder Attributes**
What is the scope and lifetime of this shared folder?

Additional attributes
☑ Enable this share
☐ Read-only
☐ Disable after this session
Disable the share the next time the virtual machine is powered off or suspended.

[< Back]  [Finish]  [Cancel]

15. You should now see the properties for the new shared folder you just created.



16. Another option to consider is to make the virtual hard drive "Nonpersistent". This setting is found in the advanced setting of the virtual hard disk. This will ensure that each time you boot up the GoogleDesktop virtual machine, it is in the exact same state. There are pros and cons to this setting. A pro is that every change you make to the virtual machine is discarded upon shutdown. Each case will start with a clean slate. This setting, of course, would not be enabled until you have all of the software installed that you need, and you have the work environment set exactly the way you want. Doing this too soon can be a headache! A con to the "Nonpersistent" setting is that each time you go through the trouble of copying a suspect's Google Desktop data to the virtual machine, all is lost when you shutdown the machine. You would have to conduct all of your work in one session, or start over each time. My personal solution to this was to create the initial virtual machine, exactly the way I wanted it, and set it to be Nonpersistent. This gives me a "template"

Google Desktop Search as an Analysis Tool                                                        29
Chris A. Poldervaart

machine. Then when I work a new case, I just clone the virtual machine (VM
☐ Clone from the toolbar) and set it to back to the default dependant state
(uncheck "Independent in the advanced Hard Disk settings). That leaves me
with a working machine specific for that case, and I don't have to worry about
losing changes each time I shut down. I then just simply delete the folder for
the case-specific virtual machine when I am done working with it (mostly for
space).

# References:

Google. Google Desktop Help Center.  2004.
    http://desktop.google.com/support/

Krishnan, Sriram. Reverse Engineering Google Desktop Search. November 22, 2004.
    http://dotnetjunkies.com/WebLog/sriram/archive/2004/11/22/33091.aspx.

# Tools Referenced by this Paper:

AccessData FTK Imager

    http://www.accessdata.com/

Encase Forensic Edition

    http://guidancesoftware.com/

Google Desktop

    http://desktop.google.com/

VMWare Workstation

    http://www.vmware.com/

Google Desktop Search as an Analysis Tool                                    31
Chris A. Poldervaart