

Evil Searching: Compromise and Recompromise of Internet Hosts for Phishing

Tyler Moore^{a)} and Richard Clayton^{b)}

^{a)} Harvard University, Center for Research on Computation and Society, USA
tmoore@seas.harvard.edu

^{b)} Computer Laboratory, University of Cambridge, UK
richard.clayton@cl.cam.ac.uk

Abstract. Attackers compromise web servers in order to host fraudulent content, such as malware and phishing websites. While the techniques used to compromise websites are widely discussed and categorized, analysis of the methods used by attackers to identify targets has remained anecdotal. In this paper, we study the use of search engines to locate potentially vulnerable hosts. We present empirical evidence from the logs of websites used for phishing to demonstrate attackers' widespread use of search terms which seek out susceptible web servers. We establish that at least 18% of website compromises are triggered by these searches. Many websites are repeatedly compromised whenever the root cause of the vulnerability is not addressed. We find that 19% of phishing websites are recompromised within six months, and the rate of recompromise is much higher if they have been identified through web search. By contrast, other public sources of information about phishing websites are not currently raising recompromise rates; we find that phishing websites placed onto a public blacklist are recompromised no more frequently than websites only known within closed communities.

1 Introduction

Criminals use web servers to host phishing websites that impersonate financial institutions, to send out email spam, to distribute malware, and for many other illegal activities. To reduce costs, and to avoid being traced, the criminals often compromise legitimate systems to host their sites. Extra files – web pages or applications – are simply uploaded onto a server, exploiting insecurities in its software. Typical techniques involve the exploitation of flaws in the software of web-based forums, photo galleries, shopping cart systems, and blogs. The security ‘holes’ that are taken advantage of are usually widely known, with corrective patches available, but the website owner has failed to bother to apply them.

The criminals use a number of techniques for finding websites to attack. The most commonly described is the use of scanners – probes from machines controlled by the criminals – that check if a remote site has a particular security vulnerability. Once an insecure machine is located, the criminals upload ‘rootkits’ to ensure that they can recompromise the machine at will [26], and then exploit

the machine for their own purposes – or perhaps sell the access rights on the black market [10]. If the access obtained is insufficient to deploy a rootkit, or the criminal does not have the skills for this, the website may just have a few extra pages added, which is quite sufficient for a phishing attack.

An alternative approach to scanners, that will also locate vulnerable websites, is to ask an Internet search engine to perform carefully crafted searches. This leverages the scanning which the search engine has already performed, a technique that was dubbed ‘Google hacking’ by Long [16]. He was interested not only in how compromisable systems might be located, but also in broader issues such as the discovery of information that was intended to be kept private. Long called the actual searches ‘googledorks’, since many of them rely upon extended features of the Google search language, such as ‘inurl’ or ‘intitle’.

In this paper we examine the evidence for the use of ‘evil searches’: googledorks explicitly intended to locate machines that can be used in phishing attacks.¹ In Section 2 we explain our methodology and give details of our datasets. Although it is widely accepted that criminals use these techniques, to our knowledge, this is the first study to document their prevalence ‘in the wild’.

We make a number of contributions. In Section 3 we clearly establish ‘cause and effect’ between the use of evil searches and the compromise of web servers and estimate the extent of evil searching. In Section 4 we study website *re*-compromise, showing that over 19% of compromised servers host a phishing website on at least one more occasion. In Section 4.3 we demonstrate a clear linkage between evil search and these recompromises. However, ‘findability’ is not necessarily bad; in Section 5 we consider the subset of websites that appear in PhishTank’s publicly available list of compromised sites and find evidence that being listed in PhishTank slightly decreases the rate of recompromise, demonstrating the positive value of this data to defenders. Our final contribution, in Section 6, is to discuss the difficulties in mitigating the damage done by evil searching, and the limitations on using the same searches for doing good.

2 Data collection methodology

We receive a number of disparate ‘feeds’ of phishing website URLs. We take a feed from a major brand owner, which consists almost exclusively of URLs for the very large number of websites attacking their company, and another feed that is collated from numerous sources by the Anti-Phishing Working Group (APWG) [3]. We fetch data from two volunteer organizations: ‘PhishTank’ [21], which specializes in the URLs of phishing websites, and ‘Artists Against 419’ [4], which mainly deals with sites designed to facilitate auction scams or complex advanced fee fraud conspiracies. We also receive feeds from two ‘brand protection’ companies who offer specialist phishing website take-down services. These companies amalgamate feeds from numerous other sources, and combine them with data from proprietary phishing email monitoring systems.

¹ While we focus on websites used for phishing, once a site is found it could be used for any malevolent purpose (e.g., malware hosting).

Type of phishing attack	Count	%
Compromised web servers	88 102	75.8
Free web hosting	20 164	17.4
Rock-phish domains	4 680	4.0
Fast-flux domains	1 672	1.4
‘Ark’ domains	1 575	1.4
Total	116 193	100

Table 1. Categorization of phishing website hosting, October 2007–March 2008.

Although by their nature these feeds have substantial overlaps with each other, in practice each contains a number of URLs that we do not receive from any other source. The result is that we believe that our database of URLs is one of the most comprehensive available, and the overwhelming majority of phishing websites will come to our attention. In principle, we could use capture-recapture analysis to estimate what proportion of sites we were unaware of, as attempted by Weaver and Collins [27]. However, the lack of independence between the various feeds makes a robust estimate of coverage impractical to achieve.

2.1 Phishing-website demographics

In this paper we consider the phishing websites that first appeared in our feeds during the six month period from October 2007 through March 2008. We can split these into a number of different categories according to the hosting method used. Table 1 summarizes their prevalence.

By far the most common way to host a phishing website is to compromise a web server and load the fraudulent HTML into a directory under the attacker’s control. This method accounts for 75.8% of phishing. It is these sites, and the extent to which they can be located by evil searches, that this paper considers.

A simpler, though less popular approach, is to load the phishing web page onto a ‘free’ web host, where anyone can register and upload pages. Approximately 17.4% of phishing web pages are hosted on free web space, but since there is no ‘compromise’ here, merely the signing up for a service, we do not consider these sites any further.

We can also distinguish ‘rock-phish’ and ‘fast-flux’ attacks, where the attackers use malware infected machines as proxies to hide the location of their web servers [19]. A further group, we dub ‘Ark’, appears to use commercial web hosting systems for their sites. All of these attackers use lengthy URLs containing randomly chosen characters. Since the URLs are treated canonically by the use of ‘wildcard’ DNS entries, we ignore the specious variations and just record canonical domain names. Collectively, these three methods of attack comprise 6.8% of phishing websites. Once again, because the exploitation does not involve the compromise of legitimate web servers, and hence no evil searching is required, we do not consider these attacks any further.

Search type	Websites	Phrases	Visits
Any evil search	204	456	1 207
Vulnerability search	126	206	582
Compromise search	56	99	265
Shell search	47	151	360

Table 2. Evil search terms found in Webalizer logs, June 2007–March 2008.

2.2 Website-usage summaries

Many websites make use of The Webalizer [24], a program for summarizing web server log files. It creates reports of how many visitors looked at the website, what times of day they came, the most popular pages on the website, and so forth. It is not uncommon to leave these reports ‘world-readable’ in a standard location on the server, which means that anyone can inspect their contents.

From June 2007 through March 2008, we made a daily check for Webalizer reports on each website appearing in our phishing URL feeds. We recorded the available data – which usually covered activity up to and including the previous day. We continued to collect the reports on a daily basis thereafter, allowing us to build up a picture of the usage of sites that had been compromised and used for hosting phishing websites.

In particular, one of the individual sub-reports that Webalizer creates is a list of search terms that have been used to locate the site. It can learn these if a visitor has visited a search engine, typed in particular search terms and then clicked on one of the search results. The first request made to the site that has been searched for will contain a ‘Referrer’ header in the HTTP request, and this will contain the terms that were originally searched for.

2.3 Types of evil search

In total, over our ten month study, we obtained web usage logs from 2 486 unique websites where phishing pages had been hosted (2.8% of all compromised websites). Of these usage logs, 1 320 (53%) recorded one or more search terms.

We have split these search terms into groups, using a manual process to determine the reason that the search had been made. Many of the search terms were entirely innocuous and referred to the legitimate content of the site. We also found that many advanced searches were attempts to locate MP3 audio files or pornography – we took no further interest in these searches.

However, 204 of the 1 320 websites had been located one or more times using ‘evil’ search terms, viz: the searches had no obvious innocent purpose, but were attempts to find machines that might be compromised for some sort of criminal activity. We distinguish three distinct types of evil search and summarize their prevalence in Table 2.

Vulnerability searches are intended to pick out a particular program, or version of a program, which the attacker can subvert. Examples of searches in this

group include ‘phpizabi v0.848b c1 hfp1’ (CVE-2008-0805 is an unrestricted file upload vulnerability) and ‘inurl: com_juser’ (CVE-2007-6038 concerns the ability of remote attackers to execute arbitrary PHP code on a server).

Compromise searches are intended to locate existing phishing websites, perhaps particular phishing ‘kits’ with known weaknesses, or just sites that someone else is able to compromise. Examples include ‘allintitle: welcome paypal’ and ‘inurl:www.paypal.com’ which both locate PayPal phishing sites.

Shell searches are intended to locate PHP ‘shells’. When attackers compromise a machine they often upload a PHP file that permits them to perform further uploads, or to search the machine for credentials – the file is termed a shell since it permits access to the underlying command interpreter (`bash`, `csh` etc.). The shell is often placed in directories where it becomes visible to search engine crawlers, so we see searches such as ‘intitle: "index of" r57.php’ which looks for a directory listing that includes the `r57` shell, or ‘c99shell drwxrwx’ which looks for a `c99` shell that the search engine has caused to run, resulting in the current directory being indexed – the `drwxrwx` string being present when directories have global access permissions.

3 Evidence for evil searching

So far, we have observed that some phishing websites are located by the use of dubious search terms. We now provide evidence of evil searches leading directly to website compromise. While difficult to attain absolute certainty, we can show that there is a consistent pattern of the evil searches appearing in the web logs at or before the time of reported compromise.

3.1 Linking evil search to website compromise

Figure 1 presents an example timeline of compromises, as reconstructed from our collections of phishing URLs and Webalizer logs. On 30 November 2007, a phishing page was reported on the `http://chat2me247.com` website with the path `/stat/q-mono/pro/www.lloydstsb.co.uk/lloyds_tsb/logon.ibc.html`.

We began collecting daily reports of `chat2me247.com`’s Webalizer logs. Initially, no evil search terms were recorded, but two days later, the website received a visit triggered by the search string ‘phpizabi v0.415b r3’. Less than 48 hours after that, *another* phishing page was reported, with the quite different location of `/seasalter/www.usbank.com/online_banking/index.html`.

Given the short period between search and re-compromise, it is very likely that the second compromise was triggered by the search. Also, the use of a completely different part of the directory tree suggests that the second attacker was unaware of the first. Figure 1 shows a screenshot from a web search in April 2008 using the same evil search term: `chat2me247.com` is the 13th result out of 696 returned by Google, suggesting a high position on any attacker’s target list.

We have observed similar patterns on a number of other websites where evil search terms have been used. In 25 cases where the website is compromised



1:	2007-11-30 10:31:33	phishing URL reported: http://chat2me247.com/stat/q-mono/pro/www.lloydstsb.co.uk/lloyds_tsb/logon.ibc.html	
2:	2007-11-30	no evil search term	0 hits
3:	2007-12-01	no evil search term	0 hits
4:	2007-12-02	phpizabi v0.415b r3	1 hit
5:	2007-12-03	phpizabi v0.415b r3	1 hit
6:	2007-12-04 21:14:06	phishing URL reported: http://chat2me247.com/seasalter/www.usbank.com/online_banking/index.html	
7:	2007-12-04	phpizabi v0.415b r3	1 hit

Fig. 1. Screenshot and timeline of a phishing website compromise using an evil search.

multiple times (as with `chat2me247.com`), we have fetched Webalizer logs in the days immediately preceding the recompromise (because we were studying the effects of the initial compromise). For these sites we are able to ascertain whether the evil search term appears before compromise, on the same day as the compromise, or sometime after the compromise.

Figure 2 (top) shows a timeline for the 25 websites with Webalizer data before and after a second compromise. For 4 of these websites, the evil search term appeared before the recompromise. For the vast majority (20), the evil search term appeared on the day of the recompromise. In only one case did the evil search term appear only after recompromise. Since most evil terms appear at or before the time of recompromise, this strongly suggests that evil searching is triggering the second compromise. If the evil searches had only occurred after website compromise, then there would have been no connection.

We also examined the Webalizer logs for an additional 177 websites with evil search terms but where the logs only started on, or after, the day of the com-

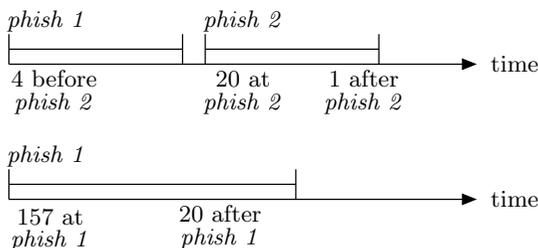


Fig. 2. Timeline of evil web search terms appearing in Webalizer logs.

promise (see Figure 2 (bottom)). Again, in most cases (157) the evil search term appeared from the time of compromise. Taken together, evil search terms were used at or before website compromise 90% of the time. This is further evidence that evil searching is a precursor to the compromise of many web servers.

3.2 Estimating the extent of evil search

We can use the incidence of phishing websites that have Webalizer logs as a sample to estimate the overall prevalence of evil search when servers are compromised and used to host phishing websites.

Recall that we have obtained search logs for 1320 phishing websites, and that 204 of these websites include one or more evil search terms in these logs. Frequently, the record shows one visit per evil search.

Unfortunately, Webalizer only keeps a record of the top 20 referring search terms. Hence, if a site receives many visitors, any rarely occurring search term will fall outside the top 20. We therefore restrict ourselves to considering just the 1085 of Webalizer-equipped hosts that have low enough traffic so that even search terms with one visit are recorded. Of these hosts, 189 include evil search terms, or approximately 17.6% of the hosts in the sample. Viewed as a sample of all compromised phishing websites, the 95% confidence interval for the true rate of evil searching is (15.3%, 19.8%).

This estimate is only valid if the hosts with Webalizer logs represent a truly random sample. A number of factors may affect its suitability:

- Running Webalizer (or programs that it may be bundled with) may affect the likelihood of compromise. We have no evidence for any such effect.
- Sites running Webalizer are not representative of the web server population as a whole. Webalizer typically runs on Unix-like operating systems. Since many compromised servers run on Windows hosts, we cannot directly translate the prevalence of evil web search terms to these other types.
- Evil searches are only recorded in the website logs if the attacker clicks on a search result to visit the site. Using automated tools such as Goolag [6], or simple cut & paste operations, hides the search terms. This leads us to underestimate the frequency of evil searches.

On balance, we feel sites with Webalizer logs are a fair sample of all websites.

3.3 Other evidence for evil searches

There is a substantial amount of circumstantial evidence for the use of evil searches by criminals seeking machines to compromise. Hacker forums regularly contain articles giving ‘googledorks’, sometimes with further details of how to compromise any sites that are located. However, published evidence of the extent to which this approach has replaced older methods of scanning is hard to find, although the topic is already on the curriculum at one university [15].

LaCour examined a quarter of the URLs in the MarkMonitor phishing URL feed, and was reported [13] as finding that, “75% had been created by using some 750 evil search terms, and the related PHP vulnerabilities”. Unfortunately, he was misquoted [14]. LaCour did collect 750 evil searches from hacker forums, but he did not establish the extent to which these were connected to actual machine compromises, whether for phishing or any other purpose.

What LaCour was able to establish from his URL data was that for the October to December 2007 period, 75% of attacks involved machine compromise, 5% were located on free web-hosting and 20% were the categories we have called rock-phish, fast-flux and Ark. These figures are roughly in line with our results in Table 1 above. He then observed, from the paths within the URLs, a strong link with PHP vulnerabilities, particularly ‘Remote File Inclusion’ (RFI) [8]. This is what led him to speculate that evil searches and subsequent RFI attacks are a key element in the creation of 75% of all phishing websites.

4 Phishing website recompromise

Removing phishing websites can be a frustrating task for the banks and other organizations involved in defending against phishing attacks. Not only do new phishing pages appear as fast as old ones are cleared, but the new sites often appear on the web servers that were previously compromised and cleaned up. This occurs whenever the sysadmin removing the offending content only treats the symptoms, without addressing the root problem that enabled the system to be compromised in the first place.

We now provide the first robust data on the *rate* of phishing-website recompromise. We show how the recompromise rate varies over time, and then provide evidence of how evil search raises the likelihood of recompromise.

4.1 Identifying when a website is recompromised

Websites may be recompromised because the same attacker returns to a machine that they know to be vulnerable. Alternatively, the recompromise may occur because a different attacker finds the machine and independently exploits it using the same vulnerability, or even a second security flaw. We think it unlikely that a single attacker would use multiple security flaws to compromise a machine when just one will do the trick.

The general nature of the security flaw that has been exploited is often quite obvious because the phishing pages have been added within particular parts of

the directory structure. For example, when a particular user account is compromised the phishing pages are placed within their filespace; when a file upload vulnerability is exploited, the pages are put in sub-directories of the upload repository. However, since it is not always possible to guess what exploit has been used, we instead consider how much time elapses between phishing reports to infer distinct compromises.

If two phishing websites are detected on the same server within a day of each other, it is more likely that the same attacker is involved. If, instead, the attacks are months apart, then we believe that is far more likely that the website has been rediscovered by a different attacker. We believe that attackers usually have a relatively small number of machines to exploit at any given moment and are unlikely to keep compromised machines ‘for a rainy day’ – this is consistent with the short delay that we generally see between detection (evil search logged) and use (phishing website report received).

Our equating of long delays with different attackers is also based on the distribution of recompromises over time. If we treat every phishing site on a particular server as a different attack, whatever the time delay, then we observe a recompromise rate of 20% after 5 weeks, rising to 30% after 24 weeks. If we insist that there is a delay of at least 3 weeks between attacks to consider the event to be a recompromise, then the rates change to 2% after 5 weeks and 15% after 24 weeks. The long term rates of recompromise vary substantially for cut-off points of small numbers of days, which we believe reflects the same attackers coming back to the machine. However, the long term rates of recompromise hardly change for cut-off times measured in weeks, which is consistent with all recompromises being new attackers.

An appropriate cut-off point, where there is only a small variation in the results from choosing slightly different values, is to use a gap of one week. We therefore classify a phishing host as recompromised after receiving two reports for the same website that are at least 7 days apart. Using a 7-day window strikes a reasonable balance between ensuring that the compromises are independent without excluding too many potential recompromises from the calculations.

As a further sanity check, we note that for 83% of website recompromises occurring after a week or longer, the phishing page is placed in a different directory than previously used. This strongly suggests that different exploits are being applied, and therefore, different attackers are involved.

4.2 Measuring website recompromise rates

The rate of website recompromise should only be considered as a function of time. Simply computing the recompromise rate for all phishing websites in the October to March sample would skew the results: websites first compromised on October 1st would have six months to be recompromised, while websites first compromised in late March would have far less time. For this reason, we consider website recompromise in four-week intervals. For instance, we test whether a website compromised on October 1st has been recompromised by October 29th,

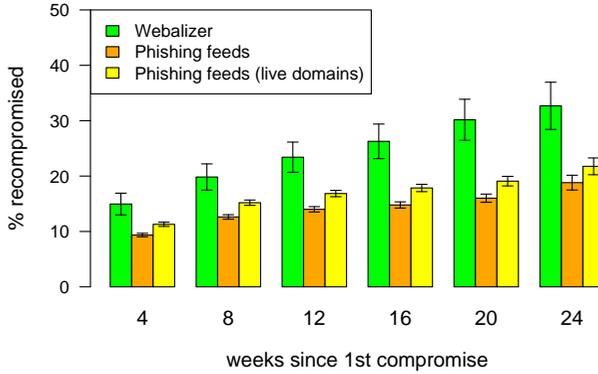


Fig. 3. Recompromise rates for phishing websites over time. The error bars show the 95% confidence interval for the true value of the recompromise rate.

November 26th, and so on. Similarly, we can only check whether a website compromised on March 1st 2008 has been recompromised by March 29th.

Figure 3 plots phishing website recompromise over time. The graph includes recompromise rates for the 1 320 Webalizer hosts, along with the 36 514 other hosts we recorded between October 2007 and March 2008. In both cases, the recompromise rate increases over time: 15% of hosts with Webalizer logs are recompromised within 4 weeks, rising steadily to 33% within 24 weeks. The recompromise rate for the other hosts is somewhat lower, but follows the same pattern: 9% are recompromised within 4 weeks, rising to 19% within 24 weeks.

What might explain the discrepancy in the recompromise rates for the Webalizer sample? One factor is that the sites with Webalizer logs, by definition, were accessible at least once shortly after being reported. This is not the case for all hosts – some phishing websites are completely removed before we are able to access them.²

Sites that quickly disappear are far less likely to be recompromised in the future. Hence, Figure 3 also plots the recompromise rates for the 29 986 websites that responded at least once. The recompromise rate for these websites is slightly higher than that for all phishing websites.

In any event, the results from this graph offer a strong indication that phishing website recompromise happens frequently. Many website administrators are not taking adequate steps to prevent recompromise following an intrusion.

² Many sites that are compromised are long-abandoned blogs and image galleries. It is not surprising that a number of these are removed altogether, rather than being cleaned up and left publicly available.

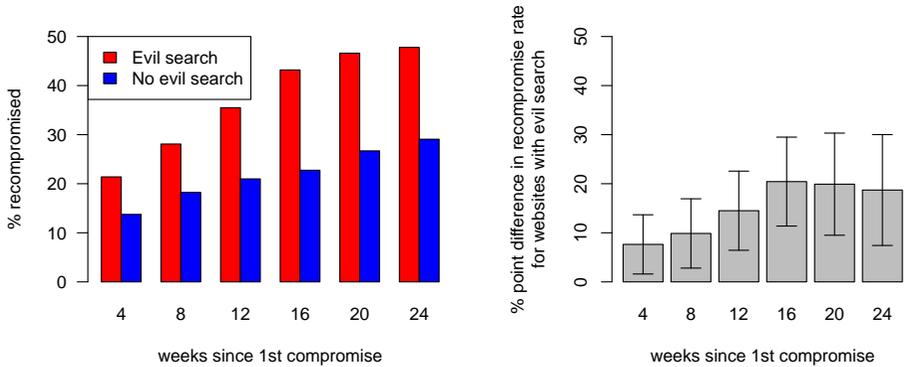


Fig. 4. Recompromise rates for phishing websites with and without evil search responses found in the Webalizer logs (left). The right graph shows the percentage point difference along with 95% confidence intervals.

4.3 Evil searching and recompromise

Section 3.1 established that evil searches can precede website compromise. We now show that the evil searches are linked to much higher rates of recompromise.

Figure 4 (left) compares the recompromise rates for hosts in the Webalizer sample. Sites with evil search terms in the logs are far more likely to be recompromised than sites without such terms. Hosts reached by evil search face a 21% chance of recompromise after 4 weeks, compared to 14% otherwise. Within 24 weeks these numbers rise to 48% and 29% respectively.

Moreover, these differences are statistically significant. Figure 4 (right) plots the percentage point difference between recompromise rates when evil and non-evil searches are present, along with 95% confidence intervals. For instance, there is a 20.4 percentage point difference in recompromise rates after 16 weeks (43.2% recompromise for evil searches compared to 22.8% for the rest). The evil search recompromise rate is nearly twice that of ordinary phishing websites for the period. What does this mean? Vulnerable websites that can be found through web search are likely to be repeatedly rediscovered and recompromised until they are finally cleaned up.

5 PhishTank and recompromise

We have shown that attackers use web search to find websites to compromise. We now consider whether they are using public phishing website blacklists as an alternative way to find sites to compromise. These blacklists provide valuable data for ‘phishing toolbars’ that block visits to fraudulent websites. Most blacklists are kept hidden: Google’s SafeBrowsing API [12] only allows users to verify suspected URLs, while the APWG’s blacklist [3] is only available to members.

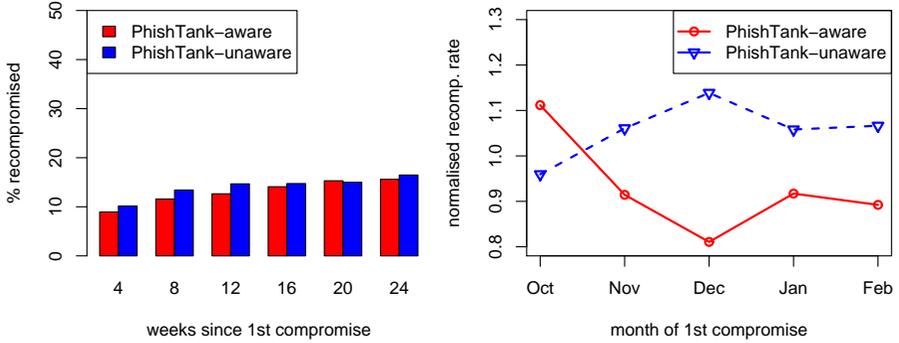


Fig. 5. Recompromise rate for phishing websites appearing on the public website PhishTank (left); normalized 4-week recompromise rates based upon the month of first compromise (right).

In contrast, ‘PhishTank’ [21] provides an open source blacklist which is generated and maintained through web-based participation. Users are invited to submit URLs of suspected phishing websites and verify each other’s entries. Consequently, PhishTank provides a permanent record of phishing websites dating back to its inception in late 2006. They also publish a more dynamic list of recently active phishing websites. We now test whether appearing in PhishTank’s public blacklist makes website recompromise more likely.

It is unfair to simply compare recompromise rates for sites PhishTank knows about with those of which it is unaware. While aiming to be comprehensive, in practice PhishTank fails in this aim, and is aware of only 48% of the phishing websites in our collection. Since some of our other URL feeds get some of their data from PhishTank, it is more accurate to view PhishTank as a subset of the phishing URLs we record. So although PhishTank has a roughly even chance of recording a particular phishing incident, there will be further chances to record the host if it is recompromised. This biases PhishTank’s record to include a disproportionate number of hosts where multiple compromises occur.

Consequently, we apply a fairer test to determine whether a host’s appearance in PhishTank makes it more likely to be recompromised. We compare the recompromise rates of new hosts following their first compromise. 9 283 hosts detected by PhishTank during their first reported compromise are compared against 15 398 hosts missed by PhishTank during the first compromise. Since we are only considering URLs reported from October 2007 to March 2008, we ignore URLs first appearing in PhishTank prior to October 2007.

The results are presented in Figure 5 (left) and show that new websites appearing in PhishTank are no more likely to be recompromised than new websites that do not appear. Website recompromise over the short term (up to 12 weeks) is less for websites publicized in PhishTank compared to those hidden from it.

Within 4 weeks, PhishTank-aware phishing websites are recompromised 8.9% of the time, compared to 10.2% for sites not reported to PhishTank. A similar trend holds for recompromised websites within 8 and 12 weeks, with recompromise rates around two percentage points lower for websites known to PhishTank. These differences are maintained with 95% confidence. However, over the longer term (16 to 24 weeks), the recompromise rates become indistinguishable.

Why might sites appearing in PhishTank be recompromised less often? It appears that defenders are paying more attention to PhishTank's lists than attackers are. By making its blacklist available free of charge, more support staff at ISPs and sysadmins are informed of compromised websites in need of cleanup. Other companies sell phishing feeds to aid ISPs in this manner, but PhishTank's free service may be more widely adopted. As more defenders become aware of PhishTank (and consequently aware of more phishing websites), we might expect PhishTank's recompromise rate to diminish further over time. To test this hypothesis in our data, Figure 5 (right) plots recompromise rates after 4 weeks for phishing websites based on the month the site is reported. The data is normalized with respect to the overall phishing activity in the relevant month. In October, the recompromise rate for websites reported to PhishTank is higher than in the set of websites of which PhishTank is unaware. However, this situation turns around thereafter, with the recompromise rates for sites in PhishTank reducing and becoming *lower* than the rising recompromise rate for the sites which PhishTank missed.³

Based on our data analysis, we conclude that the good offered by PhishTank (better information for defenders) currently outweighs the bad (exploitation of compromised websites by attackers). However, the use of PhishTank by both attackers and defenders might change dynamically over time. Consequently, we believe that continued monitoring is necessary in case attackers begin to leverage PhishTank's public blacklist.

6 Mitigation strategies

Thus far we have demonstrated clear evidence that evil searches are actively used to locate web servers for hosting phishing websites. We have also shown that server re-compromise is often triggered by evil search. Therefore, we now consider how evil searches might be thwarted, in order to make the criminals' task harder. We set out and review a number of mitigation strategies, the first two of which can be implemented locally, whereas the others require action by outside parties. Unfortunately each has drawbacks.

Strategy 1: Obfuscating targeted details Evil searches could be made less effective if identifying information such as version numbers were removed from web server

³ In October 2007 PhishTank added highly-publicized features to its website, which permit searches for phishing sites based on ASNs, and RSS feeds of new entries within an ASN; exactly meeting the requirements of an ISP that wished to keep track of any compromised customers.

applications. While this might make it a bit harder for attackers to discover vulnerable websites, it does nothing to secure them.

Damron [7] argued for obfuscation by noting that removing the version numbers from applications is easy for the defender, while adding a significant burden for the attacker. However, defenders also stand to gain from detailed application information, as the presence of a version number can assist sysadmins in keeping track of which of their users continues to run out of date software.

We note that very few of the evil search terms we examined contained explicit version numbers, but merely sought to identify particular programs. The final objection to this strategy is that obscuring version numbers still leaves users exposed to ‘shotgun’ attackers who run all of their exploits against every candidate site without worrying whether or not it is running a vulnerable version.

Strategy 2: Evil search penetration testing Motivated defenders could run evil searches to locate sites that might be compromised and then warn their owners of the risk they were running. For many evil searches, which only return a handful of exploitable sites amongst many thousands of results, this is unlikely to be an effective scheme. Furthermore, the search results are usually just hints that only indicate the potential for compromise. Confirming suspicions normally requires an active attack, which would be illegal in most jurisdictions.

Strategy 3: Blocking evil search queries An alternative approach is for the search engines to detect evil searches and suppress the results, or only provide links to law enforcement sites. Given their inherent specificity, constructing a comprehensive and up-to-date blacklist of evil searches is likely to be difficult and costly. Blocking some of the more obvious terms (e.g., those found in Long’s popular database [11]) is unlikely to be effective if the terms used by the criminals rapidly evolve. In any event, the search engines are unlikely to have any meaningful incentive to develop and deploy such a list.

Strategy 4: Removing known phishing sites from search results The low-cost option of removing currently active phishing sites from the search results has almost nothing to recommend it. Search engines suppress results for known child-pornography sites, and Google prevents users from clicking through to sites that are hosting malware [9] until they are cleaned up [17]. However, phishing presents different circumstances. Malware is placed on high traffic sites where removal from search results is a powerful incentive towards getting it removed, but phishing sites are often on semi-abandoned low traffic sites where the incentive to remove will be limited. Although the evil search will not work while the phishing site is active, the site will be findable again as soon as the fraudulent pages are removed. This approach would also prevent any use of searches by defenders, which means that it does some harm as well as doing little good.

Strategy 5: Lower the reputation of previously phished hosts discoverable by evil search terms In addition to flagging active phishing URLs, website reputation

services such as SiteAdvisor [18] already give a warning for websites that consistently host malicious content. Since we have shown that a substantial proportion of systems that host a phishing website are later recompromised, such services might mark previously compromised hosts as risky. Furthermore, it would be entirely prudent to proactively flag as a much higher risk any hosts used for phishing which can also be found by evil search terms. The magnitude of the risk should reflect our finding that about half of these sites will be recompromised within 24 weeks.

7 Related work

As indicated earlier, very little academic research has examined the use of search engines to compromise websites. However, researchers have recently begun to recognize the importance of empirically studying electronic crime. Thomas and Martin [25] and Franklin *et al.* [10] have characterized the underground economy by monitoring the advertisements of criminals on IRC chatrooms. Provos *et al.* [22] tracked malicious URLs advertising malware, finding that 1.3% of incoming Google search queries returned links to malware-distributing URLs. Moore and Clayton [19] studied the effectiveness of phishing-website removal by recording site lifetimes. Collins *et al.* [5] used NetFlow data on scanning, spamming and botnet activity to classify unsafe IP address ranges. The current work contributes to this literature by measuring the prevalence of evil search terms for compromising websites and the impact on site recompromise.

Another related area of literature is the economics of information security [2]. One key economic challenge identified by this literature is overcoming asymmetric information. Better measurement of security is needed, from the prevalence of vulnerabilities in competing software to the responsiveness of ISPs in cleaning up infected hosts. Publishing accurate data on website recompromise can identify serial underperformers and highlight opportunities for improvement. Google and StopBadware [23] publicly disclose infected websites, and it has been claimed that this disclosure encourages prompt cleanup [9]. At a policy level, Anderson *et al.* [1] have recommended that regulators collect better data on system compromise and use it to punish unresponsive ISPs.

8 Conclusion

In this paper, we have presented clear evidence that the criminals who are compromising web servers to host phishing websites are using Internet search engines to locate vulnerable machines. We have found direct evidence of these ‘evil searches’ in 18% of our collection of Webalizer logs from phishing sites, and believe the true prevalence to be even higher.

We have also shown a clear linkage with the recompromise of servers. The general population of phishing websites exhibits a recompromise rate of 19% after 24 weeks, but where evil searches are found in the logs, the rate reaches 48%. Although the use of evil searches has been known about anecdotally, this

is the first paper to show how prevalent the technique has become, and to report upon the substantial rates of recompromise that currently occur.

In contrast, phishing website URLs that are made public by the PhishTank database currently enjoy a slight, but statistically significant, reduction in their recompromise rates. This suggests that defenders are able to use the database in order to reduce criminal attacks, and that the sometimes touted benefits of keeping attack data hidden from public view may be minimal.

Other strategies for mitigating evil search that work by limiting attackers' access to information – obfuscating version numbers, filtering search results, blocking evil search queries – we also consider to be flawed. The most promising countermeasure we discuss is to incorporate a website's likelihood of recompromise into the calculation of its reputation.

References

1. R. Anderson, R. Böhme, R. Clayton, and T. Moore. *Security economics and the internal market*. European Network and Information Security Agency (ENISA), 2008.
http://enisa.europa.eu/doc/pdf/report_sec_econ_&_int_mark_20080131.pdf
2. R. Anderson and T. Moore. The economics of information security. *Science*, 314(5799):610–613, 2006.
3. Anti-Phishing Working Group. <http://www.apwg.org/>
4. Artists Against 419. <http://www.aa419.org/>
5. M. P. Collins, T. J. Shimeall, S. Faber, J. Janies, R. Weaver, M. De Shon and J. Kadane. Using uncleanliness to predict future botnet addresses. In *Proceedings of the ACM SIGCOMM Conference on Internet Measurement (IMC)*, pp. 93–104, ACM Press, New York, 2007.
6. Cult of the Dead Cow. *Goolag Scanner Specifications*. Jan 2008.
<http://goolag.org/specifications.html>
7. J. Damron. Identifiable fingerprints in network applications. *USENIX ;login*, 28(6):16–20, Dec 2003.
8. M. Dausin. *PHP File Include Attacks*. Tipping Point, Feb 2008.
<http://dvlabs.tippingpoint.com/blog/2008/02>
9. O. Day, B. Palmén and R. Greenstadt. Reinterpreting the disclosure debate for web infections. In *7th Workshop on the Economics of Information Security (WEIS)*, 2008.
10. J. Franklin, V. Paxson, A. Perrig and S. Savage. An inquiry into the nature and causes of the wealth of Internet miscreants. In *Proceedings of the 14th ACM Conference on Computer and Communications Security (CCS)*, pp. 375–388, 2007.
11. Google Hacking Database. <http://johnny.ihackstuff.com/ghdb.php>
12. Google Safe Browsing API. <http://code.google.com/apis/safebrowsing/>
13. K.J. Higgins. *Phishers Enlist Google 'Dorks'*. DarkReading, Mar 2008.
http://www.darkreading.com/document.asp?doc_id=149324
14. J. LaCour. Personal communication, 28 Mar 2008.
15. L. Lancor and R. Workman. Using Google hacking to enhance defense strategies. In *Proceedings of the 38th SIGCSE Technical Symposium on Computer Science Education*, pp. 491–495, 2007.

16. J. Long. *Google Hacking Mini-Guide*. informIT, May 2004.
<http://www.informit.com/articles/article.aspx?p=170880>
17. P. Mavrommatis. *Malware Reviews via Webmaster Tools*. Aug 2007.
<http://googlewebmastercentral.blogspot.com/2007/08/malware-reviews-via-webmaster-tools.html>
18. McAfee Inc. *SiteAdvisor* <http://www.siteadvisor.com>
19. T. Moore and R. Clayton. Examining the impact of website take-down on phishing. In *Anti-Phishing Working Group eCrime Researcher's Summit (APWG eCrime)*, pp. 1–13, ACM Press, New York, 2007.
20. Netcraft Inc. *March 2008 Web Server Survey*., 2008.
http://news.netcraft.com/archives/web_server_survey.html
21. PhishTank. <http://www.phishtank.com/>
22. N. Provos, P. Mavrommatis, M. Rajab and F. Monrose. All your iFrames point to us. In *17th USENIX Security Symposium*, pp. 1–15, 2008.
23. Stop Badware. <http://www.stopbadware.org/>
24. The Webalizer. <http://www.mrunix.net/webalizer/>
25. R. Thomas and J. Martin. The underground economy: priceless. *USENIX ;login*, 31(6):7–16, Dec. 2006.
26. D. Watson, T. Holz and S. Mueller. *Know your Enemy: Phishing*. The HoneyNet Project & Research Alliance, May 2005.
<http://www.honeynet.org/papers/phishing/>
27. R Weaver and M.P. Collins. Fishing for phishes: applying capture-recapture methods to estimate phishing populations. In *Anti-Phishing Working Group eCrime Researcher's Summit (APWG eCrime)*, pp. 14–25. ACM Press, New York, 2007.
28. Yahoo! Inc. *Yahoo! Search Web Services*. <http://developer.yahoo.com/search/>