

Empirically Characterizing Domain Abuse and the Revenue Impact of Blacklisting

Neha Chachra* Damon McCoy[◇] Stefan Savage* Geoffrey M. Voelker*

*Department of Computer Science and Engineering
University of California, San Diego

[◇]Computer Science Department
George Mason University

ABSTRACT

Using ground truth sales data for over 40K unlicensed prescription pharmaceuticals sites, we present an economic analysis of two aspects of domain abuse in the online counterfeit drug market. First, we characterize the nature of domains abused by affiliate spammers to monetize what is evidently an overwhelming demand for these drugs. We found that the most successful affiliates are agile in adapting to adversarial circumstances, and channel the full spectrum of domain abuse to advertise to customers. Second, we use contemporaneous blacklisting data to provide an economic analysis of the revenue impact of *domain blacklisting*, a technique whereby lists of “known bad” registered domains are distributed and used to filter email spam. We found that blacklisting rapidly and effectively limited per-domain sales. Nevertheless, blacklisted domains continued to monetize, likely as a result of high demand, non-universal use of blacklisting, and delay in deployment. Finally, our results suggest that increasing the number of domains discovered and using blacklists to block access to spam domains could undermine profitability more than further improving the speed with which domains are added to blacklists.

1. INTRODUCTION

Virtually every mode of mass communication in use online today — email, search, blogs, social networks, instant messaging, and VoIP — engenders some form of spam that is used to sell for products or services. In nearly all cases, this activity is monetized by driving users to click on Web links for spammer-affiliated e-commerce sites which then process conversions using standard online payment mechanisms (e.g., Visa and Mastercard).

In response, a wide range of defenses have been proposed and implemented to identify such unwanted communication and filter it out of the user’s view (typically either preventing the advertisement containing the link from appearing or preventing the link from being visited). Among the oldest and most widely used of these defenses is *domain blacklisting* — the active identification and distribution of domain names advertised in an unwanted manner. This approach is used today in a broad range of defenses, including email classification, anti-phishing toolbars and search classification, and in turn has driven spammers to a variety of countermeasures (e.g., churning through large numbers of domains, using one or more layers of domain redirection, abusing existing sites to host content or redirect traffic, etc.) However, while the technical components of this domain name arms race are widely understood, the underlying economic issues that drive them are not.

In this paper, we investigate the economics of domain name abuse and place several aspects on an empirical footing — both characterizing the economic value enabled by spam-advertised domain names and the concurrent economic impact of domain name blacklisting. In particular, we use almost two and a half years of sales data from two counterfeit pharmaceutical affiliate programs (com-

prising over 40K sites and 1.3M sales records) to characterize the economic role and extent of domain abuse to drive traffic to pharmacy sites. Using recorded Web referrer logs we infer the nature of Web services abused (e.g., free hosting sites, search results, Web-mail providers, etc.), and the dynamics of such abuse over time and the revenue afforded by different vectors. Specifically, we identify the nature of Web services abused that account for \$25M in revenue for SpamIt and \$41M for GlavMed during 2007–2010.

Against this backdrop, we then use nine months of contemporaneous data from a widely-used domain blacklisting service to quantify the impact of blacklisting on these same pharmacy store sites. Since our data set provides the precise time of every revenue event for each advertised domain name, we are able to directly investigate the extent to which domain blacklisting was successful as a strategic defense mechanism; did it undermine the fundamental business model or simply change marginal costs? Moreover, this data allows us to investigate hypothetical questions such as the extent to which improvements in blacklisting “speed” would impact profitability.

Interestingly, we found that existing blacklists rapidly identified a large percentage of spammed domains (88% within 2 days) and that additional improvements in blacklisting “speed” would, by itself, have little impact on profitability. Indeed, our findings suggest that domain discovery is a more important issue in the efficacy of domain blacklisting. To wit, over 60% of revenue for domains advertised through spam came from 12% of the sites in our data set that evaded blacklisting (either through luck or, as we observe in some cases, through careful advertising to avoid the sensors of defenders). Thus, even if blacklisting is otherwise robust, a small fraction of non-blacklisted domains may be sufficient to sustain overall profitability.

Another complex aspect of this problem is the interaction between consumer demand and how blacklisting is used. Domain blacklists are not universally used and in many cases they are only used in an advisory fashion (e.g., labeling email as “spam” that contains offending domains). However, we find strong evidence that motivated consumers are not dissuaded by such advisories. From the referrer logs in our data set, we found that *20 to 40 percent* of sales from email spam arise from users who *actively open their spam folder* and click on links to pharmacy sites. Indeed, this user behavior is one of the reasons that blacklisted domains in our data set earn *87% of their revenue after* being blacklisted. Using a simple revenue model to represent our data we establish that even if blacklists can identify all counterfeit pharmacy domains, blacklisting can make spamming unprofitable only when used to completely block access to offending domains.

While our data is limited to a particular time period (2007–2010) and a particular set of actors (GlavMed and SpamIT), we believe that the underlying conflict is largely unchanged today and that our key findings — that incomplete domain discovery and the advisory use of blacklists limits the strategic value of the approach — are likely to still hold today across a broad range of scenarios.

2. BACKGROUND

Abusive advertising, such as email spam, dates back to the origins of the Internet. The attraction of an advertising medium with virtually non-existent marginal cost is irresistible. While a combination of legal limits (e.g., the US CAN SPAM act [8]) and the creation of structured advertising vectors (e.g., sponsored search on popular search engines) have placed controls on legitimate advertisers, those who are already breaking the law by selling counterfeit or fraudulent products continue to abuse communication channels to shill their wares.

Today, virtually every form of Internet communication has an attendant form of spam: email [25, 37], search [20, 24, 49, 50], blogs and forums [33, 40], social networks [17], instant messaging [35] and so on. The prevalence of such widespread abuse of different services suggests profitability but the relative differences in the extent of abuse have not been studied previously. Over the years, different Web services have developed techniques to counter spam. By far the oldest of these is the email vector which became so prevalent that even as of 2013, spam was still the dominant form of email in transit [12]. To manage this problem, security researchers learned to classify mail as wanted or unwanted based on both its content and from where the message was sent. Thus was born Internet blacklisting.

The first blacklists focused on identifying and distributing the IP addresses of hosts known to be sending spam messages so that mail servers could know to properly drop or classify their messages [46]. A wide range of literature has focused on evaluating and improving upon such IP-based blacklisting approaches (e.g., [7, 21, 41, 42, 52, 53, 54]) but at their core they are “bot detectors” and thus their value is primarily in limiting the amount of mail that can be sent with impunity from a given host. However, if the advertiser has a large number of senders available (as with large botnets) or is able to “launder” their mail traffic through a major Web mail server [56], SMTP relay [18] or SMTP server, then this sort of blacklisting will be ineffective (i.e., one cannot blacklist the IP addresses for `hotmail.com`).

An alternative approach was designed by the anti-phishing community: URL blacklists. These systems distributed full URLs of sites known to be hosting counterfeit pages (typically representing banks or other financial institutions) and would be used either by mail servers (to classify emails containing such URLs) or by Web browsers (to block or warn users about to visit such URLs). A range of empirical studies have focused on evaluating the reaction time of such services, with results suggesting that the reaction time is short (typically a couple hours or less) [28, 39, 55]. More recently, a number of predictive approaches have been proposed, using some combination of the lexicographic features of URLs [29] or the characteristics of domain registration [16]. In practice, many high-volume URL blacklists have focused primarily on the registered domain in a URL. Feeds from such *domain blacklists*, such as the Spamhaus DBL [43] and the SURBL [44], have become standard inputs to virtually all enterprise spam filtering systems today.

In characterizing any blacklists, two questions need to be considered: how is the blacklist created, and how will it be used?

Today, since most blacklisting activity is driven by addressing the email spam vector, the blacklists are created via spam traps — open MX resolvers, honey accounts, botnet output or sometimes human-labeled spam messages [36]. By definition these lists can only detect abusive domains that are collected by these sensors. This truism is well known to spammers and “list washing” services abound (for example, <http://emallistcleaning.com/>) to remove honey and test accounts for output spam lists. Still other spammers traffic only in lists of likely customers (e.g., who have

purchased goods in the past). Finally, spammers who move on to other advertising vectors (e.g., search engine optimization) may experience no impact from blacklisting since there is no organized ecosystem to collect or distribute blacklists for that medium.¹

The second question is how the blacklist data is used. Email spam filtering software will typically use domain blacklist data as a strong feature in their classification algorithms. Thus, an email message advertising a given domain (i.e., including a URL with that domain in the message body) will be likely classified as unwanted and automatically filed in a “Spam” or “Junk” folder. In other situations, such as with anti-phishing toolbars or Web filtering software (e.g., such as offered by Websense or Cisco Ironport), users may be prevented from resolving DNS queries for domains on the blacklist even if they are allowed to click on the URLs.

Indeed, it is this last use that has generated the most controversy as governments have sought to legislate its use. For example, in China, comprehensive DNS filtering is used to prevent resolution of domains which the government deems as threatening [27]. However, this desire to use the DNS in this manner appears in democratic regimes as well. For example, the Australian Communications and Media Authority (ACMA) maintains a blacklist of Web sites and several administrations have proposed that ISP filtering of this list be mandatory (the two major Australian ISPs filter based on the blacklist on a voluntary basis) [19]. In the United States, the controversial Stop Online Piracy Act (SOPA) and Personal Information Protection Act (PIPA) would have required all ISPs to filter DNS requests to domains identified by brand holders as infringing on their copyright or trademark [15]. This last case generated tremendous opposition. However, most of the resulting arguments focused either on claims that it would have a chilling effect on innovation and potential infringement on free speech [13]. We are unaware of any academic evaluation about whether the statutes would have in fact prevented counterfeiters from still pursuing their business at a profit.

3. DATA SETS

At the core of our analysis are two data sets, originally described by the journalist Brian Krebs in his “PharmaWars” series [23] and documented more fully by McCoy et al. [31], that capture the full “back end” database for the GlavMed and SpamIt pharmaceutical affiliate programs between 2007 and 2010. As described in PharmaLeaks, these affiliate programs provided drugstore storefronts (including domain names and Web sites), drug fulfillment, payment processing, and customer service to independent affiliate advertisers who were paid on a commission basis [1, 38]. Thus an individual affiliate would be given one or more domain names to advertise and they would be paid a fraction of the revenue for every sale they brought through advertising using any vector (e.g., email spam or search engine optimization).

3.1 Authenticity and Ethics

As discussed in McCoy et al. [31], studying these leaked data sets raises concerns regarding authenticity and ethics. Here we briefly summarize the evidence that makes us confident about the authenticity of the data, and refer readers to [31] for a more detailed discussion of these concerns.² While there is no mechanism to ascertain the authenticity of this data beyond all doubt, we never found any inconsistencies in over 140 linked tables with

¹The Google Safe Browsing list contains URLs known for phishing or distributing malware.

²Excerpts from both data sets and additional discussion can also be found on Krebs’ blog [23].

over 2M sales records. We further compared the databases to the separately leaked corpora of metadata containing detailed chat logs from the program operators for both GlavMed and SpamIt and similarly found no inconsistencies. Moreover, we found these data sets accurately contain all of our past purchases [22, 25] in the database as further evidence of the authenticity of the data.

We address the ethical concerns surrounding the data using the same principle [31] of causing no additional harm in analyzing a leaked data set already in the public domain. We also reiterate that we again strictly adhere to our institution’s human subjects review process and ethical guidelines. For this study we only use anonymized data and do not mention any identifiable information about any person or institution who appear in the data other than naming the affiliate programs GlavMed and SpamIt themselves.

3.2 GlavMed and SpamIt

The data dumps of these two affiliate programs are in the form of complete, self-contained PostgreSQL databases but no other code external to the database. GlavMed and SpamIt were sister programs and therefore shared the same schema. Of the 140 tables in the database, we used four tables of which three (`shop_sales`, `shop_transactions`, `shop_affiliates`) were originally also used by McCoy et al. The `shop_sales` table contains details of every order such as timestamp, sale amount, etc. The `shop_transactions` table includes payment attempts and details of orders, and `shop_affiliates` contains information about affiliates such as when they joined the program and their user handle. Unlike McCoy et al., who focused on the nature of sales and the role of affiliates in these programs, our focus is on domain abuse. As a result, we also used the `shop_sites` table which contains domain information such as their `create_date` and the affiliate responsible for advertising the domains.

Besides basic order data, the `shop_sales` table also contains an HTTP referrer field which was previously not used by McCoy et al. For 45% of all sales in both programs combined, this field contains the URL that referred the customers to the shop storefronts. We use this field to determine how a pharmacy shop was advertised to customers: whether customers visited the Web site directly from a Webmail message (e.g., referrer domain is `hotmail.com`), a search result (e.g., referrer is `google.com` with search terms in the URL), etc. We further restrict our analysis to valid sales, i.e., sales for which all fraud checks passed, all test purchases are removed, and a valid credit card authorization is attempted (we do not perform further sanitization of sales beyond that performed by McCoy et al.).

Finally, when discussing blacklisting of SpamIt domains, we purposely omit “public shops”, domains which are shared among different affiliates (using a cookie or URL token to claim commission) and “reorder shops” (not advertised publicly, but provided to past customers for reorders) because we cannot attribute revenue to a particular affiliate or a mode of advertising. These sites account for just 0.1% of all sites in SpamIt.

3.3 URIBL

To assess the impact of blacklisting, we use the URIBL blacklist [47]. The data we extract from URIBL contains a timestamped list of spam-advertised blacklisted domains starting July 9, 2009. While URIBL is primarily reactive, it does include some predictive features and thus some domains appeared on it before they were seen in a spam trap (we confirmed our observation with URIBL). Therefore, to distinguish between the predictive listing of domains and domains that are simply reused at a much later point of time, we exclude all domains that appeared on the blacklist more than a

month before their recorded “`create_date`” (equivalent to 0.3% of all shop domains).

Moreover, we understand the inherent risk in characterizing the entire blacklisting defense mechanism using a single blacklist. Despite our efforts we were unable to acquire any other contemporaneous domain blacklist for this study that provided fine-grained blacklisting timestamps necessary for our analysis.

3.4 Spam Feeds

We also used two feeds of spam-advertised domains that we obtained from Pitsillidis et al. [36] between July 9, 2009 and March 18, 2010. We use these feeds to indicate, for instance, when spammers advertised the domains to customers in spam. The first feed consists of domains captured by MX spamtraps which are honeypot email addresses advertised to be visible only to Web scrapers searching for email addresses online. The second is a human-identified (HI) feed of domains contained in messages marked by users of a major Webmail provider as spam on the Web mail user interface. By construction the HI feed contains domains that were actually seen by a human whereas the MX feed contains domains that were indiscriminately advertised to all email addresses. The HI feed has two gaps from September 19, 2009 to October 7, 2009 and October 26 to November 12, 2010.

Since our spam feeds end on March 18, 2010, we only consider shop domains created through March 10, 2010 to allow for a week for domains to appear in the feeds. We believe this period is sufficient because over 90% of domains appear on each feed and blacklist within a week of their `create_date`.

Given the above constraints, our analysis of blacklisting only uses the overlapping subset of all these data sets (databases, blacklist, and spam feeds) between July 2009 and March 2010.

4. DOMAIN ABUSE

Our first goal is to understand how affiliates abused various domains and Web services to drive traffic to the pharmaceutical storefronts. We partition the domains in our data set into three categories. The first are domains that belonged to SpamIt and GlavMed and hosted storefronts where customers could purchase various drugs (primarily erectile dysfunction). We call these “shop sites” or “shop domains”. There are 51.6K such domains in SpamIt and 2.3K in GlavMed created between November 7, 2007 and April 30, 2010.

The second category consists of domains representing an *advertising vector*: external Web services through which customers discovered the shop domains. These include Webmail providers (e.g., Gmail, Hotmail, etc.) and Web search engines such as Google Search, Yahoo Search, etc.

The remaining domains are *infrastructure domains* that were used by affiliates to facilitate advertising via email and Web search, and to prevent exposing the shop domains directly to blacklists. These include free hosting domains (e.g., blogspot, geocities, etc.) which are legitimate sites where anyone can host free content, compromised private sites that did not belong to affiliates, and domains purchased by affiliates in bulk for the sole purpose of redirecting traffic to the shop domains.

A significant portion of GlavMed revenue (23.7% as shown in Table 2) also came from customers arriving at shop sites via traffic purchased from traffic sellers. As discussed later in this section, these services share characteristics with both advertising vectors and infrastructure domains, yet have a role distinct from the other categories and therefore we have included it as a separate category. Finally, there are some domains we were unable to classify in part

	Shop Sites	Sales	Revenue	Revenue/Sale	Affiliates
Advertisement vectors	11957	147582	\$18.05M (73.2%)	\$122.36	330
<i>Email spam</i>	11898	145041	\$17.83M (98.7%)	\$122.94	326
<i>Web search</i>	173	2541	\$0.23M (1.25%)	\$89.21	71
Infrastructure domains	1402	54351	\$6.58M (26.7%)	\$121.17	174
<i>Free hosting</i>	1282	45941	\$5.67M (86.1%)	\$123.37	154
<i>Bulk purchased domains</i>	120	7781	\$0.84M (12.8%)	\$108.36	50
<i>Compromised sites</i>	64	629	\$0.07M (1.13%)	\$119.33	27
Purchased traffic	11	199	\$0.02M (0.08%)	\$104.15	4
Uncategorized	863	4610	\$0.54M (2.20%)	\$117.91	165

Table 1: Classification of referrers used by SpamIt affiliates.

	Shop Sites	Sales	Revenue	Revenue/Sale	Affiliates
Advertisement vectors	1433	134977	\$13.8M (38.3%)	\$102.25	787
<i>Email spam</i>	615	10855	\$1.35M (9.81%)	\$124.79	578
<i>Web search</i>	1182	124122	\$12.4M (90.2%)	\$100.27	537
Infrastructure domains	1017	134832	\$13.71M (38.0%)	\$101.68	898
<i>Free hosting</i>	684	38094	\$3.91M (28.5%)	\$102.65	654
<i>Bulk purchased domains</i>	374	63639	\$6.27M (45.7%)	\$98.55	356
<i>Compromised sites</i>	456	33099	\$3.53M (25.7%)	\$106.59	393
Purchased traffic	458	86657	\$8.55M (23.7%)	\$98.68	366
Uncategorized	1047	45337	\$4.72M (13.1%)	\$104.21	890

Table 2: Classification of referrers used by GlavMed affiliates.

due to limitations on being able to find reliable contemporaneous historical data about the domains labeled as *Uncategorized*.³

While the shop sites are conveniently listed as such in the database dumps we received, we identified the advertising vectors and infrastructure domains using the HTTP referrers recorded for 30% of 690K SpamIt sales (accounting for \$25M) and 61% of 660K GlavMed sales (totaling \$41M). These referrers reflect the kind of Web site that led a customer to the shop site. For example, a customer arriving at a shop site after clicking on a URL for the shop site URL in an email message in Gmail will have a recorded referrer from `mail.google.com`. To classify referrers, we used features such as domain names, historical page content from The Wayback Machine [51], historical WHOIS information from DomainTools [14], and keywords in the referrer URLs. For some vectors, such as free hosting domains, we were able to find aggregated lists of domains online which we manually verified before using for classification of referrer URLs. Unfortunately, we do not have the entire redirection chain of URLs from a user’s click to the shop domain, but only the penultimate referrer that led the customer to the shop site in the next hop. However, contemporaneous data from Levchenko et al. [25] shows that 90% of the 8M spam-advertised domains they crawled using Firefox resulted in either zero or one redirects, suggesting that the redirection chains are likely short.

In the remainder of this section we present our analysis of subset of sales that have corresponding referrers.⁴ We start with some overall observations about the data. We then describe how we classified sales with referrers into the various categories, the sites and services that affiliates frequently targeted, and the spamming behavior of the top affiliates using each strategy. For reference, Table 3 shows example referrers in each category.

³Manually sampling these found them to be primarily bulk domains with some compromised domains.

⁴We can only speculate as to the remaining sales, but we suspect a large fraction arise from email clients that do not naturally transmit a referrer and, in some cases, from intermediate domains that explicitly strip referrers.

4.1 Overall Observations

Tables 1 and 2 breakdown the number of affiliates, sales, and revenue generated by the affiliates in each category. By intent, email spam was the dominant form of advertising used in SpamIt (Table 1). There was a moderate use of infrastructure domains (26.7% revenue dominated by free hosting) to mask the shop domains in the URLs advertised presumably also via email. In contrast, affiliates in GlavMed attracted customers mostly via Web search (Table 2) results. However, the use of various infrastructure mechanisms to facilitate traffic via Web search was more prevalent (38% revenue) and more evenly distributed than in SpamIt.

The differences in the use of infrastructure domains for SpamIt and GlavMed can be attributed to differing pressures in the dominant advertising channels (email and Web search, respectively) used by both programs. SpamIt affiliates needed to bulk advertise their domains repeatedly via email to maintain traffic volumes, while GlavMed affiliates could have placed their content on a compromised site once and, in return, received ongoing traffic from that site until it was identified and taken down by administrators. Similarly, while SpamIt affiliates had a cost structure that needed to accommodate adversarial blacklisting and filtering of email messages, GlavMed affiliates monetizing search traffic needed to maintain the rank of their shop sites for popular search terms. We discuss these differences further below in the context of individual categories.

Yet another interesting result of this classification is the revenue generated per sale. Notably, the average revenue per sale was relatively uniform at just over US\$100/sale for all categories in both Tables 1 and 2. No matter how an affiliate attracted customers, customers tended to spend the same amount of money regardless of the kind of URL they clicked on; there is little customer differentiation by strategy in terms of revenue. So the dominant goal for affiliates remained attracting as many customers as possible. Generally speaking, though, the top affiliates often used multiple infrastructure domains for redirection at a time, often emphasizing one kind over another over time in response to a dynamic environment.

Category	Referrer
Email spam	http://mail.live.com/mail/readmessagelight.aspx?action=markasnotjunk&folderid=...
Web search	http://search.yahoo.com/search?p=canadian+viagra&ei=utf-8&fr=blie7
Free hosting	http://groups.google.com/group/... http://www.umbc.edu/ddm/wiki/user:cheap_cialis http://answers.yahoo.com/my/profile?show=...
Bulk purchased	http://accutanewithoutprescription.org
Compromised	http://library.newschool.edu/askal/request/.inc/c/clomid-without-prescription.html
Traffic	http://traffic-analytics.net/tds/in.cgi?3&seoref=http://search.comcast.net/?... q=lavitra&http_referer=http://www.plantright.org/?id=49&default_keyword= http://klikcentral.com/traffic/in.cgi?11¶meter=buy%20viagra&seoref=http://www.google.com/search?q=buy+viagra&...&http_referer=www.vfcc.edu

Table 3: Example referrers for advertising vectors (email and Web search), infrastructure domains (free hosting, bulk, and compromised), and purchased traffic.

Whereas Tables 1 and 2 provide a summary overview, Figure 1 shows the temporal dynamics of the revenue from clicks on different kinds of domains over time (binned by weeks) for SpamIt and GlavMed. The dynamics in Figure 1 highlight the freedom of innovation of the affiliate program model, which provides the flexibility for different affiliates to explore different strategies for generating sales and the agility of affiliates to react to defensive pressures.

Even though the vast majority of revenue in SpamIt came from shop domains directly advertised via email, there was some use of infrastructure domains as redirection mechanisms. In July 2008, one affiliate began using free hosting sites. It was an effective strategy for a while, but gradually the free hosting providers were able to undermine the abusive practice. As free hosting dwindled, in January 2010 a small group of affiliates began using bulk domains for redirection as well, a profitable strategy for three months.

In contrast, the use of Web search for direct advertisement of shop domains was much smaller in GlavMed and there was significant use of infrastructure domains, presumably for search engine optimization (SEO). The revenue from different SEO efforts is more distributed. After a steady rise in sales throughout 2008, GlavMed experienced a jump in revenue primarily via purchased traffic and the use of bulk domains to direct traffic to shop sites. A rise in sales from direct advertising of shop domains on Web search contributed to another spike in January 2010.

4.2 Advertising Vectors

Email spam and Web search are the primary direct advertising vectors in SpamIt and GlavMed, respectively.

4.2.1 Email spam

Sales from email spam are those in the data set where users clicked on links to the shop sites advertised directly in email messages. Since SpamIt caters to email spammers, it is not surprising that email-based sales account for nearly all of its revenue.

We classified referrers in this category by matching domains of known popular Web mail providers (e.g., mail.google.com), regional Web mail providers (e.g., poczta.o2.pl), and keywords corresponding to known email clients (e.g., zimbra, squirrelmail, etc.). We also included sales from other online sites with internal message services, most notably Facebook. To validate likely but uncertain referrers, we manually inspected them by visiting the sites using the Wayback Machine [51]. The vast majority of sales came via spam to Web mail providers, with spam to Yahoo, Hotmail, and AOL accounting for 84% of the total revenue.

Historically, the goal of filtering email spam has been to prevent it from reaching the user’s inbox. To account for the possibility of false positives, though, services file messages classified as spam separately (e.g., into a spam or junk folder with a timeout) rather than deleting them immediately. Surprisingly, the sales records indicate that this filtering approach *does not* necessarily undermine revenue: despite such active filtering, users intentionally locate messages classified as spam and visit the storefront sites advertised in the messages. In effect, for some users the spam filtering, and contributing defenses such as blacklists, makes it *easier* for people to locate advertised storefronts.

Specifically, we were able to infer the folder from which an email spam message was clicked in 68% of referrers from Hotmail and 40% of referrers from Yahoo Mail. We used the well known folder names in Yahoo Mail and Hotmail to determine the folder that the user found the email in. For example, Table 3 shows a Hotmail referrer with the parameter *folderid*. A *folderid* of 5 corresponds to the spam folder while 1 corresponds to the inbox for Hotmail. Similarly, the parameter *fid* contains the name of the folder for Yahoo Mail referrers. We found that for Hotmail, over 20% of email-based sales came from customers who clicked on links in messages not in the inbox. Similarly, 39% of sales from Yahoo Mail referrers arose from non-inbox folders: 31% are from the *bulk* folder and the remaining 8% from various custom folders such as *online orders*, *cheap medication*, *viagra reorder*, etc. Such folder names clearly suggest that some people save these messages for future use and we also identified multiple referrers where users explicitly marked pharmaceutical spam as “not junk”. Table 3 shows such an example referrer for Hotmail. This evidence shows strong demand in the counterfeit pharmaceuticals market.

4.2.2 Web search

With Web search sales, customers arrived at pharmacy sites by directly clicking on shop domain URLs in results to Web search queries. Again reflecting the duality of the two affiliate programs, search-based sales are far more popular in GlavMed. Revenue from search results predominates in GlavMed at 31% of the total revenue, while it forms only a tiny fraction in SpamIt at 1.2%.

As seen above with users explicitly searching their mail folders, Web search sales again demonstrate customer demand in explicitly seeking out online pharmacy sites. We identified sales from all major search engines (Google, Yahoo, Bing, Ask and AOL), portal search sites such as search.rr.com, search.msn.com, search.orange.co.uk, as well as other sites that allow searching for

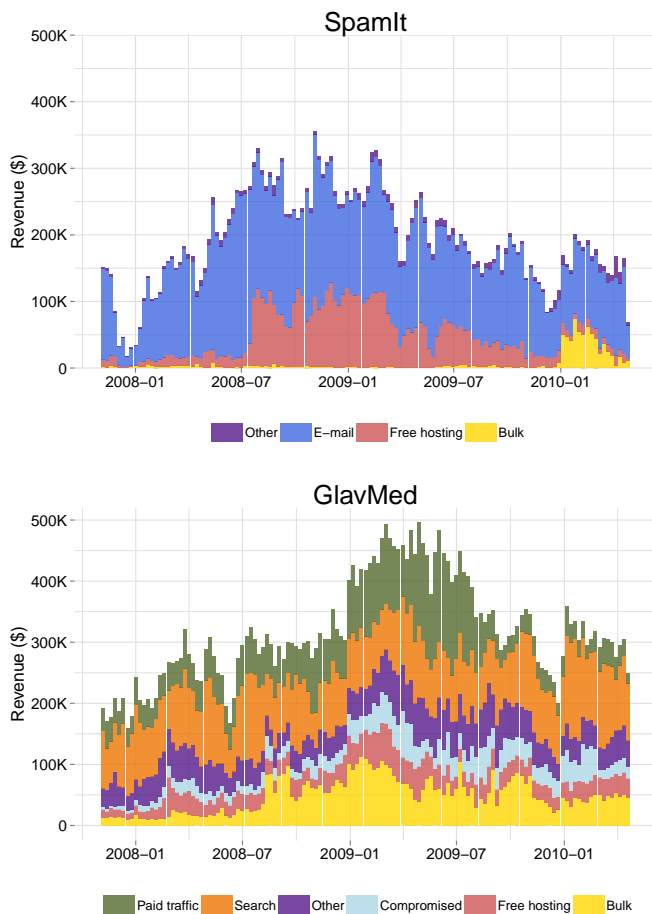


Figure 1: Revenue from clicks on different kinds of referrers.

arbitrary keywords. Referrals from the top two search engines at the time, Google and Yahoo, dominate GlavMed revenue at 78%.

Moreover, nearly all referrers include the keywords (e.g., *canadian viagra* for the URL shown in Table 3) for which the customer searched. For Google and Yahoo, the most popular keywords are *cialis* and *viagra*, respectively. These terms reflect the overwhelming demand for male enhancement products in these programs [31].

Over the period of study, GlavMed affiliates received steady sales from search results for shop sites with 4,137 sales on average per month, and all shop sites received at least one sale from search engines. The affiliate *webplanet* received the most search-based sales in GlavMed, evolving his strategy over time. *Webplanet* initially received search based sales from Yahoo and MSN Search, and did not start monetizing Google Search until April 2008. From that point, he attracted customers from both Yahoo and Google equally.

Lastly, we also observed referrers from searches on mobiles in GlavMed throughout the period of the data set. Although initially accounting for a negligible fraction of sales, monthly sales increased continuously over time — suggesting affiliates started to explore a nascent yet growing advertising vector.

4.3 Infrastructure Domains

In addition to using email spam and Web search to attract customers to their shop sites, affiliates also made use of other Web services and domains to boost traffic from both of these vectors.

4.3.1 Free hosting

Free hosting domains are sites where any user can post content for free. Spammers frequently abuse these sites by creating blogs, profiles, forums, and wiki pages, and posting comments, uploading images or other files, etc. The spammed content has links to entice potential customers to pharmacy shop sites.

We classified many of the free hosting domains in referrers using lists of such domains generally available online. Examples include *docs.google.com*, *spaces.live.com*, *imageshack.us*, etc. For domains that did not appear in our free hosting site list, in most cases we were able to notice free hosting sites when multiple referrers only differed in the profile identity string. We verified that these domains were in fact free hosting domains using Wayback Machine. We also identified forum abuse using keywords such as *viewtopic*, *discuss*, *showthread*, etc. We manually inspected referrers to distinguish between open forums and wikis used to freely post content, and forums and wikis hosted on compromised sites (Section 4.3.2). Table 3 shows 3 canonical examples of free hosting referrer URLs from our data set including a wiki hosted on *umbc.edu* that was used to create a page to advertise erectile dysfunction drugs.

We also included URL shortening services in this category, including *translate.google.com* exploited as a redirection service. Even though services such as *bitly.com* were very popular in 2009 [2], we only see a small number of sales via shorteners for structural reasons. Most popular shortening services respond with a 301 Moved Permanently HTTP status, causing the browser to resend the request to the final site using the original referrer. As a result, the referrer seen by the shop site is the site where the user clicked on a shortener link, not the shortener itself.

Free hosting was the most popular form of infrastructure domains used in SpamIt (Table 1). While free hosting abuse was less popular than bulk domain abuse in GlavMed, it still comprised 5–7% for both affiliate programs. However, the nature of free hosting abuse differs for SpamIt and GlavMed because of differing objectives with these sites. SpamIt affiliates used free hosting to host content on trusted domains to overcome blacklisting-based content filtering (i.e., blacklists do not list *google.com* as a bad domain because of some abuse on *docs.google.com*). Thus, various services on *google.com*, *live.com*, *yahoo.com*, and *imageshack.us* are the most abused free hosting services among SpamIt affiliates. Google Groups was abused most effectively, typically using bogus group profiles, and accounted for 29% of the SpamIt revenue via free hosting sites.

In contrast, the motivation for free hosting abuse among GlavMed affiliates is to attract traffic by boosting search engine ranks of their domains. Also, abusing a large range of redirection sites causes multiple results linking to the same shop site to show up when a potential customer queries for pharmacies. Thus, among GlavMed sales we notice abuse of a larger number of free hosting sites (3,956 unique domain names vs. 830 among SpamIt) and sales are spread more evenly among domains: *backpage.com* was the most abused domain but accounted for only 6% of all free hosting abuse in GlavMed.

For spammers, a disadvantage of using free hosting sites is that once the abused domain removes the offending content, the spam links break and no longer point to the affiliate’s shop site. While we do not know the time it takes sites to detect and takedown spam pages, in at least one case correlating spammer behavior with news reports of abuse suggests that takedowns by free hosting sites require months to be effective.

Further, spammers seamlessly switched targeted sites in the face of such takedowns. Figure 2 illustrates the agility of a top SpamIt

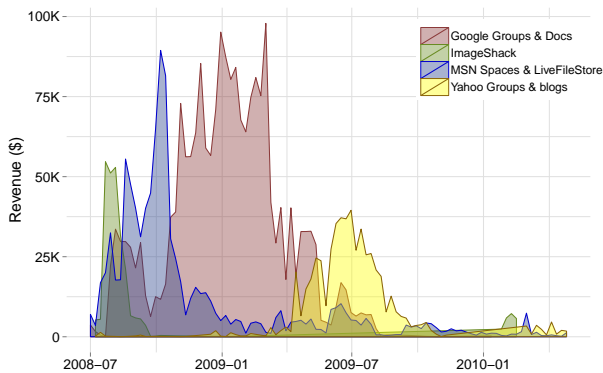


Figure 2: Spammers seamlessly switch from one free hosting site to another in the face of takedowns.

affiliate, *master*, who accounted for 93% of all sales via free hosting referrers. From July–August 2008, a large fraction of *master*’s revenue came via ImageShack. Subsequently, in August–September 2008 reports of abuse emerged of email spammers sending messages with links to Flash files hosted on ImageShack [9]. Revenue via ImageShack in SpamIt immediately declined (suggesting takedowns of advertised files by ImageShack), while *master*’s free hosting revenue from `spaces.live.com` increased from August–September 2008 using automatically created profile pages on MSN Spaces. This pattern repeats. Revenue from `live.com` almost entirely disappears in November 2008, coinciding with a Spamhaus report that ranked Microsoft as the fifth most spam friendly ISP [3]. *Master*’s sales switched to heavy and almost exclusive abuse of Google free hosting services between December 2008 and February 2009. At this point revenue from Google Groups declined significantly, once again coinciding with Spamhaus ranking Google as the fourth most spam friendly ISP [4]. *Master*’s free hosting sales then switch to Google Docs followed by a brief switch to the abuse of Yahoo in June–July 2009.

Such trends demonstrate that even when free hosting sites took action against spammers, the general prevalence and availability of such sites enabled a skilled spammer to quickly and seamlessly switch to newer services in the face of takedowns.

4.3.2 Compromised sites

We also found a large number of referrers to domains in GlavMed that appear to be compromised sites. Such sites are valuable for their search rank in poisoning search results [20, 49] for attracting traffic to store fronts: over 66% of these domains are under `.edu` or `.gov` TLDs, which purportedly have higher search engine rank.

Using DomainTools [14] and the Wayback Machine [51], we found legitimate sites hosting spam content in subdirectories. Often referrers from compromised sites contain content in either hidden directories (e.g., `.inc` as shown in Table 3), or subdirectories intended for other purposes (e.g., `css`, `images`). Some hackers added content to compromised sites in a signature style, facilitating matching. For example, eight affiliates received traffic from sites where the redirecting page had content placed in a directory named `md`.

A challenge in identifying compromised sites is distinguishing whether wiki and forum hosting software was compromised, or if spammers just created their own free pages on publicly accessible wikis and forums. We classified discussion and message board abuse as free hosting abuse (Section 4.3.1). For wikis, however,

we used the Wayback Machine to see if they were open for public editing at that time. If not, we considered them compromised.

As discussed in Section 4.4, hackers often compromise sites and install malware to direct customers to traffic buyers on demand. As a result, the relation between compromised sites and affiliates is often not one-to-one. Nearly 36% of compromised sites redirected traffic to multiple affiliates, while 65% of affiliates receiving this kind of traffic received traffic from more than one domain.

The most effective affiliates used compromised sites differently. GlavMed affiliate *glavmed2* received the most revenue (12%) from 44 compromised sites, with one site (`arkansasbaptist.edu`) accounting for 81% of his revenue. While *glavmed2* primarily monetized just one site, affiliate *grbk* received the second highest revenue more evenly distributed among 268 different sites.

SpamIt affiliates received a negligible amount of traffic from compromised sites. The primary advantage that compromised sites offer to email spammers is the reputation of the advertised site in the spam filter calculation — an advantage that free hosting sites offer as well, but at a much lower cost.

4.3.3 Bulk purchased domains

Affiliates purchase bulk domains as intermediaries for redirecting users to shop sites. Many bulk domains contain pharmacy-related keywords such as `accutanewithoutprescription.org`, `tramadol-shop24.com`, etc. The pharmacy content is typically at the root of these sites (Table 3 shows an example), distinguishing them from compromised sites where the content is on pages deeper in the name hierarchy. Furthermore, each domain redirected sales to just one affiliate, suggesting that these were owned by the affiliates themselves.

The revenue from the use of bulk domains by SpamIt affiliates to redirect to shop domains is much smaller than the use of free hosting sites (Table 1) even though bulk domains are inexpensive and can be purchased in large numbers conveniently. This small use is perhaps because bulk domains advertised in email spam are blacklisted very quickly and therefore do not offer much advantage over spamming shop domains directly.

For GlavMed affiliates, though, bulk domains offer advantages similar to compromised and free hosting sites by potentially increasing the number of results that appear on search pages and attracting more traffic. We counted 63,639 sales from 1,957 distinct bulk referrer domains, and nearly 23% of all the affiliates in GlavMed used bulk domains as a mechanism to attract buyers.

The use of bulk domains by the two highest earners once again highlights the flexibility of the affiliate program model. Venerable affiliate *webplanet* received the most revenue via bulk domains (43%), with over half of the revenue coming from three domains (`bestedmed.com`, `newedpills.com`, and `thebette rsxml.com`) and the remainder from 127 other domains.

In contrast, affiliate *andrew13plus* had the next highest number of sales but used a different strategy for monetizing bulk purchased domains. In particular, *andrew13plus* apparently purchased domains after they had expired, but also once they had accumulated useful search rank (e.g., `carrollfootball.org`, `sharonlnorris.com`, etc.). As a result, most of the sales from the domains came during the first month of use. The remaining revenue came within three months, after which the ownership of both of these sites, as well as their search potential, changed.

4.4 Purchased Traffic

Third-party advertising in general is a popular method for attracting traffic to Web sites, and affiliates of pharmaceutical sites are no exception. Purchased traffic comprises the second largest source of

revenue (24%) for GlavMed affiliates. We grouped traffic providers into three classes ranging from legitimate to suspicious services.

First are premier advertising services such as Google Ads. Early in 2008 affiliates experimented with ads with such services, but soon abandoned them presumably because of the high cost of popular pharmaceutical keywords.

The second class consists of traffic distribution systems (TDSs) such as `traffic-analytics.net`. These services act as intermediaries buying and selling traffic [45]. Frequently TDS kits are also installed on compromised sites to gather traffic, which is then monetized in a variety of ways including forwarding traffic to pharmacy sites, fake anti-virus, malware distribution, etc. [10]. A distinguishing characteristic of these referrers is that the referring domains point to several affiliates and the referrer URLs frequently have affiliate IDs. Table 3 shows a referrer for a user who searched for *levitra* on `search.comcast.net` and clicked on `plantright.org`, a site owned by Sustainable Conservation between 2007–2010. Upon clicking, the customer was redirected to `traffic-analytics.net` which then sent the customer to a GlavMed affiliate shop site.

Purchased traffic from TDSs are attractive for affiliates because their illegitimate nature enables them to be cheaper than premier ad services; for instance, the price for keywords such as *viagra* and *cialis* varies between 30–90 cents per click based on bids from the TDS RivaClick, while bidding for the same term on a mainstream advertising network costs several times more.

The third class consists of content providers who gather traffic to sell to TDS vendors. These two actors can also be the same. For example, `klikcentral.com` was a fake search engine for certain categories such as pharmaceuticals, cruise deals, degrees online, etc., and where the results are primarily ads. However, we also have evidence of it gathering traffic from compromised sites via Google Search. Another referrer in Table 3, for instance, shows a user who queried for *buy viagra* on Google Search and landed on `www.vfcc.edu`, and redirected through `klikcentral.com` to the pharmacy store site.

Some sites also use the guise of legitimate search engines for pharmaceuticals, but actually just gather and resell traffic (e.g., `viagra-prices-comparison.com`). All of these sites redirect customers to more than one affiliate, including one such site, `topmeds10.com` which redirected customers to as many as 73 affiliates. A large number of traffic buyers and sellers are merely intermediate domains (e.g., `klikcentral.com`) that specialize in search engine optimization, often using compromised sites. As such, they also act as infrastructure domains.

The top GlavMed affiliate using purchased traffic was once again `webplanet`, who received 39% of the revenue from purchased traffic. Through September 2008 `webplanet` received few sales from purchased traffic (83/month), but then invested heavily over several months in purchased traffic and averaged 3,537 sales per month.

5. BLACKLISTING

Since affiliates rely upon domains to host shop sites and intermediate sites, a common defense is to blacklist such domains. In particular, email services use domain blacklists to identify shop domains advertised via email spam, and classify incoming messages containing these domains into designated spam folders. In this section, we describe the revenue impact of blacklisting on gross revenue for spam-advertised pharmaceutical campaigns.

We use the URIBL blacklist described in Section 3 as a representative list of blacklisted domains and the leaked sales data sets as ground truth for the sellers impacted by blacklisting. As mentioned in Section 3, we restrict our analysis to the nine month period from

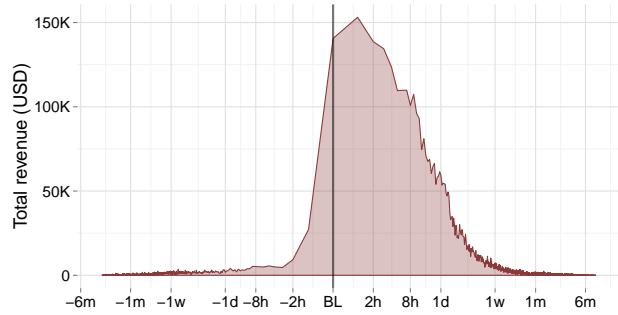


Figure 3: Revenue of domains before and after blacklisting. Note that the x -axis is non-linear.

July 9, 2009 to March 18, 2010 for which all our data sets (leaked sales, blacklist, and spam feeds) overlap.

As an email-based blacklist, URIBL identified 88% of 40K SpamIt shop domains as offending, most of them within two days of the creation of domains. Unsurprisingly, it identified only 4% of 1K GlavMed domains which are advertised predominantly via Web search. Therefore we restrict our blacklisting analysis to SpamIt domains only. This subset of data has 40K SpamIt shop domains that received 137K sales grossing \$15.6M total revenue.

In our analysis, we use four parameters — common to all blacklisting based defenses — to assess the impact of blacklisting.

5.1 Blacklisting Speed

The first aspect we consider is the time it takes for a spam domain to appear on a blacklist. Figure 3 shows the revenue distribution of the 35K blacklisted domains before and after blacklisting. We normalize domains by the time they appear on the blacklist: time zero is their blacklisted time, revenue earned before being blacklisted is negative in time, and revenue afterward is positive in time. Relative to their blacklist time, the curve shows the amount of revenue earned from customers across all domains per hour.

The figure shows a number of interesting results. First, it shows that domains receive most of their revenue *after* they are blacklisted. The revenue before blacklisting was \$740K, or just 13% of the total revenue of 5.9M (Table 4) from the blacklisted domains. One explanation is that blacklisting did not have a universal effect: customers may have received spam from email services that did not use the blacklist for spam filtering or deployed it sometime after the blacklist was updated and hence blacklisting had no effect on them. However, as found in Section 4.2.1, even when email services use blacklisting to classify spam, customers still explicitly searched their spam folders and clicked on URLs with shop domains to make purchases.

In effect, while blacklisting does not immediately stop the domains from earning revenue, it does set a lifetime to their earning potential. Revenue rises sharply just before blacklisting, peaks immediately after blacklisting, and then drops substantially over the next two days. In contrast, domains that are not blacklisted earn revenue over much longer time spans (weeks and months) With blacklisting, domains no longer have significant earning potential after a few days and affiliates have to purchase new domains to continue to earn revenue by advertising with spam.

Therefore, keeping everything else constant, placing these domains on blacklists even “faster” would not have prevented domains from being monetized.

5.2 Coverage

During the period of study, URIBL identified an impressive 88% of all SpamIt domains. However, while missing only 12% of all domains attests to the diligence of the blacklist maintainers, this minority of domains still accounted for 62% (\$9.7M as shown in Table 4) of the total revenue from all blacklisted and non-blacklisted domains combined. Evading blacklists was clearly advantageous for affiliates.

While there is evidence that affiliates made some effort to evade blacklisting by making much greater use of free hosting and bulk domains as layers of indirection (Table 4), this difference is rather small. Even the non-blacklisted domains received most of their sales from clicks directly in email messages. Thus, it appears that these domains evaded blacklisting not just because of *how* they were advertised but also because of *who* they were advertised to.

Blacklists are typically created using MX honeypot accounts that do not belong to real people. Also, McCoy et al. [31] showed that some affiliates were far more successful than others in their ability to spam effectively and earn more revenue. We found evidence suggesting that affiliates who evade blacklisting do not send their email messages to spam traps. As shown in Table 4, 88% of the blacklisted domains appeared on the MX feed that consists of URLs seen in spam traps, while only 0.5% of the non-blacklisted domains ever appeared on these feeds.

On the other hand, our Human Identified (HI) feed identified 96% of the blacklisted domains as being spam, while only 25% of the non-blacklisted domains appeared on it.⁵ The observation that non-blacklisted domains were significantly more likely to be seen by a human rather than a spam trap suggests that there was a difference in the nature of the email addresses used to advertise blacklisted and non-blacklisted domains. We attribute this difference to the sophistication of spammers advertising them.

To improve blacklist coverage going forward, one possibility is to extend the provenance of domains beyond just spam-advertised domains. For example, crawling spam advertised links can determine whether a URL for an otherwise legitimate domain might lead to a pharmacy (or other counterfeit storefront) domain. Another possibility is for services that maintain human-identified feeds (e.g., Webmail providers) to share domains from spam-advertised URLs with blacklist maintainers.

Finally, if we assume that we could have blacklisted domains that managed to avoid being listed — and that doing so would have caused them to monetize similar to the blacklisted domains — then these domains would have earned just \$168 on average as opposed to the \$2038 that they actually earned. Thus, discovering every additional domain would have reduced spammer revenue by \$1870 (92% of its original revenue).

5.3 Blacklisted Resource

As an intervention against affiliate spammers, blacklisting could use any of the uniquely identifiable resources used by spammers such as the IP addresses of hosts sending spam messages, domain names hosting storefronts, or even the bank accounts used to process transactions [25]. We next analyze the efficacy of choosing domains as the resource that is blacklisted.

In 2009–2010, the bulk price of a domain varied between 15¢ for a .cn domain [26] to \$7 for a .com domain [48]. Also, purchasing domains in bulk can also be automated, making the effort to replace a blacklisted domain negligible. Moreover, infrastructure domains

⁵The stated number of domains appearing in HI is a lower bound because four weeks of data is missing in the feed.

	Blacklisted	Non-blacklisted
Shop-sites	34959	4751
Sales	56K	80K
Revenue	\$5.9M	\$9.7M
Sales/Site	1.6	16.9
Affiliates	119	144
Sites seen in feeds	34771 (99%)	1193 (25%)
<i>MX</i>	30647	27
<i>HI</i>	33701	1185
Sales with referrers	6076	28576
<i>Email Spam</i>	4798 (78.9%)	18206 (63.7%)
<i>Purchased Traffic</i>	-	168 (0.59%)
<i>Free hosting</i>	284 (6.31%)	4507 (15.8%)
<i>Compromised sites</i>	8 (0.13%)	124 (0.44%)
<i>Web search</i>	7 (0.11%)	22 (0.07%)
<i>Bulk Purchased Domains</i>	709 (11.7%)	4375 (15.3%)
<i>Uncategorized</i>	170 (2.79%)	1172 (4.10%)

Table 4: Statistics differences between blacklisted and non-blacklisted domains.

such as free hosting domains can be useful for evading blacklisting, yet are abundantly available at low prices as well.

During the nine months we consider for blacklisting, affiliates used 40K domains and 88% of them were blacklisted. Assuming that blacklisting forced affiliates to replace the listed domains, the total cost of replacement for domains was \$245K when conservatively assuming \$7 per domain. This cost is only 1.6% of the total revenue from this period. Again, given the low costs of domains and the relatively much higher revenue earned per domain, blacklisting did not impose a serious cost of replacement. In Section 6 we estimate the cost per domain that would have made blacklisting prohibitively expensive for spammers.

5.4 Blacklisting Penalty

The last aspect we consider is the penalty of having domains blacklisted. During 2009–2010, blacklisting was used to identify and filter spam messages into designated *spam* or *junk* folders, where they would remain for up to a month for most Webmail providers. In the absence of demand, such spam filtering would have hidden unwanted ads away from user inboxes. As discussed above, though, despite blacklisting affiliates continued to receive sales at least in part as a result of market demand: 20–40% of customers accessed messages in their spam folder, searched for drugs by name, etc.

Thus, for the spam-advertised unlicensed prescription drug market, the penalty imposed by a classification-based defense was overshadowed by the demand for these drugs. As discussed in Section 4.2.1, classifying messages as spam effectively made it easier for people to find these storefronts and thereby enabling a domain, at least to some extent, to earn 87% of its total revenue after blacklisting (Figure 3). In the next section, we also consider the effects of increasing the blacklisting penalty.

6. DISCUSSION

While the previous section describes the lackluster role played by blacklisting in our data set, it begs the larger question of whether these results might change if blacklisting were deployed more aggressively. To reason about this issue, we parametrize a simple model of blacklisting impact and then explore the general implications for blacklisting email advertised domains.

6.1 A Simple Revenue Model

In general, prediction and extrapolation can be highly error-prone, and this is only enhanced by the presence of an intelligent adversary. Thus, we use a very simple model to capture the impact of blacklisting on profitability. As shown in Section 5, blacklisting does have a significant opportunity cost for the spammer, but our assumption is that the central goal of domain blacklisting is to hurt spammer revenue sufficiently to make such unsolicited email based domain advertising unprofitable altogether. Thus, we parametrize a model to determine the conditions under which blacklisting can achieve this goal.

To formulate our model, we first assume that there is a sufficiently good detection capability that all abusively-advertised domains will eventually be blacklisted (an assumption we will return to in Section 6.3). We then model the marginal revenue \mathcal{R} from a domain as the composition of the revenue *before* blacklisting and the revenue *after* blacklisting. To describe the pre-blacklisting revenue in a parsimonious fashion, we use a single parameter α to represent the mean revenue per unit time and we assume that this revenue remains constant from the time a domain is first advertised until it is eventually blacklisted. Once the domain is blacklisted we assume that, as in our data set, sales decline swiftly and thus all additional revenue can be captured by a single parameter β .⁶

Thus, we describe the marginal revenue per domain as:

$$\mathcal{R} = \alpha * t + \beta$$

Finally, assuming that blacklisting causes spammers to replace a domain, we estimate the marginal cost for every domain is c . We estimate this cost as only the cost of purchasing and registering a domain name. We assume that all other costs associated with replacing a blacklisted domain name such as creating a new Web site for the pharmacy storefront or even sending out spam emails to advertise a domain are negligible because these processes can be automated and do not require appreciable resources for every additional domain at scale. Even the cost of attempting to evade blacklisting (through the use of “list washing services”) is amortized over all domains because the same email address lists can be used for all domains. Thus, such efforts do not impose any marginal cost for a new domain acquired due to blacklisting. Past work in the area also suggests that domain registration is the dominant per-domain cost [31].

Combing these, the marginal profit per domain for a spammer is:

$$\mathcal{P} = \alpha * t + \beta - c$$

This simple linear relation does not capture the variation in revenue earned per domain (either before or after blacklisting). Since we are dealing with aggregates we believe this approximation is sufficient to examine gross effects.

Using the blacklisted domains in our data set we can thus empirically calculate the values for these parameters, with the average pre-blacklisting revenue per domain per day α as \$1.14 and the average post-blacklisting revenue β as \$104.19. These values reflect some combination of contemporaneous demand for pharmaceutical products and the intensity of the advertising effort and are by no means universal. Domain registration cost varies considerably with time and TLD, ranging from a low of roughly \$0.15 for .cn domains circa 2008—2009, to \$7 for retail .com domains during roughly the same time period as mentioned in Section 5.3 (in practice, bulk domains for selected TLDs are readily available today

⁶Recall from Section 5.1 that 87% of income for blacklisting domains occurs *after* blacklisting, due to some combination of demand, delays in using blacklist data, or non-universal deployment of blacklist-based filtering.

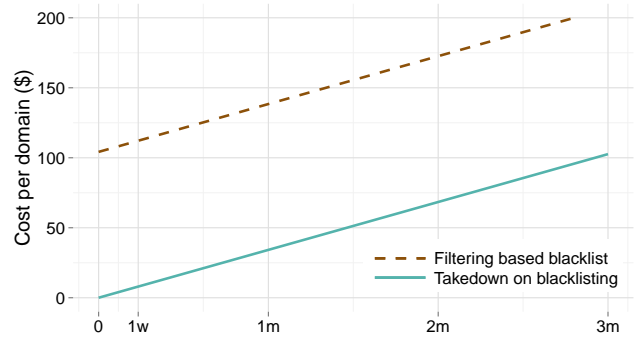


Figure 4: The highest cost of domain a spammer can afford (y -axis) against the time delay (x -axis) in blacklisting.

for \$2–3 each from a broad range of resellers). At any price, however, it is clear that the marginal revenue per domain from our data set is at least an order of magnitude greater than the marginal cost imposed on the spammer due to blacklisting of a domain.

6.2 Changing Blacklisting Penalty

Given our model, there are two factors that we now consider: how quickly a newly advertised domain is blacklisted and the regime in which blacklisting is used to undermine the advertising vector, either filtering (e.g., as in anti-spam) or blocking (e.g., as in DNS blocking or registrar takedown). In the nomenclature of our model this corresponds to varying t and β (setting it to zero in the case of blocking). We capture both of these effects in Figure 4 which plots the minimal cost per domain that a spammer can afford (i.e., the break-even point) for both regimes. The dotted line corresponds to a filtering regime, like spam filtering, in which revenue is acquired (β) even after a domain appears in a blacklist. As per the empirical parameter values described earlier, even if a domain were blacklisted *instantly*, the post-blacklist revenue is such that per-domain costs would need to be greater than \$100 to make advertising unprofitable.

The solid line, however, reflects a regime in which a domain ceases to generate revenue once it has been blacklisted (e.g., because the domain name is shut down by the registrar or, *a la* the proposed SOPA legislation, because DNS resolvers refuse to lookup the associated A records). Thus $\beta \rightarrow 0$, and the break-even point is represented simply by $c = \alpha * t$. In this case, there is a meaningful interplay between the time to blacklist and the practical cost of domains. Even the nominal cost of \$2.28 per domain (in line with current prices for cheap bulk registration) is sufficient to undermine the profitability of blacklisted domains.

These results suggest that even large reductions in blacklisting latency would not have made costs prohibitive for spammers, whereas increasing the penalty of being listed on a blacklist could have had more severe consequences for the domains that were identified.

6.3 Increasing Coverage

The discussion above neglects the small number of non-blacklisted domains in our data set altogether. We do not consider them here since, absent a blacklisting date, we cannot reason about how their revenue changes before and after blacklisting. However, these domains constitute almost two-thirds (Table 4) of total revenue (even for the email spam vector) and hence were a major source of revenue for spammers.

Putting these findings together suggests that while blacklisting did have an impact on spammer revenue, characteristics such as high consumer demand, sophistication of spammers, and the reactive use of classification based blacklisting made it far less effective. Our blacklisting analysis only focused on SpamIt during 2009–2010 when counterfeit pharmaceuticals dominated email spam [11]. However, after the shutdown of the SpamIt (possibly due to its inability to process MasterCard payments [5]), the global spam volume has been steadily declining from 87% in 2009 to 70% in 2013 [12]. Qualitatively, spam is now dominated by phishing messages instead of ads for pharmaceuticals and other counterfeit goods. Counterfeit pharmaceuticals are now almost solely advertised through non-email vectors discussed in Section 4.

But despite these changes in the global spam trends, filtering based blacklisting continues to be the primary form of intervention for email spam. Spam still dominates global email volume suggesting that even the current monetization methods from email spam are profitable. Thus blacklisting remains an important intervention mechanism. Consequently, our findings that further improvements in sensors to identify domains and better data sharing on spam domains is necessary to defend against email spam remain applicable today. Our quantitative analysis of the revenue impact of domain blacklisting on email spam is limited to the counterfeit pharmaceuticals market because features such as consumer demand and conversion rates vary for different markets. However, our results remain applicable for revenue from counterfeit pharmaceuticals sold through non-email vectors. Indeed, improved domain based classification tools for these vectors (such as social networks, ads, search) still have the potential to significantly affect the business of online counterfeit pharmaceuticals.

7. RELATED WORK

Various aspects of the online pharmaceutical industry have been discussed in the research literature including its use of abusive advertising and defense mechanisms against the same [6, 34, 50, 53]. For example, Leontiadis et al. [24] studied the use of compromised sites to drive traffic to online pharmacy storefronts. Our work differs both in providing a wider analysis of domain abuse (email, search, free hosting, traffic sellers, etc.) and, more critically, in examining the relative *revenue* provided by such traffic. We provide analysis of ground truth data and explain the most successful abuse strategies for spammers in the face of takedowns. Similarly, Moore et al. [32] also look at the potential impact of blacklisting pharmaceutical domains (in particular for search engine result traffic). By comparison, our work is narrower (restricted to only those sites of a particular set of pharmaceutical affiliate programs) but is comprehensive within that set and we are able to analyze the true economic activity in dollars, rather than using proxies such as site popularity.

More broadly, the use of blacklisting for filtering email spam has been a popular topic for many years. However, most of the work in this space is aimed at evaluating and improving the mechanical aspects of blacklisting-based defenses such as speed and coverage [16, 29, 39, 42, 55]. By contrast, our paper has focused on the larger question of the extent to which blacklisting efforts impact the profitability of the underlying business enterprise. Closer to our work in motivation are recent efforts focused on “payment intervention”, an intervention which seeks to undermine the profitability of abusive businesses by blocking their ability to obtain consumer payments [25, 30]. Our paper explores similar questions, but focuses on a different point of intervention (domain names).

Finally, our work builds on and supplements that of McCoy et al. [31] which uses overlapping, but different aspects of the same data set. McCoy et al. focused on analyzing the nature of global

demand for counterfeit pharmaceuticals, the role of third-party affiliates in the industry, and the cost structure of such businesses, while our work is concerned with the role played by domain names within this business model and uses external data sets to measure the impact of by blacklisting defenses.

8. CONCLUSION

There are as many ways to spam as there are to communicate, yet virtually all are Web-centric and require user clicks to convert. This commonality makes domain blacklisting a highly attractive mechanism for managing unwanted ads. Indeed, all evidence suggests that blacklisting is a quick and largely comprehensive process (at least for email spam, which has an active blacklisting ecosystem). However, the success of domain blacklisting has done little to stem the tide of email spam (let alone other abusive advertising practices). That a defense can simultaneously achieve its goal, yet not appreciably bother the adversary, is counter intuitive yet this fairly describes the current state of affairs today.

Our study of thousands of online pharmaceutical sites demonstrates that a combination of appreciable demand for counterfeit pharmaceuticals (indicated by the large fraction of revenue arising from the email spam folders and Web search queries leading users to pharmacy sites), the ability of sophisticated spammers to evade blacklisting and heavily monetize a small number of domains, and the existence of multiple vectors for traffic ensures that online pharmaceuticals business remains profitable. As such in this context of high demand and significant returns from successful evasion, domains alone do not impose a significant resource cost to the attacker who utilizes the whole spectrum of advertising strategies and is agile in the face of takedowns.

Finally, our results suggest that changing the blacklisting penalty for email spam from simply filtering messages in spam folders to completely blocking access to domains is necessary to grossly undermine the profitability of the blacklisted domains. Even then blacklist evasion and the other traffic vectors remain possible and lucrative alternatives for the attacker.

Acknowledgements

We thank Brian Krebs for providing the Glavmed and SpamIt data sets used in this study. We are also in debt to Andreas Pitsillidis and Kirill Levchenko for their assistance creating, interpreting and using the affiliate databases (as well as useful comments on this work). We are also grateful to `uribl.com` for maintaining the blacklist data and for answering our queries, and to Mark Felegyhazi for managing the blacklist data on our end. We also thank the anonymous reviewers for their valuable feedback. This work was supported in part by National Science Foundation grant NSF-1237264, by the Office of Naval Research MURI grant N000140911081, and by generous support from Yahoo, Google, Microsoft, and the UCSD Center for Networked Systems (CNS).

9. REFERENCES

- [1] Behind Online Pharma. From Mumbai to Riga to New York: Our Investigative Class Follows the Trail of Illegal Pharma. <http://behindonlinepharma.com>, 2009.
- [2] Ben Parr. Bit.ly is Eating Other URL Shorteners for Breakfast [Stats]. <http://mashable.com/2009/10/12/bitly-domination/>, 2009.
- [3] Brian Krebs. Spamhaus: Microsoft Now 5th Most Spam Friendly ISP. http://voices.washingtonpost.com/securityfix/2008/11/spamhaus_microsoft_now_5th_mos.html, 2008.

- [4] Brian Krebs. Spamhaus: Google Now 4th Most Spam-Friendly Provider. http://voices.washingtonpost.com/securityfix/2009/01/google_now_4th_most_spam-frien.html, 2009.
- [5] Brian Krebs. Spam Affiliate Program Spamit.com to Close. <http://krebsonsecurity.com/2010/09/spam-affiliate-program-spamit-com-to-close/>, 2010.
- [6] Chris Kanich and Christian Kreibich and Kirill Levchenko and Brandon Enright and Geoffrey M. Voelker and Vern Paxson and Stefan Savage. Spamalytics: An Empirical Analysis of Spam Marketing Conversion. *cacm*, 52(9):99–107, Sept. 2009.
- [7] Christian Kreibich and Chris Kanich and Kirill Levchenko and Brandon Enright and Geoffrey M. Voelker and Vern Paxson and Stefan Savage. Spamcraft: An Inside Look at Spam Campaign Orchestration. In *Proceedings of the USENIX Workshop on Large-scale Exploits and Emergent Threats (LEET)*, pages 4:1–4:9, Boston, MA, Apr. 2009.
- [8] Controlling the Assault of Non-Solicited Pornography and Marketing Act of 2003. <http://www.gpo.gov/fdsys/pkg/PLAW-108publ187/pdf/PLAW-108publ187.pdf>, 2003.
- [9] Dancho Danchev. Malware and Spam Attacks Exploiting Picasa and ImageShack. <http://www.zdnet.com/blog/security/malware-and-spam-attacks-exploiting-picasa-and-imageshack/1852>, 2008.
- [10] Daniel Cid. Large Scale Compromises Leading to Traffic Distribution System. <http://blog.sucuri.net/2013/02/large-scale-compromises-leading-to-tds.html>, 2013.
- [11] Darya Gudkova. Spam Evolution: January – March 2009. https://www.securelist.com/en/analysis/204792061/Spam_evolution_January_March_2009, 2009.
- [12] Darya Gudkova. Kaspersky Security Bulletin. Spam Evolution 2013. http://www.securelist.com/en/analysis/204792322/Kaspersky_Security_Bulletin_Spam_evolution_2013, 2014.
- [13] Derek Broes. Why Should You Fear SOPA and PIPA? <http://www.forbes.com/sites/derekbroses/2012/01/20/why-should-you-fear-sopa-and-pipa/>, 2012.
- [14] DomainTools. <http://www.domaintools.com/>.
- [15] Electronic Frontier Foundation. SOPA/PIPA: Internet Blacklist Legislation.
- [16] M. Felegyhazi, C. Kreibich, and V. Paxson. On the Potential of Proactive Domain Blacklisting. In *Proceedings of 3rd USENIX LEET*, 2010.
- [17] C. Grier, K. Thomas, V. Paxson, and M. Zhang. @spam: The Underground on 140 Characters or Less. In *Proc. of 17th ACM CCS*, 2010.
- [18] P. H. C. Guerra, D. Guedes, W. M. Jr., C. Hoepers, M. H. P. C. Chaves, and K. Steding-Jessen. Spamming Chains: A New Way of Understanding Spammer Behavior. In *Proc. of 6th CEAS*, 2009.
- [19] Jillian York. Australia Heads Down the Slippery Slope, Authorizes ISPs to Filter. <https://www.eff.org/deeplinks/2011/06/australia-heads-down-slippery-slope-authorizes>, 2011.
- [20] J. P. John, F. Yu, Y. Xie, A. Krishnamurthy, and M. Abadi. deSEO: Combating Search-result Poisoning. In *Proceedings of the 20th USENIX conference on Security*, 2011.
- [21] J. Jung and E. Sit. An Empirical Study of Spam Traffic and the Use of DNS Black Lists. In *Internet Measurement Conference*, Taormina, Italy, 2004.
- [22] C. Kanich, N. Weaver, D. McCoy, T. Halvorson, C. Kreibich, K. Levchenko, V. Paxson, G. M. Voelker, and S. Savage. Show Me the Money: Characterizing Spam-advertised Revenue. In *Proceedings of the USENIX Security Symposium*, San Francisco, CA, Aug. 2011.
- [23] B. Krebs. SpamIt, Glavmed Pharmacy Networks Exposed. Krebs on Security Blog, <http://www.krebsonsecurity.com/category/pharma-wars/>, 2011.
- [24] N. Leontiadis, T. Moore, and N. Christin. Measuring and Analyzing Search-Redirection Attacks in the Illicit Online Prescription Drug Trade. In *Proc. 20th USENIX Security*, 2011.
- [25] K. Levchenko, A. Pitsillidis, N. Chachra, B. Enright, M. Félegyházi, C. Grier, T. Halvorson, C. Kanich, C. Kreibich, H. Liu, D. McCoy, N. Weaver, V. Paxson, G. M. Voelker, and S. Savage. Click Trajectories: End-to-End Analysis of the Spam Value Chain. In *Proceedings of the IEEE Symposium and Security and Privacy*, pages 431–446, Oakland, CA, May 2011.
- [26] H. Liu, K. Levchenko, M. Félegyházi, C. Kreibich, G. Maier, G. M. Voelker, and S. Savage. On the Effects of Registrar-level Intervention. In *Proceedings of the USENIX Workshop on Large-scale Exploits and Emergent Threats (LEET)*, pages 1–8, Boston, MA, Mar. 2011.
- [27] G. Lowe, P. Winters, and M. L. Marcus. The Great DNS Wall of China. <http://cs.nyu.edu/Epcw216/work/nds/final.pdf>, 2007.
- [28] C. Ludl, S. Mcallister, E. Kirda, and C. Kruegel. On the Effectiveness of Techniques to Detect Phishing Sites. In *Proceedings of the 4th DIMVA*, 2007.
- [29] J. Ma, L. K. Saul, S. Savage, and G. M. Voelker. Identifying Suspicious URLs: An Application of Large-Scale Online Learning. In *Proc. of 26th ICML*, 2009.
- [30] D. McCoy, H. Dharmdasani, C. Kreibich, G. M. Voelker, and S. Savage. Priceless: The Role of Payments in Abuse-advertised Goods. In *Proceedings of the ACM Conference on Computer and Communications Security*, Raleigh, NC, Oct. 2012.
- [31] D. McCoy, A. Pitsillidis, G. Jordan, N. Weaver, C. Kreibich, B. Krebs, G. M. Voelker, S. Savage, and K. Levchenko. PharmaLeaks: Understanding the Business of Online Pharmaceutical Affiliate Programs. In *Proceedings of the USENIX Security Symposium*, Bellevue, WA, Aug. 2012.
- [32] T. Moore, N. Leontiadis, and N. Christin. Fashion Crimes: Trending-term exploitation on the web. In *Proceedings of the 18th ACM conference on Computer and communications security*, pages 455–466. ACM, 2011.
- [33] Y. Niu, Y.-M. Wang, H. Chen, M. Ma, and F. Hsu. A Quantitative Study of Forum Spamming Using Context-based Analysis. In *Proc. of 14th NDSS*, 2007.
- [34] G. Orizio, P. Schulz, S. Domenighini, L. Caimi, C. Rosati, S. Rubinelli, and U. Gelatti. Cyberdrugs: a Cross-sectional Study of Online Pharmacies Characteristics. *The European Journal of Public Health*, 19:375–377, 2009.
- [35] L. Paulson. Spam hits Instant Messaging. *Computer*, 37(4), 2004.
- [36] A. Pitsillidis, C. Kanich, G. M. Voelker, K. Levchenko, and S. Savage. Taster’s Choice: A Comparative Analysis of Spam Feeds. In *Proceedings of the ACM Internet Measurement Conference*, Boston, MA, Nov. 2012.
- [37] A. Ramachandran and N. Feamster. Understanding the Network-Level Behavior of Spammers. In *Proc. of ACM SIGCOMM*, 2006.
- [38] D. Samosseiko. The Partnerka — What is it, and why should you care? In *Proc. of Virus Bulletin Conference*, 2009.
- [39] S. Sheng, L. F. Cranor, J. Hong, B. Wardman, G. Warner, and C. Zhang. An Empirical Analysis of Phishing Blacklists. In *Proceedings of the 6th CEAS*, 2009.
- [40] Y. Shin, M. Gupta, and S. Myers. The Nuts and Bolts of a Forum Spam Automator. In *Proceedings of the 4th USENIX LEET*, 2011.
- [41] S. Sinha, M. Bailey, and F. Jahanian. Shades of Grey: On the Effectiveness of Reputation-based Blacklists. In *Proc. of 3rd MALWARE*, 2008.
- [42] S. Sinha, M. Bailey, and F. Jahanian. Improving SPAM Blacklisting through Dynamic Thresholding and Speculative Aggregation. In *Proc. of 17th NDSS*, 2010.
- [43] Spamhaus. <http://www.spamhaus.org/dbl/>.
- [44] SURBL. <http://www.surbl.org/>.
- [45] Symantec. Web-Based Malware Distribution Channels: A Look at Traffic Redistribution Systems. <http://www.symantec.com/connect/blogs/web-based-malware-distribution-channels-look-traffic-redistribution-systems>, 2011.

- [46] The Anti-Abuse Project. DNS Blacklists. <http://www.anti-abuse.org/dns-blacklists/>.
- [47] URIBL. <http://www.uribl.com/>.
- [48] Verisign Announces Increase in .com/.net Domain Name Fees. <https://investor.verisign.com/releasedetail.cfm?releaseid=431292>, 2009.
- [49] D. Wang, S. Savage, and G. M. Voelker. Juice: A Longitudinal Study of an SEO Campaign. In *Proceedings of the Network and Distributed System Security Symposium (NDSS)*, San Diego, CA, Feb. 2013.
- [50] Y.-M. Wang, M. Ma, Y. Niu, and H. Chen. Spam Double-Funnel: Connecting Web Spammers with Advertisers. In *Proc. of 16th WWW*, 2007.
- [51] The Wayback Machine. <http://web.archive.org>.
- [52] Y. Xie, F. Yu, K. Achan, E. Gillum, M. Goldszmidt, and T. Wobber. How Dynamic are IP Addresses. In *Proceedings and 2007 SIGCOMM*, 2007.
- [53] Y. Xie, F. Yu, K. Achan, R. Panigrahy, G. Hulten, and I. Osipkov. Spamming Botnets: Signatures and Characteristics. In *Proceedings of ACM SIGCOMM*, 2008.
- [54] J. Zhang, P. Porras, and J. Ullrich. Highly predictive blacklisting. In *Proceedings of the 17th USENIX Security*, pages 107–122, 2008.
- [55] Y. Zhang, S. Egelman, L. Cranor, and J. Hong. Phinding Phish: Evaluating Anti-phishing Tools. In *Proceedings of the 14th NDSS*, 2007.
- [56] Y. Zhao, Y. Xie, F. Yu, Q. Ke, Y. Yu, Y. Chen, and E. Gillum. BotGraph: Large Scale Spamming Botnet Detection. In *Proceedings of the 6th USENIX NSDI*, 2009.