

# Hacking is not random: a case-control study of webserver-compromise risk

Marie Vasek, John Wadleigh, Tyler Moore

**Abstract**—We describe a case-control study to identify risk factors that are associated with higher rates of webserver compromise. We inspect a random sample of around 200 000 webserver and automatically identify attributes hypothesized to affect the susceptibility to compromise, notably content management system (CMS) and webserver type. We then cross-list this information with data on webserver hacked to serve phishing pages or redirect to unlicensed online pharmacies. We find that webserver running WordPress and Joomla are more likely to be hacked than those not running any CMS, and that servers running Apache and Nginx are more likely to be hacked than those running Microsoft IIS. We also identify several WordPress plugins and Joomla extensions that associated with compromise. Furthermore, using a series of logistic regressions, we find that a CMS's market share is positively correlated with website compromise. Surprisingly, we find that webserver running outdated software are less likely to be compromised than those running up-to date software. We present evidence that this is true for core WordPress software (the most popular CMS platform) and many associated plugins. Finally, we examine what happens to webserver following compromise. We find that under 5% of hacked WordPress websites are subsequently updated, but those that do are recompromised about half as often as those that do not update.

**Keywords:** Content-management systems, webserver security, case-control study, cybercrime, security economics

## 1 INTRODUCTION

EACH month many thousands of websites are compromised by criminals and repurposed to host phishing websites, distribute malware, and peddle counterfeit goods. Despite the substantial harm imposed, the number of infected websites has remained stubbornly high. While many agree that the current level of Internet security is unacceptably low, there is no consensus on what countermeasures should be adopted to improve security or where limited resources should be focused. One key reason we are in such a sorry state is that measuring security outcomes (and what factors drive them) is hard. In part, this is because those who fall victim to cybercrime often prefer not to speak out. But it is also because security mechanisms are deployed in the wild, where it can be impossible to design a randomized controlled experiment isolating the effect of a particular countermeasure to evaluate effectiveness.

However, even when controlled experiments are not feasible, other techniques may still be usefully applied. In this paper, we apply a widely-used method from epidemiology, called a *case-control study*, in order to better understand the factors driving webserver insecurity. Working backwards from data on security incidents and a control sample, we can identify risk factors associated with compromise. This in turn can help defenders better allocate scarce defensive resources to do the most good.

We investigate many observable characteristics of webserver that may affect the likelihood of compromise. Chief among them is whether or not they run a content management system (CMS), an application that simplifies the creation of web content. Some of the more popular CMSes, such as Joomla and WordPress, are consistently exploited to give a miscreant control over the webserver. Additional characteristics include the server type (e.g., Apache), the hosting country, and whether or not the webserver has demonstrated savviness in secure administration practices.

We find the following:

- We apply a novel case-control study design to identify risk factors associated with webserver compromise. We use two compromise datasets (phishing and search-redirectation attacks) and construct a control dataset of uncompromised servers, automatically extracting characteristics suspected of affecting the risk of compromise (Section 2).
- We establish a link between running particular content management systems and an elevated risk of compromise (Section 3). We present evidence that compromise risk is linked to both customizability and popularity. We show that CMSes that are both highly customizable and popular (e.g., WordPress and Joomla) face the biggest risk.
- We further explore the link between popularity and compromise and find that the greater the CMS market share, the more at risk webserver running the CMS are (Section 3).
- We further explore the link between customizability and compromise and find that the more WordPress plugins or Joomla extensions a website has, the greater the odds of compromise (Section 3.3).
- We discover that servers running more up-to-date software are at greater risk of compromise (Section 4). We find this to be true for WordPress and many popular WordPress plugins and Joomla extensions. We hypothesize this is because the most up-to-date software is both more customizable and more popular than previous versions.
- We examine what happens to servers following an initial compromise to distribute malware (Section 5). We find that webserver running outdated software are compromised again, or *recompromised*, most often when compromised and are not subsequently updated. Recompromise rates are lowest when webserver are outdated at the time of compromise but brought up-to-date following compromise.

• Darwin Deason Institute for Cyber Security, Southern Methodist University, Dallas, TX. {mvasek,jwadleigh,tylerm}@smu.edu

## 2 METHODOLOGY

We begin by setting out the key research questions in Section 2.1, then outline the case-control study design in Section 2.2. We discuss the data collection and classification approach in Section 2.3. The collected data and analysis scripts are publicly available for replication purposes at [doi:10.7910/DVN/25608](https://doi.org/10.7910/DVN/25608).

### 2.1 Research Questions

We investigate three categories of research questions about factors that may influence webserver compromise: software type, software market share, and webserver hygiene.

Most generally, we hypothesize that there are measurable differences in compromise rates according to the type of software run on webserver.

- H0:** Running a CMS is a positive risk factor<sup>1</sup> for compromise.
- H0b:** (*corollary*) Some CMS types are risk factors for compromise.
- H0c:** (*corollary*) Some CMS add-ons are risk factors for compromise.
- H1:** Some server types are risk factors for compromise.

There are several reasons why servers running CMSes may be compromised more often. First, CMSes simplify configuration by reducing technical barriers, which means that they are often administered by non-experts. This could lead to a greater chance for server misconfiguration. Second, CMS platforms are a form of software monoculture, exhibiting common vulnerabilities in both the underlying code and the default configurations. We also expect some CMS platforms to be more secure than others. For similar reasons, we expect some CMS add-on software, which adds customization to the standard platform, will also increase the likelihood of infection.

We also anticipate that there will be differences in compromise rates based on the type of server software used. This is because there are different amounts of exploitable vulnerabilities present in the underlying code bases. Additionally, some applications (including CMSes) run only or primarily on particular server types, and each application has its own susceptibility to compromise.

Furthermore, we suspect that a key driving force behind the variation in compromise rates across software types is the software’s market share. When more webserver run a particular type of software, they collectively become a more attractive target for miscreants. The cost of crafting new exploits can be amortized over many more infections for more popular software. While many would agree with such logic on software types, we hypothesize that the same logic also applies to different versions of the same software: more popular software versions tend to be targeted more often than less popular ones. We suspect this is true even when the less popular version is more outdated and has more vulnerabilities.

- H2:** CMS market share is a positive risk factor for webserver compromise.
- H2b:** (*corollary*) Outdated software with limited market penetration is a negative risk factor for compromise.
- H2c:** (*corollary*) The number of exploits available for a type of software is a positive risk factor for compromise.

1. In this paper, a *positive* risk factor is actually a *bad* thing, as it indicates greater odds of compromise. By contrast, a negative risk factor indicates lower odds of compromise.

Our final group of hypotheses involve the individual security practices of webserver administrators. We believe that, independent of the software running on a webserver, adopting security best practices that improve server “hygiene” can influence the likelihood of compromise.

- H3:** Actively hiding detailed software version information is a negative risk factor for compromise.
- H4:** Running a webserver on a shared hosting platform is a positive risk factor for compromise.
- H5:** Setting the HTTPONLY cookie, which protects against cross-site scripting attacks, is a negative risk factor for compromise.

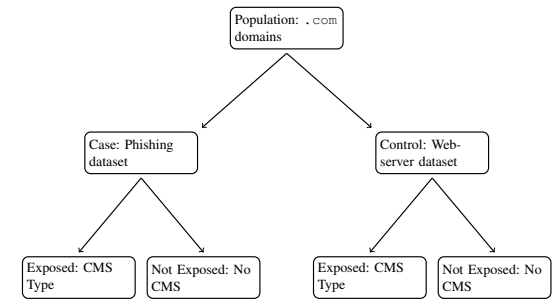
We note that there are other reasons why a webserver could be put at greater risk of being hacked than just the factors discussed above. For example, administrator competence (not captured by the hygiene indicators) certainly plays a role. Security policies also matter: lax password policies or practices could lead to compromise. Finally, the value of the target influences what gets hacked: high-reputation websites, for instance, are targeted for compromise more frequently in search-redirection attacks [1].

We have chosen not to examine the impact of these additional factors in the present study. We decided to focus on CMSes, server software, and webserver hygiene indicators for three reasons. First, as explained above, there is substantial evidence that these factors strongly affect compromise rates (e.g., the large number of exploits available that target CMSes). Second, we have restricted ourselves to factors that could manageably be observed directly and in an automated fashion. By contrast, many of the factors that we chose not to study are not directly observable, such as a company’s password policy. Factors that require extensively crawling or fuzzing a domain to observe, such as inferring firewall policies, are also excluded because they cannot be carried out at sufficient scale. Third, we have restricted ourselves to factors that appear in our sample population with sufficient frequency. In particular, we investigated many of the risk factors from [2] and found the vast majority of them to occur too infrequently to include in our study. It is our view that the methods of analysis presented here could in fact be applied to additional factors, but we defer the task to future work.

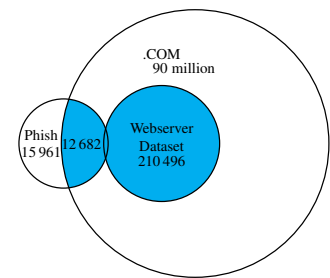
### 2.2 Case-Control Study Design

In a case-control study typically used in epidemiology, data on those afflicted with a disease are compared against as similar a population as possible of those not afflicted [3]. For example, in the seminal case-control study that uncovered the link between smoking and lung cancer, Doll and Arthur surveyed British doctors about their smoking habits, then compared it against data collected subsequently on doctors’ mortality rates [4]. They found that doctors who smoked were much more likely to die than doctors who did not. In general, case-control studies work by comparing two populations, one with a condition (the ‘case’) to one without who are otherwise similar (the ‘control’). Researchers can then work backwards to identify important risk factors by comparing the relative incidence of different characteristics in the case and control populations.

Similarly, we sample a population of webserver and compare them to other populations of webserver that have been compromised. Figure 1a demonstrates the design for the phishing dataset. We start with a comparable webserver population – domains



(a) Case-control study design, demonstrated for phishing dataset and CMS type as risk factor.



(b) Venn diagram demonstrates how we join webserver and phishing datasets.

Fig. 1: We join the webserver and compromise datasets to compare risk factors with outcomes.

registered in `.com`. We then assign the `.com` domains from the phishing dataset as the case and the domains from the webserver dataset as the control. We can then treat characteristics such as CMS type, server type, and hosting country as potential risk factors. (We explain how each of these datasets and risk factors are collected in the next subsection below.) Figure 1b shows a Venn diagram that explains how the phishing and webserver datasets are joined. A similar approach is used for the search-redirection attacks dataset and the webserver dataset.

Note that with case-control data, we do not make any claims about the overall incidence of compromise in the population. This is because we compare two different samples (the compromised and broader samples). Instead, we analyze the prevalence of compromise relative to the occurrence of risk factors such as CMS type.

## 2.3 Data Collection Overview

### 2.3.1 Control Population: Webserver Sample

To answer our research questions, we need a random sample of webserver; however, obtaining a perfectly representative sample of all webserver is not possible since there is no global list available from which to sample. According to Verisign, there are over 252 million registered domains [5], but most zone files listing domains are not made public. Instead, we take a random sample of domains listed in the `.com` zone file. While limited to a single TLD, it is worth noting that `.com` comprises nearly half of all registered domains, and it is used by websites in many countries. Furthermore, `.com` domains include websites from a wide range of popularities. Thus, we feel that sampling from `.com` is broad enough to be representative of all webserver online.

We sampled webserver over a period of 9 days, obtaining information on 210496 domains selected at random from the `.com` zone file downloaded January 15, 2013. We chose this sample size to ensure that it would likely include enough webserver running CMSes with at least 1% market share. This, in turn, improves the chances of obtaining statistically significant results.

We remove all free hosting and URL shortening services (where the URLs are likely set up purposely by the criminals) from our collection. Finally, we refer to the trimmed sample of `.com` domains as the webserver dataset.

### 2.3.2 Case Populations: Compromised Webserver

We consider two sources of data on hacked webserver. First, we examine an amalgamated “feed” of phishing URLs, comprising real-time reports from two firms that remove phishing websites

on behalf of banks, a large brand owner, the crowdsourced list from PhishTank [6], and the Anti-Phishing Working Group’s community feed [7]. We examined 97788 distinct URLs from 29682 domains impersonating 1098 different brands reported between November 20, 2012 and January 7, 2013 in the phishing dataset. According to [8], 94% of domains used for phishing during this period were compromised webserver. Since nearly all of the remainder are highly-ranked sites that we excluded as described below, we are confident that the domains used in our study were compromised to serve phishing pages.

The second dataset on webserver compromise came from webserver observed to be engaging in search-redirection attacks. Here, webserver with high reputation are hacked and reconfigured to surreptitiously channel traffic from search engines to unlicensed pharmacies. Notice that we do not gather information on pharmacies set up by criminals, but rather on domains that were set up by legitimate Internet users and compromised to redirect to criminal webserver. We obtained the dataset gathered by the authors of [1], who updated their system to detect advanced forms of cookie-based redirection as described in [9]. The dataset includes web search results from 218 pharmaceutical-related search terms. Webserver are included in the list if they are observed to redirect to a third-party website and subsequently found to engage in cloaking. The search-redirection attacks dataset includes 58516 distinct URLs gathered between October 20, 2011 and December 27, 2012. These correspond to 10677 unique domains, 6226 of which have a `.com` TLD.

### 2.3.3 Extracting Webserver Risk Factors

The head of an HTML webpage often contains metadata about the webpage in so-called meta tags. One piece of information that many content management system (CMS) authors (and text editors) include is a “generator” tag. This optional tag generally contains the text editor type, content management system, version number and/or any special CMS themes used. For example, a website running WordPress version 3.2.1 might contain the tag `<meta name=‘generator’ content=‘WordPress 3.2.1’>`. We downloaded a copy of the HTML for the top-level webpage on a given domain, and then parsed the HTML to extract the tag.

We then attempted to identify the CMS, if any, along with the version information if included. We used manually crafted regular expressions to complete the task. We focused on the top 13 CMSes with at least 1.0% of CMS market share as of January 2013 according to W<sup>3</sup>Techs [10]. These 13 CMSes collectively comprise 88.4% of all webserver using CMSes. We could identify

CMS type for 9 of the top 13 (84.6% of all CMSes). We also included 3 more CMSes, each with less than 1.0% of market share.

However, we cannot solely rely on generator tags to classify websites by CMS. For instance, most websites running Drupal, one of the most popular CMSes, do not display generator information in their metadata. Consequently, in addition to gathering generator information, we ran a number of regular expressions corresponding to 3 of the 4 most popular CMSes against the dataset. Appendix A compares our custom approach to several off-the-shelf tools for CMS identification.

To identify server software, we collected the packet headers along with the HTML code. In each header was a line specifying the server such as `Server: Microsoft-IIS/7.5`. From this we extracted the server type and version number. We also fetched the IP address of the server and mapped this to the country of origin using MaxMind [11].

We also detected the presence of add-on software commonly used in CMSes, called plugins by WordPress and extensions by Joomla. It has been frequently argued that this third-party add-on software is a natural target for attackers, especially as the core CMS software becomes hardened. We focus on the top fifty most popular WordPress plugins and Joomla extensions (100 in total), as determined by their occurrence in the control dataset.

We identified WordPress plugins by scanning each website's stored HTML files for paths beginning with `/wp-content/plugins/`. The following directory indicates the corresponding plugin, e.g., a website using the WP eCommerce plugin has the `/wp-content/plugins/wp-e-commerce/` path.

We detected Joomla extensions in a similar manner. Extensions are comprised of components, modules, plugins, templates, and languages. We used regular expressions to identify each plugin, such as `/components/com_\w*/` for finding components.

We also tried to find versioning information for WordPress plugins. As there is no standard way to convey version information in plugins, from manual inspection we successfully identified plugin information for 19 of the top 50. Some WordPress plugins broadcast their version in a parameter handed to their scripts. For example, a website running version 6.1 of Google Analyticator would contain `wp-content/plugins/google-analyticator/external-tracking.min.js?ver=6.1`. Unfortunately not all plugins are so transparent with versions, and those that are may not be specifying the plugin's version (instead, a version of the script or the WordPress installation).

### 2.3.4 Reducing False Positives in the Infection Datasets

Not all of the URLs in the compromise datasets are from hacked webpages. For the phishing dataset, we deem any URL to be a *false positive* if the URL does anything other than impersonate another website. For the search-redirection attacks dataset, we classify any URL as a false positive if the destination website following redirection appears related to the source website (e.g., `ilike.com` redirects to `myspace.com`, which bought the company).

Since the false positive rates for phishing are consistently higher than for search-redirection attacks, we developed automated techniques to discard websites that were errantly placed on these lists. We removed all FQDNs that redirected to legitimate US-based banks<sup>2</sup> and other known non-banks frequently

targeted by phishing, such as `paypal.com`, `amazon.com` and `facebook.com`. We also generated a sequence of regular expressions that detected Microsoft Outlook Web Applications and coupon websites and checked them against the HTML we downloaded previously. These initial steps reduced our overall false positive rate for the phishing dataset from 9.4% to 5.0%. To further improve, we manually inspected all URLs in the Alexa top million sites and excluded any false positives from further consideration, yielding final false positive rates of 2.3% for phishing and 4.3% for search-redirection attacks. These false positive rates were calculated by inspecting a stratified random sample by Alexa rank.

## 3 IDENTIFYING RISK FACTORS FOR COMPROMISE

### 3.1 Odds Ratios

Odds are defined by the ratio of the probability that an event will occur to the probability it will not occur. For example, if  $p = 0.2$ , then the odds are  $\frac{p}{1-p} = \frac{0.2}{0.8} = 0.25$ . Odds express relative probabilities. *Odds ratios* compare the odds of two events, each occurring with different probabilities.

In case-control studies, odds ratios compare the odds of a subject in the case population exhibiting a risk factor to the odds of a subject in the control population exhibiting a risk factor. Consider the four cases:

	Case (afflicted)	Control (not afflicted)
Has risk factor	$p_{\text{CaseRF}}$	$p_{\text{CtrlRF}}$
No risk factor	$p_{\text{Case}\overline{\text{RF}}}$	$p_{\text{Ctrl}\overline{\text{RF}}}$

The odds ratio, then, is the following product of probabilities:

$$\text{odds ratio (OR)} = \frac{p_{\text{CaseRF}}/p_{\text{Case}\overline{\text{RF}}}}{p_{\text{CtrlRF}}/p_{\text{Ctrl}\overline{\text{RF}}}} = \frac{p_{\text{CaseRF}} * p_{\text{Ctrl}\overline{\text{RF}}}}{p_{\text{Case}\overline{\text{RF}}} * p_{\text{CtrlRF}}}$$

An odds ratio of 1 means that there is no difference in proportions of the risk factor among the case and control groups. An odds ratio greater than 1 indicates that those in the case group are more likely to exhibit the risk factor (so-called *positive* risk factors). By contrast, an odds ratio less than 1 indicates that those in the case group are less likely to exhibit the risk factor (indicating a *negative* risk factor).

### Results

Table 1 reports odds ratios for different CMS and server types for both compromise datasets. We computed odds ratios for web-servers running each of the major CMSes compared to web-servers not running any CMS. We denoted statistically significant positive risk factors in red and statistically significant negative risk factors in green. We can interpret these results as web-servers running WordPress were more than four times as likely to be compromised to serve phishing pages than a webserver running no CMS, but a webserver running Blogger was only about a third as likely to be compromised to serve phishing pages than a webserver running no CMS. For the phishing dataset, some less popular CMSes fare better than not using a CMS, but the more popular CMSes are positive risk factors. WordPress, Joomla and Zen Cart had increased odds of compromise, while Blogger, TYPO3 and Homestead reduced risk.

This supports hypothesis **H0b**, but partially refutes hypothesis **H0** that using *any* CMS increases the odds of compromise. For search-redirection attacks, CMSes are either as bad or worse than not using a CMS, supporting **H0**. Notably, the odds ratios for Joomla and WordPress are even higher than for phishing. The

2. Found on the FDIC website [12].

	Content Management System (CMS) Type									
	Risk factor	Odds ratio	Phishing dataset		Risk factor	Odds ratio	Search-redirection attacks dataset			
			95% CI	# Phish	# Not phish			95% CI	# Redir.	# Not redir.
No CMS		1.00		8 747	190 305		1.00		2 260	190 314
WordPress	+	<b>4.44</b>	(4.24, 4.65)	2 673	13 101	+	<b>17.18</b>	(16.20, 18.22)	2 674	13 106
Joomla	+	<b>7.11</b>	(6.62, 7.63)	1 106	3 384	+	<b>23.96</b>	(22.05, 26.04)	963	3 385
Drupal		0.79	(0.58, 1.04)	46	1 279	+	<b>6.59</b>	(5.33, 8.07)	100	1 279
Zen Cart	+	<b>4.84</b>	(3.26, 6.96)	33	149		2.35	(0.71, 5.56)	4	149
Blogger	-	<b>0.28</b>	(0.13, 0.52)	8	637		1.08	(0.49, 2.02)	8	637
TYPO3	-	<b>0.14</b>	(0.03, 0.37)	3	481	+	<b>4.23</b>	(2.72, 6.24)	24	481
Homestead	-	<b>0.04</b>	(0.00, 0.18)	1	607	-	<b>0.16</b>	(0.01, 0.69)	1	607

	Server Type									
	Risk factor	Odds ratio	Phishing dataset		Risk factor	Odds ratio	Search-redirection attacks dataset			
			95% CI	# Phish	# Not phish			95% CI	# Redir.	# Not redir.
Microsoft IIS		1.00		1 002	60 495		1.00		193	60 497
Apache	+	<b>5.44</b>	(5.10, 5.81)	10 549	117 017	+	<b>14.12</b>	(12.26, 16.36)	5 276	117 031
Nginx	+	<b>2.24</b>	(2.01, 2.50)	507	13 649	+	<b>8.63</b>	(7.26, 10.30)	376	13 649
Yahoo	-	<b>0.62</b>	(0.41, 0.89)	27	2 634		1.57	(0.85, 2.64)	13	2 634
Google		0.63	(0.35, 1.03)	14	1 359		1.88	(0.84, 3.57)	8	1 359

TABLE 1: Odds ratios for varying CMS and server types. Statistically significant results are listed in bold. Odds ratios less than one are negative risk factors, while odds ratios greater than one are positive risk factors for compromise.

WordPress odds ratio jumps from 4.4 phishing to 17 for search-redirection attacks; for Joomla, the jump is from 7 to nearly 24!

For some smaller CMSes, the evidence for phishing and search-redirection attacks is mixed. Homestead has a negative risk factor for phishing and search-redirection attacks dataset. TYPO3 and Blogger are negative for phishing, but TYPO3 has a positive risk factor for search-redirection attacks, whereas Blogger is not statistically significant.

We note that the larger CMSes tend to be the strongest positive risk factors for compromise, according to both datasets. This supports hypothesis **H2** that CMS market share is positively correlated with compromise, but more analysis is needed.

For server software type, we compute risk factors relative to Microsoft IIS, the second-most popular server software. Apache and Nginx are positive for both phishing and search-redirection attacks. Note that we are not making any claims about the relative security levels of the different software classes. All software contains vulnerabilities, and we are not taking sides on the debate over whether open- or closed-source software has fewer unpatched holes [13]. Instead, our results simply show that, relative to software popularity, criminals tend to use Apache and Nginx more for perpetrating their crimes than Microsoft IIS.

### 3.2 Logistic Regressions

We now present logistic regressions to study why websites are compromised. We run four regressions in all: two for webserver running a CMS (one each for the phishing and search-redirection attacks datasets) and two for webserver not running any CMS (one for each compromise dataset). We run the additional regressions because some explanatory variables only apply to CMSes, but many of the variables measuring security signals apply regardless of whether or not a webserver uses a CMS.

We group the following explanatory variables into three categories: CMS market share, webserver hygiene and server attributes.

#### CMS Market Share

**# Servers:** We took market share for each CMS from [10] as of January 1, 2013 and multiplied it by population of registered .com domains (106.2 million) and estimated server response rate (85%) [5]. This variable was omitted for non-CMS regressions.

#### Webserver Hygiene

**HTTPONLY cookie:** We checked the header for an HTTPONLY

cookie used to protect against cross-site-scripting attacks. We interpret setting this cookie as a positive signal of overall server hygiene. Checking for this cookie was one measure of server hygiene also used in [2].

**Server Version Visible:** We analyzed the server headers for any version information regarding the server, whether it be Apache 2 or Apache 2.2.22. This is a Boolean variable which is true if the server gave any potentially valid version information.

**Shared Hosting:** We counted the number of times we observed an IP address in the combined webserver and compromised datasets. We deem a domain to be part of a shared host if 10 domains resolve to the same IP address. A recent Anti-Phishing Working Group report presents evidence that some attackers target shared hosting in order to simultaneously infect many domains [8].

#### Server Attributes

**Country:** We took the top ten countries from the combined dataset and compared each of them the domains hosted in all the other countries in the dataset.

**Server Type:** This categorical variable looks at the type of server software a webserver is running. We only consider the 5 most popular types: Apache, Microsoft IIS, Nginx, Google, and Yahoo.

The model takes the following form:

$$\log \frac{p_{comp}}{1 - p_{comp}} = c_0 + c_1 \lg(\# \text{ Servers}) + c_2 \text{ HTTPONLY} \\ + c_3 \text{ Server Vsn} + c_4 \text{ Shared Hosting} \\ + c_5 \text{ Country} + c_6 \text{ Server type} + \varepsilon$$

Table 2 shows the results from these four regressions. CMS popularity is positively correlated with compromise in the phishing dataset. Each doubling of the number of webserver running the CMS increases the odds of compromise by 9%, supporting hypothesis **H2**. The result is inconclusive for search-redirection attacks, but the trend is similar. Also, Appendix B studies the link between market share and exploitability. The analysis in Appendix B shows that the number of exploits is also a positive risk factor for being hacked to serve phishing pages, which supports **H2c**.

We consider hygiene variables next. We do not observe any consistent evidence that hiding server information promotes or inhibits compromise, so we can neither refute nor support **H3**. Setting an HTTPONLY cookie appears to be a negative risk factor

	CMS						No CMS					
	Phish			Search-redirect attacks			Phish			Search-redirect attacks		
Intercept	coef.	odds	p-value	coef.	odds	p-value	coef.	odds	p-value	coef.	odds	p-value
lg # Svrs	0.09	<b>1.09</b>	< 0.0001	0.02	1.02	0.16	-4.11	<b>0.02</b>	< 0.0001	-5.99	<b>0.00</b>	< 0.0001
HTTPONLY	0.22	1.25	0.06	-0.83	<b>0.44</b>	< 0.0001	-0.87	<b>0.42</b>	< 0.0001	0.15	1.17	0.12
No Svr Vsn	-0.15	<b>0.86</b>	0.0001	0.10	<b>1.11</b>	0.01	0.04	1.04	0.09	0.32	<b>1.38</b>	< 0.0001
Shared Host	0.95	<b>2.58</b>	< 0.0001	-1.58	<b>0.21</b>	< 0.0001	0.28	<b>1.32</b>	< 0.0001	-1.27	<b>0.28</b>	< 0.0001
Apache	1.49	<b>4.45</b>	< 0.0001	1.48	<b>4.38</b>	< 0.0001	1.80	<b>6.06</b>	< 0.0001	1.37	<b>3.94</b>	< 0.0001
Nginx	0.59	<b>1.80</b>	0.003	1.37	<b>3.93</b>	< 0.0001	0.70	<b>2.00</b>	< 0.0001	1.43	<b>4.19</b>	< 0.0001
Yahoo	-0.34	0.72	0.59	2.72	<b>15.12</b>	< 0.0001	-0.54	<b>0.58</b>	0.009	-0.02	0.98	0.97
Google	-1.50	<b>0.22</b>	0.0003	-0.81	0.44	0.10	-0.36	0.70	0.35	0.25	1.29	0.67
Other	1.92	<b>6.84</b>	< 0.0001	0.83	<b>2.30</b>	0.0009	0.81	<b>2.24</b>	< 0.0001	0.96	<b>2.62</b>	< 0.0001
Model fit:	$\chi^2 = 1353, p < 0.0001$			$\chi^2 = 1825, p < 0.0001$			$\chi^2 = 5937, p < 0.0001$			$\chi^2 = 2113, p < 0.0001$		

TABLE 2: Table of coefficients for logistic regressions comparing rate of compromise to many explanatory variables. Statistically significant odds ratios are listed in bold. Odds ratios less than one are negative risk factors, while odds ratios greater than one are positive risk factors for compromise.

for being compromised, but we need more data to support the associated hypothesis **H5**.

Running on a shared host is a positive risk factor for being hacked to serve phishing pages, which supports **H4** and findings from [8]. However, we note that it is a negative risk factor for being hacked for search-redirect attacks. It appears that cyber-criminals engaged in phishing have adopted different techniques for infecting web servers than those carrying out search-redirect attacks. Further investigation shows that there is a correlation between being on a shared host and having a low or no Alexa rank: 13% of the top 10M, 26% of the next 10M, and 55% of websites without an Alexa rank are hosted on a shared host (from our combined webserver and search-redirect attacks dataset). This result could signal that search-redirect attacks attackers target higher ranked pages, which makes sense in light of [1], which showed that compromised websites with a higher PageRank stay in search results longer.

Previous results from web servers in Section 3.1 are similar to those in this regression – notably that Apache and Nginx web servers remain positive risk factors compared to Microsoft IIS in all cases.

Finally, we note that there is more consistency between the regressions examining CMSes and no CMSes than there is between regressions for phishing and search-redirect attacks. The results for the shared host variable are the same, regardless of whether a CMS is used, as are the results for server types and most countries. Only the practice of hiding detailed server version information was very inconsistent, being a negative risk factor for phishing on CMSes and a negative risk factor for search-redirect attacks when no CMS is used.

### 3.3 The Role of CMS Add-on Software in Compromise

Even if all the exploitable holes in core CMS software are patched, a large weakness in the form add-on software remains. As described in Section 2, we identified the 100 most popular WordPress plugins and Joomla extensions.

We first examine whether running one of these popular plugins increases the odds for compromise. 57.8% of web servers running WordPress and 60.1% running Joomla run at least one top 100 add-on. Indeed, WordPress servers running a top-50 plugin are at 21.9% greater odds of compromise, and Joomla servers running a top-50 extension are at 54.3% greater odds of compromise. Running a popular add-on software, regardless of what it is, is a positive risk factor for compromise.

We next take a closer look at how individual add-ons can affect the odds of compromise. Table 3 reports odds ratios for WordPress plugins and Joomla extensions found to be statistically significant. Here we compare the odds of websites running the CMS that also run the plugin against those running the CMS but not the plugin under study. We find 15 of the top-50 WordPress plugins are positive risk factors for compromise. MM Forms Community leads the way, with a 26-times greater odds of compromise compared to running WordPress without the plugin. This plugin is rife with problems: there was a SQL Injection vulnerability reported in 2011<sup>3</sup>, an arbitrary file upload vulnerability reported in 2012<sup>4</sup>, and was last under active development in 2012. Notably, 4 top-50 WordPress plugins are negative risk factors for compromise, including the popular script TimThumb, which had been subject to a high-profile vulnerability in August 2011 [14] and again in June 2014<sup>5</sup>. Meanwhile, 17 top-50 Joomla extensions are positive risk factors for compromise, and just two are negative risk factors for compromise.

Figure 2 plots odds ratio for compromise as a function of the total number of plugins and extensions observed on a given website. The odds are compared against web servers not running any add-on software. For both WordPress and Joomla, each additional add-on creates an increasing risk for compromise. For WordPress, websites running two plugins have a 1.6 times greater odds of being compromised than those not running any plugins. This rises to 2 times greater odds for websites running 10 or more plugins. Adding new Joomla extensions, by contrast, has a more clearly additive impact on compromise risk. Web servers running 3 Joomla extensions face 1.86 times greater odds of compromise compared to those running Joomla without extensions. The odds increase steadily as more extensions are installed, rising to 4.8 times greater odds of compromise for web servers running at least 10 Joomla extensions. Finally, we ran a logistic regression to measure the impact of software add-ons controlling for CMS version, server software version and hosting infrastructure. This quantifies the incremental risk of each additional software add-on. The results indicate that, all else equal, each new add-on increases the odds of compromise by 11%.

3. <http://www.exploit-db.com/exploits/17725/>

4. <http://wpsecure.net/2012/06/mm-forms-community/>

5. <http://seclists.org/fulldisclosure/2014/Jun/117>

WordPress Plugin	Odds	95% CI	Joomla Extension	Odds	95% CI
MM Forms Community	<b>25.99</b>	(5.09, 634.31)	JomComment <sup>C</sup>	<b>7.80</b>	(5.27, 11.94)
Dynamic Content Gallery	<b>7.07</b>	(5.47, 9.23)	Autson Slide Show <sup>M</sup>	<b>2.22</b>	(1.36, 3.68)
Audio Player	<b>2.23</b>	(1.80, 2.76)	RokStories <sup>C</sup>	<b>2.17</b>	(1.53, 3.08)
WPaudio MP3 Player	<b>1.87</b>	(1.29, 2.69)	Social Media Links <sup>M</sup>	<b>2.04</b>	(1.28, 3.27)
Easing Slider	<b>1.85</b>	(1.22, 2.79)	Frontpage SlideShow <sup>C</sup>	<b>1.94</b>	(1.31, 2.87)
WordPress Popular Posts	<b>1.72</b>	(1.24, 2.36)	JComments <sup>C</sup>	<b>1.92</b>	(1.41, 2.61)
WP-Polls	<b>1.70</b>	(1.37, 2.10)	RokAjaxSearch <sup>C</sup>	<b>1.86</b>	(1.42, 2.43)
Digg Digg	<b>1.63</b>	(1.21, 2.17)	JA Tabs	<b>1.84</b>	(1.22, 2.77)
WP-reCAPTCHA	<b>1.52</b>	(1.11, 2.07)	News Show Pro GK4 <sup>M</sup>	<b>1.72</b>	(1.22, 2.42)
WP-PostRatings	<b>1.50</b>	(1.07, 2.10)	Frontpage SlideShow <sup>M</sup>	<b>1.64</b>	(1.17, 2.30)
MailChimp	<b>1.40</b>	(1.05, 1.86)	AVReloaded	<b>1.64</b>	(1.24, 2.16)
Viper's Video Quicktags	<b>1.39</b>	(1.09, 1.76)	Vinaora Visitors Counter <sup>M</sup>	<b>1.58</b>	(1.17, 2.14)
Sociable	<b>1.30</b>	(1.06, 1.60)	YOOsearch <sup>M</sup>	<b>1.56</b>	(1.01, 2.42)
Jetpack	<b>1.28</b>	(1.18, 1.45)	K2 <sup>C</sup>	<b>1.54</b>	(1.28, 1.85)
Google Analyticator	<b>1.20</b>	(1.03, 1.38)	RokBox <sup>C</sup>	<b>1.41</b>	(1.19, 1.68)
TimThumb	<b>0.81</b>	(0.68, 0.95)	YOOeffects	<b>1.37</b>	(1.00, 1.85)
Custom Contact Forms	<b>0.63</b>	(0.44, 0.88)	MTUpgrade	<b>1.31</b>	(1.04, 1.65)
Gravity Forms	<b>0.63</b>	(0.39, 0.97)	Joom!Fish <sup>C</sup>	<b>0.56</b>	(0.41, 0.74)
IE SiteMode	<b>0.08</b>	(0.04, 0.12)	Languages <sup>M</sup>	<b>0.42</b>	(0.23, 0.74)

TABLE 3: Odds ratios for varying plugin types (all statistically significant).

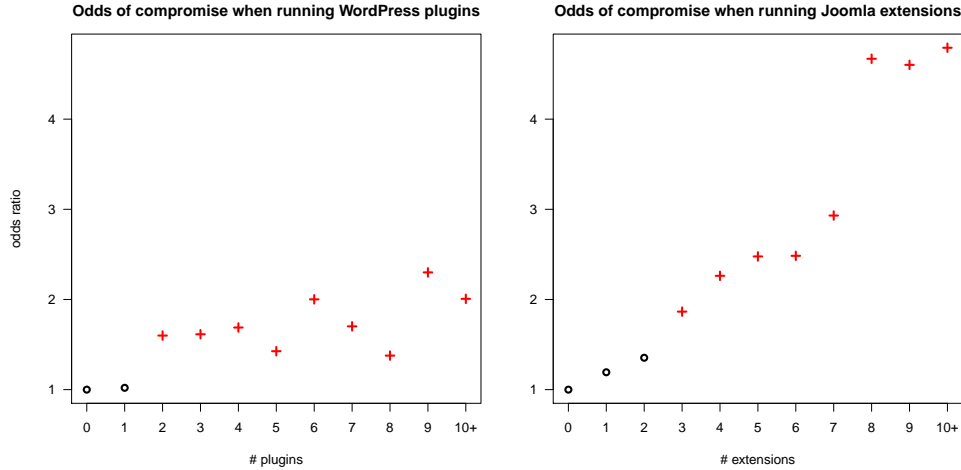


Fig. 2: Odds of compromise based on the number of top-50 WordPress plugins (left) and top-50 Joomla extensions (right). Statistically significant positive risk factors are indicated by red plus signs.

## 4 DOES OUTDATED SOFTWARE GET HACKED MORE OR LESS OFTEN?

A best practice for webserver security is to run the most recent version of software available, as updates tends to plug security holes as well as add new features. For instance, Google notifies webmasters via its Webmaster Tools when it detects outdated server software as a way to improve security [15]. However, updating server software can be a nuisance, due to cross-dependencies, poor interfaces and the demands of maintaining uptime. Consequently, many webserver run software that is many months, or even years, out of date. The security firm Sucuri Labs even runs a website [16] that names and shames websites running woefully outdated CMS or server software.

But we wondered whether or not servers running outdated software actually do get compromised more often than those that do not. We hypothesize that the opposite is usually true: that outdated webserver are compromised less often provided that most other webserver are already upgraded. To test this and related hypotheses, we first restrict ourselves to the servers running WordPress. This is for two reasons: WordPress is the most popular content management system and, by default, WordPress installs provide detailed version information ordered straightforwardly.

### 4.1 Comparing Compromise Across Major WordPress Subversions

#### Odds Ratios for Major Version Differences in WordPress

First, we investigated whether servers running WordPress that hid version information were at less risk of compromise (to test hypothesis **H3**). The results are shown in the first row of the table in Figure 3c. In fact, hiding WordPress version is a positive risk factor for being hacked for phishing pages. This contradicts the frequently held view that hiding detailed version information improves security, and it instead lends credence to the view that publishing information helps defenders more than attackers. For instance, WordPress and Google send out reminder emails to server administrators to update their software, but those who obscured their generator version for security reasons do not receive the reminders. We also note that even though we looked at version information through the generator tag, attackers oftentimes try their hack on any server running WordPress, regardless of what version it says it is. We see no statistically significant effect for search-redirection attacks, though the trend is similar.

There are differing degrees of outdated software. For servers with version information, we first compared the risk facing servers at the most recent version (3.5.1 during our collection time) to running any other version of WordPress. Running the most up-to-

date version is a positive risk factor for being hacked for search-redirectation attacks. This too goes against conventional wisdom, and indirectly supports hypothesis **H2** since the most recent version is also the most popular one.

We also looked at the difference in major versions, ignoring version 1 since we only had 7 instances in our combined datasets. We compared all of WordPress 2.\* and WordPress 3.\* against WordPress installs with no version information. We see that WordPress 3.\* installs face more risk of being hacked to serve phishing pages than WordPress 2.\*. We observe similar but statistically insignificant results for search-redirectation attacks.

### Chi-squared Test for Risk Across Subversions

The odds ratios just discussed offer initial evidence that being out of date reduces the risk of infection for webserver running WordPress, at least when comparing major versions. We now drill down and investigate differences across WordPress subversions (e.g., WordPress 3.3.\*). Figure 3a plots the relative frequency of servers in our webserver and compromise datasets running each WordPress subversion. Note the different scales to the vertical axes – the left axis tracks the frequency in the webserver dataset while the right axis is used for the two compromise datasets. We first observe that more outdated subversions are indeed less popular compared to the most recent subversions. We also see that the compromise rate roughly follows the popularity of the subversion, but with substantial variation and lower compromise rates for more outdated versions.

But are the differences in compromise rates statistically significant? We can answer that using a  $\chi^2$  test, but first, we can inspect the differences visually using the mosaic plot in Figure 3b. The vertical axis shows for each version the proportion of compromised webserver (either phishing or search-redirectation attacks) compared to the proportion of uncompromised webserver (from the webserver dataset). The horizontal axis is scaled so that the area of each cell matches the frequency of each category. For instance, the dark blue cell in the bottom right corner shows the proportion of webserver running WordPress Version 3.5.\* that have been compromised. This plot shows that the fraction compromised falls steadily as the subversions grow more outdated. It also shows that the collective proportion of outdated servers is still quite substantial.

Finally, the cells are lightly shaded if the difference in proportion for being compromised is statistically significant at the 95% confidence interval according to the  $\chi^2$  test, and over 99% confidence interval if darkly shaded. Red cells are underrepresented and blue cells are overrepresented. We can see that most of the WordPress 2.\* versions are statistically overrepresented in the webserver dataset and underrepresented in the compromise datasets. WordPress 3.0 and 3.3 are also overrepresented in the compromise datasets and underrepresented in the webserver dataset. The most recent, WordPress 3.5, is the only subversion overrepresented in the phish dataset and underrepresented in the webserver dataset. These findings support hypothesis **H2b** that unpopular outdated CMSes are negative risk factors for compromise. It is also consistent with our findings from the odds ratios that the most recent version is the most at risk of compromise.

### Logistic Regressions

The final check we make comparing compromise rates in WordPress versions is to run a simple logistic regression comparing the popularity of a version to the compromise rate in the phishing dataset.

WordPress plugin	% up-to-date compromised	% out-of-date compromised	%-pts. difference for up-to-date	Odds ratio
WP-Table Reloaded	48.28	24.71	23.57	<b>2.83</b>
The Events Calendar	48.84	28.30	20.54	<b>2.39</b>
WP eCommerce	40.43	22.70	17.73	<b>2.30</b>
WP jQuery Lightbox	37.14	21.74	15.40	2.07
Theme My Login	37.93	25.00	12.93	1.82
Contact Form 7	33.91	24.47	9.44	<b>1.58</b>
Google Analyticsator	38.26	29.03	9.23	<b>1.51</b>
WP-Polls	43.72	36.88	6.84	1.33
MailChimp	42.12	35.79	6.32	1.31
Audio Player	47.77	41.94	5.84	1.26
Easing Slider	46.67	41.27	5.40	1.24
Lightbox Plus Colorbox	33.33	28.96	4.37	1.30
Digg Digg	40.52	36.84	3.68	1.16
WPaudio MP3 Player	43.43	42.11	1.33	1.05
NextGEN Gallery	28.57	30.59	-2.06	0.95
Gravity Forms	17.65	22.58	-4.93	0.74
WooCommerce	23.68	28.81	-5.13	0.77
cforms	25.00	31.33	-6.33	0.80
WP-Paginate	29.70	39.13	-9.43	0.66

TABLE 4: Comparing compromise rates for webserver running up-to-date versus outdated WordPress plugins (statistically significant odds ratios in bold).

**# Servers:** We took the market share for each WordPress subversion from [10] as of January 1, 2013 and multiplied it by population of registered .COM domains (106.2 million) and the estimated server response rate (85%) from [5].

$$\log \frac{p_{comp}}{1 - p_{comp}} = c_0 + c_1 \lg(\# \text{ Servers}) + \epsilon.$$

The logistic regression yields the following results:

	coef.	Odds Ratio	95% conf. int.	Significance
Intercept	-5.60	<b>0.00</b>	(0.00, 0.01)	$p < 0.0001$
$\lg(\# \text{ Servers})$	0.19	<b>1.20</b>	(1.17, 1.24)	$p < 0.0001$
Model fit:	$\chi^2 = 200.31, p < 0.0001$			

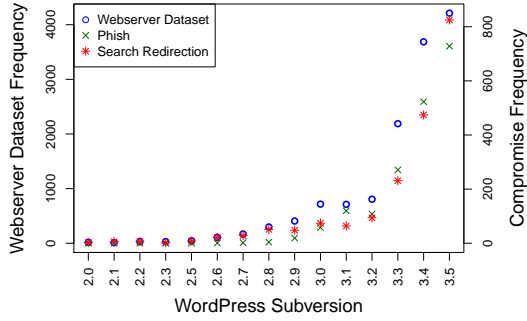
These results show that each time the number of servers running the same subversion of WordPress doubles, the risk of the server being hacked to serve phishing pages increases by 20%. This offers further evidence supporting **H2**.

## 4.2 Comparing Compromise Across Plugin Versions

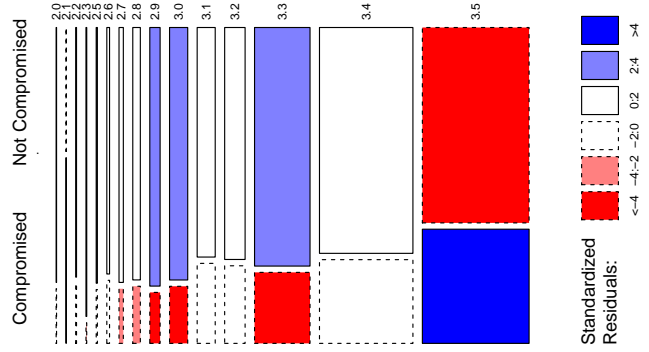
As described in Section 2, we were able to automatically extract detailed version information for 19 of the top-50 WordPress plugins. We now examine whether we see a similar pattern as already found for WordPress, namely that more up-to-date software is in fact hacked more often. For each plugin, table 4 reports the percentage of websites running an up-to-date version of the plugin that are compromised, compared against the percentage running outdated software that are compromised.

To be consistent with the findings just presented for WordPress, we would expect to see a greater incidence of compromise among up-to-date plugins than outdated ones. In fact, this is usually what happens. For 14 of the 19 plugins, the up-to-date version is compromised more often than servers running outdated versions of the plugin. In five cases, the odds ratio is statistically significant, indicating that the higher rate of compromise cannot be attributed to chance. Notably, in all cases with statistically significant odds ratios, running up-to-date plugin software is a positive risk factor for compromise. This lends further empirical support to the surprising finding that running up-to-date software may actually put webserver at greater risk of being hacked.





(a) Incidence of compromise by WordPress version, along with the popularity of WordPress version.



(b) Mosaic plot of WordPress version popularity and incidence of compromise (red cells indicate statistically significant underrepresentation, blue cells overrepresentation).

	Risk factor	Odds ratio	95% CI	Phishing dataset # Phish	# Not phish	Risk factor	Odds ratio	95% CI	Search-redirection attacks dataset # Redir.	# Not redir.
Version Found		1.00		1 834	9 676		1.00		1 936	9 680
No Version	+	<b>1.29</b>	(1.18, 1.41)	839	3 425		1.08	(0.98, 1.18)	738	3 426
Other WordPress versions		1.00		1 606	8 599		1.00		1 440	8 601
WordPress 3.5.1		1.13	(0.97, 1.32)	228	1 077	+	<b>2.75</b>	(2.43, 3.09)	496	1 079
No Version		1.00		839	3 425		1.00		738	3 426
WordPress 2.*	-	<b>0.12</b>	(0.08, 0.17)	26	918		0.88	(0.73, 1.05)	173	918
WordPress 3.*	-	<b>0.84</b>	(0.77, 0.92)	1 808	8 751		0.93	(0.85, 1.03)	1 762	8 755

(c) Odds ratios by WordPress versioning.

Fig. 3: Exploring the relationship between WordPress version and the incidence of webserver compromise.

## 5 WHAT HAPPENS TO CMS SOFTWARE FOLLOWING A COMPROMISE?

Up to now, our analysis has focused on identifying which web-servers are more likely to be compromised. We have also collected additional data to find out what happens *after* a compromise takes place: what actions webmasters take and how that affects recompromise rates. Cleaning up after an infection can be tricky. Because many webmasters are unaware of how the attackers got in [17], they do not always fix the hole that the attacker used to get in the first time. We hypothesize that the steps a webmaster takes after being hacked can greatly influence whether or not the website is hacked later on.

To that end, we set up an ongoing collection using malware URL data from StopBadware’s Data Sharing Project. This includes data from ESET, Fortinet, Internet Identity, Malware Must Die, Sophos, among others. We access each FQDN associated with a malware URL within an hour of the URL being submitted to the feed. We then follow up 1, 2, 5, 9, and 15 days after the initial access to check whether or not the CMS software gets updated. We accessed 464 737 FQDNs from unique domains from February 2014 through October 2014. Of these, 50 179 ran WordPress, the sole CMS we investigate here due to its popularity and the ease of tracking versions. We whittle this down to 44 712 of those which had easily discernible versioning

Of the WordPress sites, 88.7% had the same WordPress version throughout the observational period; 21.6% of these were up to date on the last day of measurement. Of the WordPress sites that updated their CMS, 56.4% did so to stay up to date with a new WordPress version being released. The other 43.6% who were out of date on the first day of measurement took 5 days, on average, to update their website an average of 7 subversions.

	% of hacked websites	% recomp.	Odds ratio for recompromise
Up-to date when compromised	29.1%	22.6%	1
Never updated	66.2%	33.5%	1.72 (1.63, 1.81)
Updated after compromise	4.7%	17.0%	0.70 (0.62, 0.79)

TABLE 5: Observed updating behavior WordPress websites hacked with drive-by-downloads, and the effect on recompromise rates.

Table 5 examines how often WordPress gets updated after a breach. 29.1% of WordPress websites were up-to-date when they were hacked. Of the remaining sites, 4.7% updated after a compromise and 66% continued to run the outdated software. But what impact does the decision whether or not to update have on recompromise? We can see that 22.6% of already updated WordPress websites are recompromised, compared to 33.5% of sites that never updated. Surprisingly, the group that fares best of all are those websites that are updated only after a compromise takes place: they are recompromised just 17% of the time. Thus, the outdated websites that do get hacked would be wise to update to the latest version, likely plugging the hole that the attacker used to get in the first time.

## 6 DISCUSSION

We now sum up the results of the prior sections by first revisiting the original hypotheses and second discussing how the results can be leveraged by security engineers.

**Evaluating Research Questions** We summarize the analysis of the previous section by returning to the original research questions.

- H0** (*Running a CMS pos. RF*) Supported for search-redirect attacks, not uniformly for phishing
- H0b** (*Some CMS types are RFs*) Broadly supported
- H0c:** (*Some CMS add-ons are RFs*) Broadly supported.
- H1** (*Some server types are RFs*) Broadly supported
- H2** (*CMS market share pos. RF*) Broadly supported, across all CMSes and across WordPress subversions
- H2b** (*Outdated unpopular software neg. RF*) Supported across WordPress subversions
- H2c** (*# exploits pos. RF*) Supported
- H3** (*Hiding version info neg. RF*) Contradicted
- H4** (*Shared hosting pos. RF*) Supported for phishing, contradicted for search-redirect attacks
- H5** (*HTTPONLY cookie pos. RF*) Inconclusive

Many hypotheses are broadly supported, especially that server type and CMS market share are positive risk factors. We find less support for hypothesis **H0** that *all* CMSes exhibit higher rates of compromise; instead, *most* CMSes, especially the popular ones, are positive risk factors for compromise. Similarly, *many* CMS add-on software components are also positive risk factors for compromise. Finally, it does not appear that hiding version information is a negative risk factor in most circumstances, but it is unclear how often it may be a positive risk factor.

**Making the Results Actionable** So what can be made of these results? At a high level, the findings can help reduce information asymmetries regarding security outcomes for different webserver configurations [18]. By making security outcomes such as compromise incidents more directly comparable across of platforms, we can help others make more informed decisions about the relative risks posed. Publishing such data can also motivate software developers to improve the security of their code.

We have seen, however, that not all “name-and-shame” policies are consistent with empirical observation. Notably, efforts to call out websites running outdated software are misguided, since they obscure our finding that up-to-date servers tends to be hacked more often. Instead, relative metrics such as odds ratios can be used to identify the worst offenders and apply peer pressure to improve. They can also be used as positive reinforcement by encouraging everyone to improve compared to others.

For the system administrator, our results can be applied in two ways. First, the results can be used to make better choices when choosing among available software types and configuration. Second, after systems have been deployed, the findings can be used to manage heterogeneous configurations (e.g., environments with multiple CMSes and server software types). Here, administrators can prioritize how defensive countermeasures such as attack detection should be deployed. Security policies could even be set in accordance with the observed relative risk.

More broadly, we have demonstrated a general method of studying how webserver characteristics affect the risk of compromise. The methods presented here can be applied to other characteristics if the data can be collected. Furthermore, odds ratios help to identify relationships that should be tested further using experimental methods.

## 7 RELATED WORK

While often challenging to carry out, substantial progress has been made over the past several years in conducting large-scale measurements of cybercrime. The most relevant work to ours

is from Soska and Christin who use website features to predict whether a webserver will be hacked in the future [19]. Where we explicitly parse out webserver features, Soska and Christin let their algorithm determine the relevant features.

Some work is particularly relevant due to the results from studying the security of webserver. Doupe et al. describe a state-aware fuzzer in which they evaluate vulnerabilities in CMS platforms [20]. Scholte et al. study vulnerabilities in CMS platforms, though they do not relate vulnerabilities to exploits or observed compromise [21]. Wardman et al. analyze phishing URLs to find common substrings; their method inadvertently finds vulnerable CMS plugins [22]. Nikiforakis et al. crawl many webpages on top webserver to measure the quality of third-party JavaScript libraries running on the webserver [2]. John et al. create “heat-seeking” honeypots for attackers, some running commonly exploited CMSes, to observe attacker behavior [23]. Moore and Clayton measure the recompromise of webserver abused for phishing and find that attackers using targeted, “evil” Google search queries to discover the webserver is a positive risk factor for reinfection [24].

Another series of papers are relevant to the compromise datasets we study. For example, Wang et al. performed a large-scale study of cloaking, which is often caused by search-redirect attacks [25]. Notably, the authors dealt with false positives using clustering. While our data source on search-redirect attacks focuses exclusively on redirections to unlicensed pharmacies [1], the attack technique is general [26]. Provos et al. study drive-by-download URLs [27]. They find more outdated versions of server software on malware landing pages than up-to-date versions.

Various techniques from epidemiology have been applied to cybersecurity, not just case-control studies. For example, Moore et al. and Zou et al. constructed analytical models to model the spread of the the computer worm Code Red as well as efforts to contain it [28] [29]. Holt et al. and Bossler et al. use routine activity theory to look at factors affecting cybercrime victimization [30] [31].

A number of studies deploy methods in common with our own. Notably, Lee describes the use of a small case-control study to identify characteristics that predispose academics to spear-phishing attempts [32]. This paper was expanded by Thonnard et al. to identify characteristics that predispose business people to targeted attacks [33]. Allodi and Massacci analyze the risk of a vulnerability being exploited based on the CVSS score using a case-control study [34]. Additionally, Carlinet et al. use a case-control study to quantify which ISP customer computers are more at risk to generate malicious traffic [35].

We adopt one of the signals of security hygiene used by [2], while Pitsillidis et al. measure the purity of spam feeds in a manner consistent with how we detect false positives in our compromise datasets [36].

Many studies have been primarily descriptive in nature, though some have managed to tease out the factors affecting the prevalence and success of attacks. For instance, Ransbotham connected vulnerability information with intrusion detection system data to show that open-source software tends to be exploited faster than closed-source software following vulnerability disclosure [37].

Our work is distinguished from prior work in two ways. First, we focus extensively on the relationship between webserver characteristics, notably CMS type and market share, and compro-

mise. Second, we use the case-control method to understand the characteristics of large cybercrime datasets.

## 8 CONCLUDING REMARKS

We have presented a case-control study identifying several web-server characteristics that are associated with higher and lower rates of compromise. We joined two infection datasets on phishing and search-redirectation attacks with a large sample of web-servers, then automatically extracted several characteristics of these web-servers hypothesized to affect the likelihood the webserver will be compromised. Notably, our approach is data-driven and our analysis has focused on *security outcomes*, not security levels. By studying compromise data, we have reported on what factors affect the likelihood of actually being hacked, not merely what makes a system vulnerable.

Supported by statistical methods of odds ratios and logistic regression models, we found that certain server types (notably Apache and Nginx) and content management systems (notably Joomla and WordPress) face higher odds of compromise, relative to their popularity. We also found dozens of WordPress plugins and Joomla extensions that were positive risk factors for compromise. We established that a key driving factor behind which CMSes are targeted most is the underlying popularity of the platform. This in turn led to perhaps our most surprising finding: outdated CMS software is hacked less often than up-to-date software. We showed this to be true for the core WordPress software, as well as for 14 of the 19 plugins where we could obtain reliable versioning information. In many respects, this finding can be thought of as a webserver-based corollary to the old truism for desktop operating systems that Macs are more secure than PCs because they have less market share.

By inspecting a supplementary dataset of 50 000 WordPress websites that were hacked to distribute malware, we found that less than 5% updated their software after being compromised. Those that did, however, were much less likely to be re-compromised later on: 16% were re-compromised, compared to 32% for those that never updated. Hence, while our results suggest that websites need not rush to update their CMS software prior to compromise, they definitely should do so once a hack has occurred, since it may close the hole that allowed the attacker to get in.

We are optimistic that the case-control method employed here may be applied to many other contexts of cybercrime measurement. We note that, when possible, the findings of case-control studies should be complemented by other forms of experimentation that directly isolate explanatory factors. It is our hope that both retrospective case-control studies and prospective experimentation will be more widely adopted by cybercrime researchers, which will in turn yield deeper understanding of the issues defenders should prioritize.

## ACKNOWLEDGMENTS

This work was partially funded by the Department of Homeland Security (DHS) Science and Technology Directorate, Cyber Security Division (DHS S&T/CSD) Broad Agency Announcement 11.02, the Government of Australia and SPAWAR Systems Center Pacific via contract number N66001-13-C-0131. This paper represents the position of the authors and not that of the aforementioned agencies.

## REFERENCES

- [1] N. Leontiadis, T. Moore, and N. Christin, "Measuring and analyzing search-redirectation attacks in the illicit online prescription drug trade," in *Proceedings of USENIX Security 2011*, San Francisco, CA, Aug. 2011.
- [2] N. Nikiforakis, L. Invernizzi, A. Kapravelos, S. V. Acker, W. Joosen, C. Kruegel, F. Piessens, and G. Vigna, "You are what you include: Large-scale evaluation of remote JavaScript inclusions," in *ACM Conference on Computer and Communications Security*, 2012, pp. 736–747.
- [3] J. Schlesselman, *Case-control studies: design, conduct, analysis*. Oxford University Press, USA, 1982, no. 2.
- [4] R. Doll and A. Hill, "Lung cancer and other causes of death in relation to smoking; a second report on the mortality of british doctors," *British Medical Journal*, vol. 2, pp. 1071–1081, Nov. 1956.
- [5] Verisign, "The domain name industry brief," Apr. 2013, <https://www.verisigninc.com/assets/domain-name-brief-april2013.pdf>. Last accessed May 1, 2013.
- [6] "PhishTank," 2014, <https://www.phishtank.com/>.
- [7] "Anti-Phishing Working Group," 2014, <http://www.antiphishing.org/>.
- [8] APWG, "Global phishing survey: Trends and domain name use in 2H2012," 2013, [http://docs.apwg.org/reports/APWG\\_GlobalPhishingSurvey\\_2H2012.pdf](http://docs.apwg.org/reports/APWG_GlobalPhishingSurvey_2H2012.pdf). Last accessed May 5, 2013.
- [9] N. Leontiadis, T. Moore, and N. Christin, "Pick your poison: pricing and inventories at unlicensed online pharmacies," in *ACM Conference on Electronic Commerce*, 2013.
- [10] W3techs, "Market share trends for content management systems," [http://w3techs.com/technologies/history\\_overview/content\\_management/](http://w3techs.com/technologies/history_overview/content_management/). Last accessed May 3, 2013.
- [11] "MaxMind GeoIP," 2013, [https://www.maxmind.com/en/geolocation\\_landing](https://www.maxmind.com/en/geolocation_landing).
- [12] "FDIC institutions," 2013, <http://www2.fdic.gov/idasp/Institutions2.zip>.
- [13] J.-H. Hoepman and B. Jacobs, "Increased security through open source," *Communications of the ACM*, vol. 50, no. 1, pp. 79–83, 2007.
- [14] M. Maunder, "Zero day vulnerability in many WordPress themes," 2011, <http://markmaunder.com/2011/08/01/zero-day-vulnerability-in-many-wordpress-themes/>. Last accessed October 7, 2014.
- [15] P. Chapman, "'New software version' notifications for your site," 2009, <http://googlewebmastercentral.blogspot.com/2009/11/new-software-version-notifications-for.html>.
- [16] "URLFind," 2013, <http://urlfind.org/>.
- [17] StopBadware and Commtouch, "Compromised websites: An owners perspective," 2012, <https://www.stopbadware.org/files/compromised-websites-an-owners-perspective.pdf>.
- [18] R. Anderson and T. Moore, "The economics of information security," *Science*, vol. 314, no. 5799, pp. 610–613, Oct. 2006.
- [19] K. Soska and N. Christin, "Automatically detecting vulnerable websites before they turn malicious," in *Proceedings of the 23rd USENIX Security Symposium (USENIX Security '14)*, San Diego, CA, Aug. 2014, pp. 625–640.
- [20] A. Doupe, L. Cavedon, C. Kruegel, and G. Vigna, "Enemy of the State: A State-Aware Black-Box Vulnerability Scanner," in *Proceedings of the USENIX Security Symposium*, Bellevue, WA, August 2012.
- [21] T. Scholte, D. Balzarotti, and E. Kirda, "Quo vadis? A study of the evolution of input validation vulnerabilities in web applications," in *Financial Cryptography and Data Security*. Springer, 2012, pp. 284–298.
- [22] B. Wardman, G. Shukla, and G. Warner, "Identifying vulnerable websites by analysis of common strings in phishing URLs," in *Proceedings of the Fourth eCrime Researchers Summit*. IEEE, 2009, pp. 1–13.
- [23] J. P. John, F. Yu, Y. Xie, A. Krishnamurthy, and M. Abadi, "Heat-seeking honeypots: Design and experience," in *Proceedings of the 20th International Conference on the World Wide Web*. ACM, 2011, pp. 207–216.
- [24] T. Moore and R. Clayton, "Evil searching: Compromise and re-compromise of internet hosts for phishing," in *13th International Conference on Financial Cryptography and Data Security*, Barbados, February 2009.
- [25] D. Wang, S. Savage, and G. Voelker, "Cloak and dagger: Dynamics of web search cloaking," in *Proceedings of the 18th ACM Conference on Computer and Communications Security*. ACM, 2011, pp. 477–490.
- [26] Z. Li, S. Alrwais, Y. Xie, F. Yu, and X. Wang, "Finding the linchpins of the dark web: A study on topologically dedicated hosts on malicious web infrastructures," in *34th IEEE Symposium on Security and Privacy*, 2013.
- [27] N. Provos, P. Mavrommatis, M. Rajab, and F. Monrose, "All your iFrames point to us," in *Proceedings of the 17th USENIX Security Symposium*, Aug. 2008.

- [28] D. Moore, C. Shannon, and J. Brown, “Code-Red: a case study on the spread and victims of an internet worm,” in *Proceedings of 2nd ACM/USENIX Internet Measurement Workshop*, Marseille, France, Nov. 2002, pp. 273–284.
- [29] C. C. Zou, W. Gong, and D. Towsley, “Code red worm propagation modeling and analysis,” in *Proceedings of the 9th ACM Conference on Computer and Communications Security*. ACM, 2002.
- [30] A. M. Bossler and T. J. Holt, “On-line activities, guardianship, and malware infection: An examination of routine activities theory,” *International Journal of Cyber Criminology*, vol. 3, no. 1, pp. 400–420, 2009.
- [31] T. J. Holt and A. M. Bossler, “Examining the applicability of lifestyle-routine activities theory for cybercrime victimization,” *Deviant Behavior*, vol. 30, no. 1, pp. 1–25, 2008.
- [32] M. Lee, “Who’s next? identifying risks factors for subjects of targeted attacks,” in *Proceedings of the Virus Bulletin Conference*, 2012, pp. 301–306.
- [33] O. Thonnard, L. Bilge, A. Kashyap, and M. Lee, “Are you at risk? Profiling organizations and individuals subject to targeted attacks,” in *Financial Cryptography and Data Security*, 2015.
- [34] L. Allodi and F. Massacci, “Comparing vulnerability severity and exploits using case-control studies,” *ACM Transactions on Information and System Security*, vol. 17, no. 1, pp. 1:1–1:20, 2014.
- [35] L. Carlinet, L. Mé, H. Debar, and Y. Gourhant, “Analysis of computer infection risk factors based on customer network usage,” in *Conference on Emerging Security Information, Systems and Technologies*. IEEE, 2008, pp. 317–325.
- [36] A. Pitsillidis, C. Kanich, G. Voelker, K. Levchenko, and S. Savage, “Taster’s choice: A comparative analysis of spam feeds,” in *ACM SIGCOMM Conference on Internet Measurement*, 2012, pp. 427–440.
- [37] S. Ransbotham, “An empirical analysis of exploitation attempts based on vulnerabilities in open source software,” in *Proceedings (online) of the 9th Workshop on Economics of Information Security*, Cambridge, MA, Jun. 2010.
- [38] “BlindElephant web application fingerprinter,” <http://blindelephant.sourceforge.net/>.
- [39] “WhatWeb,” 2013, <http://whatweb.net/>.
- [40] “Plecost,” 2013, <https://code.google.com/p/plecost/>.
- [41] “Exploit database,” 2013, <http://www.exploit-db.com>.

## APPENDIX A COMPARISON OF METHODS TO IDENTIFY CMS TYPE

While a number of tools provide CMS detection as part of more general-purpose web-service fingerprinters (e.g., BlindElephant [38], WhatWeb [39] and the WordPress-specific Plecost [40]), we opted to build the custom CMS detector described above to improve efficiency and accuracy over existing tools. Both BlindElephant and Plecost issue many HTTP requests to characterize each server. We ruled these tools out because we needed a lightweight solution that could quickly detect CMS type and version for hundreds of thousand webservers. Like our method, WhatWeb issues a single HTTP request per server (at its lowest “aggressiveness” level). Combined with its multi-threaded design, WhatWeb should offer fast identification of CMS versions. We therefore decided to evaluate its performance and accuracy compared to our own system.

We selected 2000 random URLs from the webserver dataset and attempted to identify the CMS type using our system and WhatWeb’s. In terms of efficiency, we were surprised to find that WhatWeb took nearly twice as long to finish, despite being multithreaded. We speculate that the difference in speed can be attributed to its general-purpose nature. We also found that our system was substantially more accurate, identifying the correct CMS on more websites and having far fewer inaccurate classifications. We manually inspected all disagreements between WhatWeb and our tool in order to establish the following detection, false positive and false negative rates:

Method	FN Rate	FP Rate	TN Rate	TP Rate	# Results
WhatWeb	40.7%	6.1%	74.3%	59.3%	1297
Our Method	5.4%	0.1%	99.0%	92.2%	1674

Based on these findings, we conclude that our custom method is best-suited to the task of identifying CMS type.

## APPENDIX B DOES CMS POPULARITY AFFECT EXPLOITABILITY?

Results from the Subsection 3.1 showed that the some of the most popular CMS platforms, notably WordPress and Joomla, are compromised disproportionately often. We now dig a bit deeper to see if there is a statistically robust connection between CMS popularity and compromise. Before inspecting the compromise rates directly, we first compare CMS popularity to the number of readily-available exploits targeting the CMS platform.

For this analysis, we considered many more CMSes than in other sections. We consider all 52 CMS platforms tracked in [10]. These additional CMSes all have very small market shares, and so not enough registered in our datasets to include in the other analysis. For each CMS we collected the following two indicators: **# Servers:** We took the market share for each CMS from [10] as of January 1, 2013 and multiplied it by population of registered .com domains (106.2 million) and the estimated server response rate (85%) from [5].

**# Exploits:** The Exploit Database [41] is a search engine that curates working and proof-of-concept exploits from a variety of sources, including the popular penetration-testing tool Metasploit. We searched the Exploit Database for each CMS and recorded the number of hits as a measure of how “exploitable” each CMS is. We discarded any results not matching the searched-for CMS. We deem this to be a more accurate measure of attacker interest in and the “hackability” of a content management system than would be counting the vulnerabilities reported for a CMS. Unlikely many vulnerabilities, exploits provide directly actionable information to compromise machines.

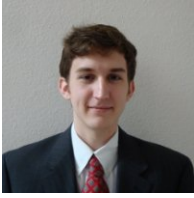
We hypothesize that the number of exploits available for a CMS depends directly on the number of servers in use. Because both variables are highly skewed, we apply a log transformation to each. Here is the statement of the linear regression:

$$\lg(\# \text{ Exploits}) = c_0 + c_1 \lg(\# \text{ Servers}) + \varepsilon.$$

The regression yields the following results:

	coef.	95% conf. int.	Significance
Intercept	<b>-8.53</b>	(-3.37, -13.69)	$p = 0.002$
$\lg(\# \text{ Servers})$	<b>0.64</b>	(0.33, 0.95)	$p = 0.0001$
Model fit: $R^2 = 0.23$			

Indeed, this simple linear model has a reasonably good fit. While there is additional unexplained variation, this lends indirect support to **H2**. Due to the collinearity of these variables, we only use one of them (# Servers) in our regressions in this paper.



**John Wadleigh** is a research assistant in the Department of Computer Science and Engineering at Southern Methodist University. He is pursuing an MS in computer science, and his research has investigated the prevalence of web-based malware and the sale of counterfeit goods.



**Marie Vasek** is a PhD student in the computer science department at Southern Methodist University and the research scientist at Stop-Badware. Her research interests include security economics and cybercrime, particularly web-based malware.



**Tyler Moore** is an Assistant Professor of Computer Science and Engineering at Southern Methodist University. Also at SMU, he serves as Director of the Economics and Social Sciences program at the Darwin Deason Institute for Cyber Security and a Fellow at the John Goodman Tower Center for Political Studies. His research focuses on security economics, cyber-crime measurement, and cybersecurity policy.