

# Lurking Malice in the Cloud: Understanding and Detecting Cloud Repository as a Malicious Service

Xiaojing Liao<sup>1</sup>, Sumayah Alrwais<sup>2</sup>, Kan Yuan<sup>2</sup>, Luyi Xing<sup>2</sup>,  
XiaoFeng Wang<sup>2</sup>, Shuang Hao<sup>3</sup>, Raheem Beyah<sup>1</sup>  
<sup>1</sup>Georgia Institute of Technology, <sup>2</sup>Indiana University Bloomington,  
<sup>3</sup>University of California, Santa Barbara

{xliao, rbeyah}@gatech.edu, {salrwais, kanyuan, luyixing, xw7}@indiana.edu, shuanghao@cs.ucsb.edu

## Abstract

The popularity of cloud hosting services also brings in new security challenges: it has been reported that these services are increasingly utilized by miscreants for their malicious online activities. Mitigating this emerging threat, posed by such “bad repositories” (simply *Bar*), is challenging due to the different hosting strategy to traditional hosting service, the lack of direct observations of the repositories by those outside the cloud, the reluctance of the cloud provider to scan its customers’ repositories without their consent, and the unique evasion strategies employed by the adversary. In this paper, we took the first step toward understanding and detecting this emerging threat. Using a small set of “seeds” (i.e., confirmed Bars), we identified a set of collective features from the websites they serve (e.g., attempts to hide Bars), which uniquely characterize the Bars. These features were utilized to build a scanner that detected over 600 Bars on leading cloud platforms like Amazon, Google, and 150K sites, including popular ones like [groupon.com](http://groupon.com), using them. Highlights of our study include the pivotal roles played by these repositories on malicious infrastructures and other important discoveries include how the adversary exploited legitimate cloud repositories and why the adversary uses Bars in the first place that has never been reported. These findings bring such malicious services to the spotlight and contribute to a better understanding and ultimately eliminating this new threat.

## 1. INTRODUCTION

Cloud hosting service today is serving over a billion users worldwide, providing them stable, low-cost, reliable, high-speed and globally available resource access. For example, Amazon Simple Storage Service (S3) is reported to store over 2 trillion objects for web and image hosting, system backup, etc. In addition to storing data, these services are moving toward a more active role in supporting their customers’ computing missions, through sharing the repositories (a.k.a. *bucket* for Google Cloud [?]) hosting various dynamic content and programming tools. A prominent example is Google’s Hosted Libraries [?], a content distribution network (CDN) for disseminating the most popular, open-source JavaScript resources, which web developers can easily incorporate into their

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

CCS’16, October 24–28, 2016, Vienna, Austria

© 2016 ACM. ISBN 978-1-4503-4139-4/16/10...\$15.00

DOI: <http://dx.doi.org/10.1145/2976749.2978349>

websites through a simple code snippet. In addition to benign users, the popularity of these services has also attracted cybercriminals. Compared with dedicated underground hosting services, repositories on legitimate commercial clouds are more reliable and harder to blacklist. They are also much cheaper: for example, it is reported that 15 GB on the dark net is sold at \$15 per month [?], which is actually offered for free by Google to every Google Driver user. Indeed, it has been reported [?] that malware distributors are increasingly using the commercial clouds to process and deploy malicious content.

**Understanding bad cloud repositories: challenges.** Although there have been indications of cloud hosting misuse, understanding how such services are abused is challenging. For the service providers, who are bound by their privacy commitments and ethical concerns, they tend to avoid inspecting the content of their customers’ repositories in the absence of proper consent. Even when the providers are willing to do so, determining whether a repository involves malicious content is by no means trivial: nuts and bolts for malicious activities could appear perfectly innocent before they are assembled into an attack machine; examples include image files for Spam and Phishing as shown in Figure 1. Actually, even for the repository confirmed to serve malicious content like malware, today’s cloud providers tend to only remove that specific content, instead of terminating the whole account, to avoid collateral damage (e.g., compromised legitimate repositories). Exploring the issue becomes even more difficult for the third party, who does not have the ability to directly observe the repositories and can only access them through the websites or sources that utilize the storage services. Further adding to the complexity of finding such a repository is the diverse roles it may play in attack infrastructures (e.g., serving malware for one attack and serving Phishing content for another), due to the mixed content a single repository may host: e.g., malware together with Phishing images. As a result, existing techniques (e.g., those for detecting dedicated malicious services [?][?]) cannot be directly applied to capture the repository, simply because their original targets often contain more homogeneous content (e.g., just malware) and contribute to different campaigns in the same way. So far, little has been done to understand the scope and magnitude of malicious or compromised repositories on legitimate clouds (called *Bad Repository* or simply *Bar* in our research) and the technical details about their services to the adversary, not to mention any effort to mitigate the threat they pose.

**Finding “Bars” online.** In this paper, we present the first systematic study on the abuses of cloud repositories on the legitimate cloud platforms as a malicious service, which was found to be highly pervasive, acting as a backbone for large-scale malicious web campaigns (Section 4). Our study was bootstrapped by a set of “seeds”: 100 confirmed malicious or compromised buckets [?], each of which

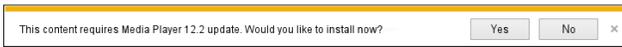


Figure 1: Example of deceptive images in Amazon S3 bucket *cicloudfront* used for malvertising. The image was shown at the bottom of a webpage as an update notification to lure visitors to download malware.

is a cloud resource repository with stored objects (often of different types) organized under a unique identification key. These buckets were collected from Spam messages or the malicious URLs cached by a popular malware scanner. Comparing them with those known to be legitimate, we found that despite various roles each bucket plays in different types of attacks (due to the diversity in the content it serves), still the websites connecting to those buckets exhibit prominent common features (see Section 3.1), particularly, the presence of “gatekeeper” sites that cover the Bars (a valuable asset for the adversary) and remarkably homogeneous redirection behavior (i.e., fetching repository resources indirectly through other sites’ references) and sometimes similar content organizations, due to the same attack payload the compromised sites upload from their backend (i.e., the Bars), or the templates the bucket provides to the adversary for quick deployment of her attack sites. By comparison, a legitimate bucket (e.g., reputable jQuery repository) tends to be *directly* accessed by the websites with highly diverse content.

Based on this observation, we developed *BarFinder*, a scanner that automatically detects Bars through inspecting the topological relations between websites and the cloud bucket they use, in an attempt to capture Bars based on the external features of the websites they serve. More specifically, for all the sites connecting to a repository, our approach correlates the domains and URLs (particular those related to cloud repositories) across their redirection chains and content features across their DOM structures to identify the presence of gatekeepers and evading behavior, and also measure the diversity of their content organization. A set of new *collective features* generated in this way, including *bucket usage similarity*, *connection ratio*, *landing similarity* and others (Section 3.1), are further utilized by a classifier to find out suspicious buckets. Running the scanner over all the data collected by the Common Crawl [?], which indexed five billion web pages, for those associated with all major cloud storage providers (including Amazon S3, Cloudfront, Google Drive, etc.), we found around 1 million sites utilizing 6,885 repositories hosted on these clouds. Among them, *BarFinder* identified 694 malicious or compromised repositories, involving millions of files, with a precision of 95% and a coverage of 90% against our ground-truth set.

**Our discoveries.** Looking into the Bars identified by our scanner, we are surprised by the scope and the magnitude of the threat. These buckets are hosted by the most reputable cloud service providers. For example, 13.7% of Amazon S3 repositories and 5.5% of Google repositories that we inspected turned out to be either compromised or completely malicious<sup>1</sup>. Among those compromised are popular cloud repositories such as Groupon’s official bucket. Altogether, 472 such legitimate repositories were considered to be contaminated, due to a misconfiguration flaw never reported before, which allows arbitrary content to be uploaded and existing data to be modified without proper authorization. The impact of these Bars is significant, infecting 1,306 legitimate websites, including Alexa top 300 sites like [groupon.com](http://groupon.com), Alexa top 5,000 sites like [space.com](http://space.com), etc. We reported our findings to Amazon and leading organizations affected by the infections. Groupon has already confirmed the compromise we discovered and awarded us for our help.

<sup>1</sup>We have manually examined and confirmed all those instances.

When it comes to malicious buckets, our study brings to light new insights into this new wave of repository based cyber-attacks, including the importance of Bars to malicious web activities and the challenges in defending against this new threat. More specifically, we found that on average, one Bar serves 152 malicious or compromised sites. In one of the large campaigns discovered in our research, the Bar [cloudfront\\_file.enjin.com](http://cloudfront_file.enjin.com) hosts a malicious script that was injected into at least 1,020 websites (Section 4.1). These Bars sit right at the center of the attack infrastructure, supporting and coordinating other malicious actors’ operations at different stages of a campaign. Interestingly, we found that they could be strategically placed on different cloud platforms, making them hard to block (due to the popularity of their hosting clouds like Google) and detect (scattered across different providers), and easy to share across multiple campaigns. As an example, the Potentially Unwanted Programs (PUP) campaign we found first loads a redirection script from a Bar on Akamaihd (the world’s largest CDN platform) to lead the victim to the attack website, then fetches Phishing pictures from an Amazon S3 Bar, and finally delivers the malware stored on Cloudfront to the target systems (Section 4.4). In the presence of such meticulously planned attacks, the cloud service providers apparently are inadequately prepared, possibly due to the privacy constraints in touching their customers’ repositories. We found that many *Bars* remain active during our study, and survive a much longer lifetime than that of the malicious content hosted on websites (Section 4.3). Further complicating the mission of Bar identification are other evasion techniques the adversary employs, including code obfuscation and use of a redirection chain and cloaking techniques to avoid exposing malicious payloads to a malware scanner.

**Contributions.** The contributions of the paper are as follows:

- *New understanding.* We performed the *first* systematic study on cloud repositories as a malicious service, an emerging security threat. For the first time, our study reveals the scope and magnitude of the threat and its significant impact, particularly on the infrastructures of illicit web activities. These findings bring to the spotlight this important yet understudied problem and lead to a better understanding of the techniques the adversary employs and their weaknesses. This will contribute to better defense against and ultimate elimination of the threat.
- *New technique.* Based on our understanding of bad cloud repositories, we take a first step toward automatically detecting them. The technique we developed relies on the topological relationship between a cloud repository and the websites it serves, which are difficult to change and effective at capturing malicious or compromised buckets. Our evaluation over a large number of popular websites demonstrates the potential of the technique, which could be utilized by both cloud providers and third parties to identify the threats posed by *Bars*.

**Roadmap.** The rest of the paper is organized as follows: Section 2 provides the background information and adversary model for our research; Section 3 describes our findings from the ground-truth dataset and the design and implementation of *BarFinder*; Section 4 provides the details of the discoveries made in our large-scale measurement study; Section 5 discusses the limitations of our work and possible future research; Section 6 compares our work with related prior research and Section 7 concludes the paper.

## 2. BACKGROUND

**Cloud hosting.** Cloud hosting is a type of infrastructure as a service (IaaS), which is rented by the cloud user to host her web assets (e.g., HTML, JavaScript, CSS, and image files). These web assets are organized into cloud repositories referred to as *buckets* which are identi-

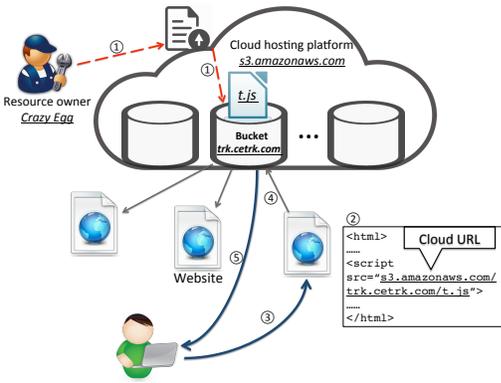


Figure 2: Overview of the cloud hosting process.

fied by unique<sup>2</sup>, user-assigned keys, that are mapped as sub-domains. For example, the subdomain `aws-publicdatasets.s3.amazonaws.com` identifies Amazon S3 as the cloud platform and `aws-publicdatasets` as the user’s cloud bucket and repository. Such name assignment is labeled as `s3.amazonaws.com_aws-publicdatasets` throughout this paper. Also, each bucket is protected by an access control list generated by the user to authorize requests for her resources.

In recent years, we have seen an increase in popularity of these services. A key feature of cloud hosting is built-in site publishing [?], where the web assets in the bucket can be served directly to users via file names in a relative path in the bucket (i.e., *cloud URL*). For instance, JavaScript files hosted in the cloud bucket can be directly run in the browser. Also, the pay-as-you-go hosting is well received as an economic and flexible computing solution. As an example, Google Drive today offers a free web hosting service with 15GB of storage, and an additional 100GB for \$1.99/month, and GoDaddy’s web hosting starts at \$1/month for 100GB.

Besides such front-end websites, mainstream cloud providers today (Amazon S3, Microsoft Azure, Google Drive, etc.) all allow their customers to store different kinds of web content and other resources in their cloud buckets, serving as back-end repositories that can be easily accessed by front-end applications (like the website) and shared across different parties. Figure 2 illustrates an example, in which the resource owner creates a bucket on the cloud hosting platform and uploads a script there (①); this resource (i.e., the script) is made public through a cloud URL, which can be embedded into any website (②); whenever the site is visited (③), requests will be generated for fetching the script (④) and delivering it to the visitor’s browser (⑤). The bucket in the example is typical of a service repository, whose resources can be fetched and updated through a cloud URL: for example, the visitor statistics of a website can be collected through a link (`s3.amazonaws.com/trk.cetrk.com/t.js`), which downloads a tracking script from `s3.amazonaws.com_trk.cetrk.com`, a bucket owned by the tracking website Crazy Egg. This is different from a “self-serving” bucket, whose resources can only be accessed by the bucket owner’s sites. Note that our study focuses on abuses of this type of cloud repositories, regardless of the additional functionalities they may have (e.g., CDNs, DDoS protection, etc.), since these functionalities do not affect the way the repositories are used by either legitimate or malicious parties.

**Adversary model.** In our research, we consider the adversary who tries to use cloud buckets on legitimate cloud platforms as service repositories for illicit activities. For this purpose, the attacker could build her own malicious bucket or compromise legitimate ones, and store various attack vectors there, including Spam, Phishing,

<sup>2</sup>The terms repositories and buckets are used interchangeably throughout this paper.

malware, click-hijacking and others. These buckets are connected to front-end websites, which could be malicious, compromised or legitimate ones contaminated only by the Bar.

### 3. FINDING BARS ONLINE

In this section, we elaborate on our analysis of a set of known Bars (the seed set) and the features identified for differentiating benign repositories and Bars. These features are utilized in our research to build a simple web scanner, *BarFinder*, for detecting other malicious or compromised high-profile, previously-unknown repositories and the malicious campaigns in which they serve.

#### 3.1 Features of Bad Repositories

Our study is based on a small set of confirmed good and bad repositories and their related domains, which we analyzed to find out how Bars (bad repositories) differ from legitimate repositories. In the absence of direct access to these buckets, good or bad, all we can do is to *infer* their legitimacy from who use them and how they are used (by different domains), that is, the features of the domains and their interactivities on the redirection paths leading to the cloud repository. Of particular interest here are a set of *collective* properties identified from the resource fetching chains (a.k.a., redirection chains) for serving the content of Bars, which is hard to change by the adversary, compared with the content features of individual Bars. Below, we elaborate on the way such data was collected and the salient features discovered in our research, which describe how the adversary attempts to hide Bars or use them to cover other attack assets, a redirection pattern never observed on legitimate repositories.

**Data collection.** To build the seed set, we collected a set of confirmed malicious or compromised buckets (called *Badset*) and legitimate buckets (called *Goodset*) as well as their related domains, as illustrated in Table 1.

- *Badset.* We utilized two feeds as the ground truth for gathering bad cloud buckets: the *Spamtrap* feed and the *CleanMX* feed [?]. The former comes from a Spam honeypot we constructed [?] that receives around 10K Spam emails per day, from which cloud URLs promoted by the emails were extracted which may include spam resources such as HTML, images, and scripts. The latter includes the historical data of CleanMX, a popular domain scanning engine, from which cloud-related URLs were collected. For both feeds, we further validate them by VirusTotal [?] and manual inspections (e.g., looking for Phishing content) to ensure that they were indeed bad (to avoid contaminating the dataset with legitimate buckets used in malicious activities). Using the collected set of malicious cloud URLs from both feeds, we extracted their repositories, which led to 100 confirmed Bars.

- *Goodset.* The good buckets were gathered from the Alexa top 3K websites, which are considered to be mostly clean. To this end, we visited each website using a crawler (as a Firefox add-on) to record the HTTP traffic triggered by the visit, including network requests, responses, browser events, etc. From the collected traffic, we extracted the HTTP cloud request URLs corresponding to 300 cloud buckets hosted on 20 leading cloud hosting services like Amazon S3, Google Drive, etc. (see Appendix Table 7 for the complete list). Note that even though some of them provide CDN service or DDOS protection, they are all provided hosting service to act as cloud repository.

- *Bucket-served sites and their HTTP traffic.* We collected HTTP traffic using the crawler mentioned above to visit a list of websites using buckets for feature extraction. Rather than blindly crawling the web to find those sites, we adopted a more targeted strategy by

Table 1: Summary results of the seed dataset.

|         | # of buckets | # of linked websites | # of average linked website | # of redirection paths |
|---------|--------------|----------------------|-----------------------------|------------------------|
| Badset  | 100          | 12,468               | 133                         | 468,480                |
| Goodset | 300          | 128,681              | 864                         | 2,659,304              |

crawling the sites found to contain links to the cloud in the past. We built the site list with the help of *Common Crawl* [?], a public big data project that crawls about 5 billion webpages each month through a large-scale Hadoop-based crawler and maintains lists of the crawled websites and their embedded links. Searching the *Common Crawl* [?] dataset, collected in February 2015, for the websites loading content from the 400 clean and malicious buckets identified above, we found 141,149 websites, were used by our crawler.

**Topological features.** We first inspected the *topology* of the redirection infrastructure associated with a specific bucket. Such an infrastructure is a collection of redirection paths, with each node being a *Fully Qualified Domain Name* (FQDN). On each path, the bucket is either a node when it directly participates in a redirection (e.g., its cloud URL delivers a redirection script to the visitor’s browser) or simply a passive repository providing resources like pictures to other domains. Figure 3 illustrates examples of redirection paths leading to two real-world repositories, one for a legitimate bucket [cloudfront.net\\_d24n15hnbwhuhn](https://cloudfront.net_d24n15hnbwhuhn) and the other for a Bar [s3.amazonaws.com\\_cicloudfront](https://s3.amazonaws.com_cicloudfront).

A key observation from our study is that *the redirection infrastructure leading to a Bar tends to include the features for protecting the Bar from being detected by web scanners*, presumably due to the fact that the repository is often considered to be a valuable asset for the adversary. Specifically, we found that typically, there are a few *gatekeeper* nodes sitting in front of a Bar, serving as an intermediary to proxy the attempts to get resources from the Bar. Examples of the gatekeepers include [fp125.mediaoptout.com](https://fp125.mediaoptout.com) and its downstream nodes in Figure 3(b). On the topology of such an infrastructure, these gatekeepers are the hubs receiving a lot of resource-access connections from entry sites (the first node on a redirection path, see Figure 3). Also interestingly, our research shows that some gatekeepers can access the Bar through multiple paths. For example, in Figure 3(b), [krd.semantichelper.com](https://krd.semantichelper.com) can either go straight to [s3.amazonaws.com\\_cicloudfront](https://s3.amazonaws.com_cicloudfront) or take a detour through [p306.atemada.com](https://p306.atemada.com). This structure could be caused by the cloaking of the gatekeeper for hiding the Bar, or constructed to maintain access to the repository even when nodes (like [1.semantichelper.com](https://1.semantichelper.com)) are down (detected, cleaned, etc.). Note that such a protection structure does not exist on the paths to a benign repository (Figure 3(a)): normally, the resources hosted in a repository (e.g., jQuery) is directly fetched by the website using it, without going through any redirection; even in the presence of redirections, there will not be any gatekeeper, not to mention attempts to cloak or build a backup path.

To identify this unique “protection” structure, we utilize two *collective features*: *bucket usage similarity* (BUS) that captures the topology involving hubs (gatekeepers) and *connection ratio* (CR) that measures the interactivities across different redirection paths (which point to the existence of cloaking behavior or the attempts to maintain back-up paths to the Bar). Specifically, consider a redirection graph  $G = (V, E)$  (as illustrated in Figure 3), where  $V$  is the set of nodes (the FQDNs involved in a redirection) and  $E$  is a set of edges from one node to the next one on individual paths:  $E = \{e_{i,j} | \text{node } i \text{ precedes node } j \text{ on a path}\}$ . The BUS is

measured by  $1 - \frac{i}{s}$ , where  $i$  is the number of immediate predecessor nodes to a repository (the domains connecting to the repository) and  $s$  is the total number of entries of the repository’s redirection graph. To find out the CR, we first remove the bucket  $b$  and all the edges to which it is attached (if they exist) to get another graph  $G' = G - G_b$ , where  $G_b = (\{b\}, E_b)$  and  $E_b = \{e_{b,j}\}$ . Note that each graph  $G'$  is associated with one bucket. Then, from  $G'$ , we find out the number of connected components  $n$  and calculate  $CR = 1 - \frac{n}{|V|}$  (see Figure 3 for an example).

Both collective features were found to be discriminative in our research. Figure 4(a) and 4(b) compare the cumulative distributions (CDF) of the ratios between Bad and Good sets. As we can see from the figures, Bars tend to have higher ratios than benign ones: the average BUS is 0.87 for the Bars and 0.79 for the legitimate repositories and the CR is 0.85 for the bad repositories and 0.67 for the good one. As mentioned earlier, this is caused by the fact that a small set of gatekeepers nodes are often placed there for protecting the Bars while the redirection chains towards the good repositories are much more direct and independent: different organizations typically do not go through an intermediary to indirectly access the public repository like jQuery, and even within the same organization, use of such a resource is often direct. Although there can be exceptions, our measurement study shows that in general, the structural differences between malicious and legitimate repositories are stark.

Also, we found that occasionally, a Bar itself may serve as a gatekeeper, running scripts to hide more valuable attack assets, such as the attack server or other malicious landing sites. When this happens, almost always the Bar leads to a small set of successors on redirection paths (e.g., attack servers, land sites). This is very different from the redirection performed by the script from a benign repository, for example, [cloudfront.net\\_d24n15hnbwhuhn](https://cloudfront.net_d24n15hnbwhuhn). In such cases, the targets of redirections are often very diverse. Based on this observation, we further measure the *landing similarity*,  $LS = 1 - \frac{l}{s}$ , where  $l$  is the number of the unique last nodes on the redirection paths associated with a repository. Again, as illustrated in Figure 4(c), our study shows that redirection paths involving Bars share fewer end nodes than legitimate ones, and therefore, the related redirection graphs (for Bars) have a higher landing similarity (0.94 vs 0.88).

**Content and network features.** In addition to their distinctive topological features, we found that the nodes on the redirection paths attached to a Bar often exhibit remarkable homogeneity in their content and network properties. Particularly, for the websites directly connecting to the repository, we found that they typically use a small set of templates (like WordPress) to build up their web pages, include similar DOM positions for script injection, carrying similar IP addresses or even having the same content management system (CMS) vulnerabilities, etc. These properties turn out to be very diverse among those utilizing a legitimate cloud repository. For example, all websites linking to a Google Drive Bar have their malicious cloud URL (for injecting a script) placed at the bottom of the DOM of each website. In another example, we found that the front-end sites using a Cloudfront Bar actually all include a vulnerable JCE Joomla extension.

To better understand the diversity of such websites, we try to compare them according to a set of content and network properties. In our research, we utilized the properties extracted by *WhatWeb* [?], a popular webpage scanner. WhatWeb is designed to identify the web technologies deployed, including those related to web content and communication: e.g., CMS, blogging platforms, statistic/analytics packages, JavaScript libraries, social media plugins, etc. For example, from the content

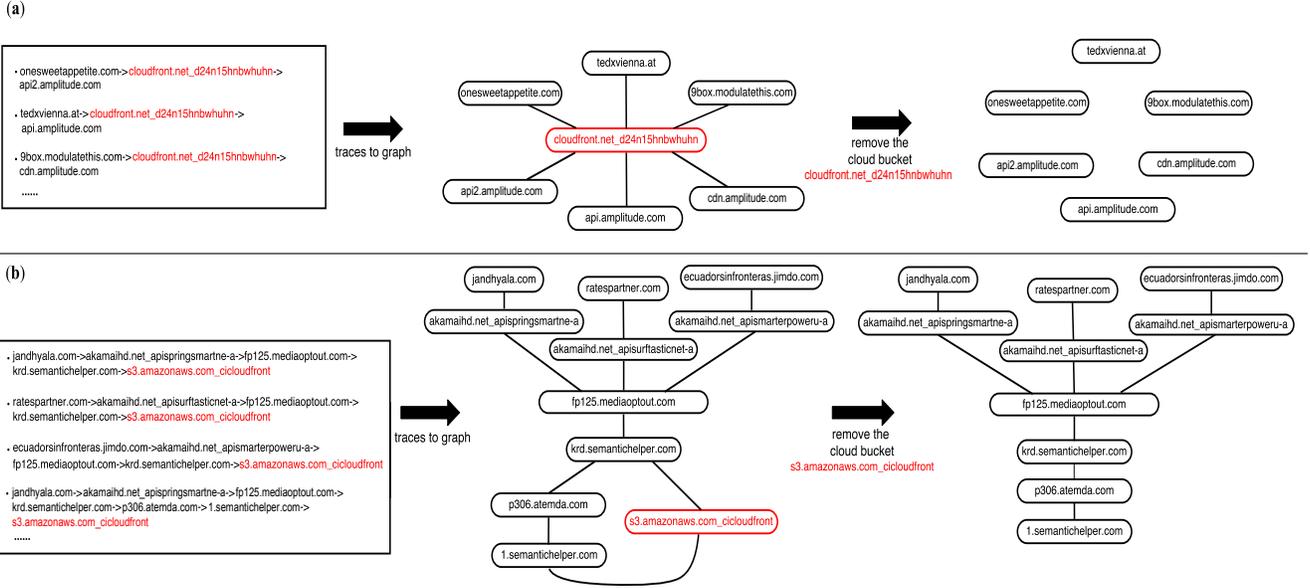


Figure 3: Example of the redirection infrastructure leading to the legitimate bucket `cloudfront.net_d24n15hnbwhuhn` (a) and the Bar `s3.amazonaws.com_cicloudfront` (b), which are in RED color.

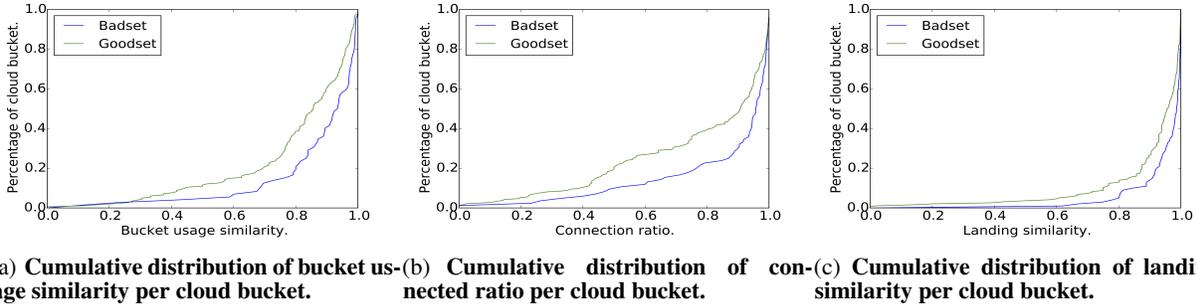


Figure 4: Bars show smaller topological diversity.

```
<link rel="search"
  type="application/opensearchdescription+xml"
  href="https://wordpress.com/opensearch.xml"
  title="WordPress.com" />
```

we obtain the property  $p$  as a key-value pair  $p = (k, v) = (\text{wordpress}, \text{opensearch})$ , which indicates the website using wordpress plugin opensearch.

From our seed dataset, the scanner automatically extracted 372 keys of 1,596,379 properties, and then we clustered the keys into 15 classes such as Analytics and tracking, CMS and plugin, Meta-data information, etc., following the categories used by *BuiltWith*, a web technology search engine [?]. Some examples of these properties are presented in Table 2. In addition to these properties extracted by *WhatWeb*, we added the following properties to characterize cloud URLs, including the position of the URL, the order in which different buckets appear in the web content and the number of cloud platforms used in a page.

Based on these properties, again we utilized a topological metric to measure the overall similarity across sites. Specifically, the relations among all the sites (connecting to the same bucket) in the same category (Analytics and tracking, CMS and plugin, etc.) are modeled as a graph  $G' = (V', E', P)$ , where  $V'$  is the set of the websites, which are characterized by a collection of properties  $P$ , and

$E'$  is the set of edges:  $E' = \{e_{i,j} | \text{website } i \text{ and } j \text{ share } p \in P\}$ , that is, both sites having a common property. Over this graph, the *site similarity* is calculated as  $SiS = 1 - \frac{n}{|V'|}$ . Here  $n$  is the number of connected components in the graph.

In our research, we computed SiS across all the categories summarized from the seed dataset, and compared those with Bars against those with the legitimate buckets. Again, the sites using Bars are found to share many more properties and therefore achieve a much higher similarity value than those linking to a good bucket. This is likely caused by mass-production of malicious sites using the same resources (templates, pictures, etc.) provided by a Bar or utilization of the same exploit tool stored in a Bar for compromising the sites with the same vulnerabilities. Therefore, such similarity is inherent to the attack strategies and can be hard to change.

### 3.2 BarFinder

**Design.** The design of BarFinder includes a web crawler, a feature analyzer, and a detector. The crawler automatically scans the web for cloud buckets (embedded in web content) and then clusters websites according to the buckets they use. From each cluster, the analyzer constructs a redirection graph and a content graph as described earlier (Section 3.1), on which it further calculates the values for a set of collective features including disconnection ratio

Table 2: Examples of content and network features.

| Category | Feature                                   | Example                           |
|----------|---|-----------------------------------|
| Content  | CMS platform information and their plugin | (wordpress, all in one SEO pack)  |
|          | Meta-data information                     | (metagenerator, drupal7)          |
|          | CloudURL information                      | (position, bottom)                |
|          | Advertising                               | (adsense, asynchronous)           |
|          | Javascript library                        | (jQuery, 1.9.1)                   |
|          | Analytics and tracking                    | (Google-Analytics, UA-2410076-31) |
|          | Widget                                    | (addthis, welcome bar)            |
|          | DocInfo technologies                      | (open graph protocol, null)       |
| Network  | Identity                                  | (IP, 216.58.216.78)               |
|          | Cookie                                    | (Cookie, harbor.session)          |
|          | Server framework version                  | (Apache, 2.4.12)                  |
|          | Custom HTTP header                        | (X-hacker, If youâ€™re..)         |

( $D$ ), bucket usage similarity ( $B$ ), landing similarity ( $L$ ) and a series of content property/network property similarities ( $S_1 \cdots S_n$ ) for  $n$  web-technology categories (e.g., analytics and tracking, CMS and plugin, meta-data information, etc.). The output of this feature analysis is then passed to the detector, which maintains a model (trained on the seed dataset) to determine whether a bucket is malicious, based on its collective features.

Specifically, the crawler visits each website, inspecting its content, triggering events, recording the redirection paths it observes and parsing URLs encountered using the patterns of known cloud platforms to recognize cloud buckets. For example, the repository on Amazon S3 is accessed through the URL formatted as  $w + .s3\{-w+\}[?].amazonaws.com$ , and Amazon CloudFront produces resource URLs in the form of  $w + .cloudfront.net$ . In our research, 20 cloud platforms were examined to identify the buckets they host. At the feature-analysis stage, for each bucket, BarFinder inspects all its redirection paths, converts every node into an FQDN to compute their topological features, and then connects different nodes according to their content and network properties to find out their site similarities, as described in Section 3.1.

Next, each cloud bucket  $i$  is uniquely characterized by a vector:  $\langle D_i, B_i, L_i, S_{i,1} \cdots S_{i,n} \rangle$ , with each element a collective feature. Individual features have different power in differentiating good and bad buckets, which we measured using the F-Score [?] (see Table 3). Note that the feature with a large score can better classify these vectors than the one with a small value. Therefore, a binary classifier with a model for weighing the features and other parameters can be used to classify the vector set and determine whether individual buckets are legitimate or not. Such a model is learned from the seed dataset. In our research, we utilized a Support Vector Machine (SVM) as the classifier, which showed the best performance among other classification algorithms (see Table 4). Its classification model is built upon the F-Scores for the collective features ( $D$ ,  $B$ , etc.) and a threshold set according to the false positive and negative discovery expected to achieve. For each bucket classified, the SVM can also report the confidence of the classification.

**Implementation.** This simple design was implemented in our study into a prototype system. The web crawler was built as a Firefox add-on. In total, 20 such crawlers were deployed. We further developed a tool in Python to recover cloud URLs from the web content gathered by Common Crawl. The feature analyzer includes around 500 lines of Python code for processing the data collected by the crawler and

Table 3: F-score of features.

| Feature                 | Label | Metric               | F-score |
|-------------------------|-------|----------------------|---------|
| Connection ratio        | D     | $1 - \frac{n}{ V }$  | 0.084   |
| Bucket usage similarity | B     | $1 - \frac{i}{s}$    | 0.076   |
| Landing similarity      | L     | $1 - \frac{l}{s}$    | 0.072   |
| CMS information         | $S_1$ | $1 - \frac{n}{ V' }$ | 0.037   |
| Meta-data information   | $S_2$ | $1 - \frac{n}{ V' }$ | 0.033   |
| Analytics and tracking  | $S_3$ | $1 - \frac{n}{ V' }$ | 0.032   |
| Widget                  | $S_4$ | $1 - \frac{n}{ V' }$ | 0.031   |
| CloudURL information    | $S_5$ | $1 - \frac{n}{ V' }$ | 0.024   |

Table 4: Performance comparison under five-fold cross validation.

| Classifier          | Precision | Recall |
|---------------------|-----------|--------|
| SVM                 | 0.94      | 0.89   |
| Decision Tree       | 0.9       | 0.83   |
| Logistic Regression | 0.91      | 0.87   |
| Naive Bayes         | 0.9       | 0.79   |
| Random Forest       | 0.85      | 0.82   |

computing the collective features (Section 3.1). Each feature in the vector is normalized using the L1 norm before passed to the SVM classifier. In our system, we incorporated the SVM provided by the scikit-learn open-source machine learning library [?].

### 3.3 Evaluation

Here we report our evaluation of BarFinder on both the ground truth and the Unknown sets. All the experiments were conducted within an Amazon EC2 C4.8xlarge instance equipped with Intel Xeon E5-2666 36 vCPU and 60GiB of memory.

**Evaluation on the seed set.** We tested the effectiveness of BarFinder over our ground-truth dataset (i.e., the seed set) through the standard five-fold cross validation: that is, 4/5 of the data was used for training the SVM and the remaining 1/5 for evaluating the accuracy of Bar detection. Specifically, we randomly chose 80 Bars (out of 100) from the Badset and 240 (out of 300) legitimate buckets from the Goodset, together with the related websites (out of 141,149). These data were first processed by our prototype to adjust the weights and other parameters for its model. Then we tested the model on the remaining dataset (20 Bars, 60 legitimate buckets). The process is then repeated 5 times. BarFinder achieved both a low false discovery rate (FDR: 1- precision) and a high recall in detection: only 5.6% of reported Bars turned out to be legitimate (i.e., 1.6% of false positive rate), and over 89.3% of the Bars were detected. We further show the Area Under Curve (AUC) of the Receiver Operating Characteristics (ROC) graph, which comes very close to 1 (0.96), demonstrating the good balance we strike between the FD rate and the coverage. This preliminary analysis shows that the collective features of the sites connecting to cloud repositories are promising in detecting Bars.

**Evaluation on the Unknown set.** We now use BarFinder to scan an unknown set. This unknown set contains HTTP traffic collected using a crawler as described in Section 3.1 to visit a list of websites. This list of websites is also extracted from common crawl [?] by searching for websites that have loaded some content in the past from the cloud platforms listed in Table 7. As a result, the unknown data set contained HTTP traffic generated from dynamically visiting 1M websites loading content from 20 cloud platforms and 6,885 cloud buckets.

To validate our evaluation results, we employ a methodology that combines anti-virus (AV) scanning, blacklist checking, and manual analysis. Specifically, for the Bars flagged by our system, we first scan their cloud URLs with VirusTotal for malware and

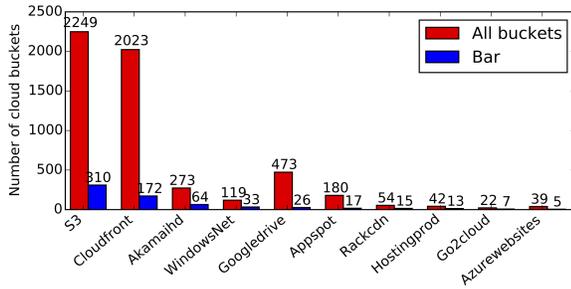


Figure 5: Top 10 cloud platforms with most Bars, compared with their total number of cloud buckets in our dataset.

check them against the list of suspicious cloud URLs collected from our Spamtrap honeypot for Spam, Phishing, blackhat Search Engine Optimization (SEO), etc. In the case of VirusTotal, a URL is considered to be suspicious if at least two scanners raise the alarm. All such suspicious URLs (from either VirusTotal or the Spamtrap list) are cross-checked against the blacklist of CleanMX. Only those also found there are reported to be a true positive. Once a URL is confirmed malicious, its corresponding bucket is labeled as bad. Those unlabeled but flagged (by BarFinder) buckets are further validated manually.

In the experiment, BarFinder reported a total of 730 Bars, about 10.6% of the 6,885 buckets. Among them, the AV scanning and blacklist verification confirmed that 502 buckets were indeed bad. The remaining 228 were manually analyzed through, e.g., inspecting the resources in the buckets for phishing or scam content, running scripts in the VM to capture binary code download. This validation further confirmed 192 Bars. The FDR was found to be at most 5% (assuming those not confirmed to be legitimate), in line with the finding from the seed set.

#### 4. MEASUREMENT AND DISCOVERIES

Based on the discoveries made by BarFinder, we further conducted a measurement study to better understand the fundamental issues about Bar-based malicious services, particularly *how the cloud repositories help facilitate malicious activities, how the adversary exploited legitimate cloud buckets and why the adversary uses Bars in the first place*. Our research shows that on the infrastructure, Bars play a pivotal role, compared with the content kept on other malicious or compromised sites, possibly because they are hosted on popular cloud services, and therefore hard to blacklist and also easy to share across different campaigns. Also, in a malicious campaign, the adversary may take advantage of multiple Bars, at different attack stages, to construct a complicated infrastructure that supports her mission (Section 4.1). More importantly, we discovered that the adversary effectively exploited misconfigured legitimate buckets to infect a large number of their front-end web services (Section 4.2), and the cloud providers have not done much to counteract the threat, often leaving Bars there for a long time (Section 4.3), possibly due to the privacy constraints and limited means to detect individual components of a malicious activity. Such observations, together with the challenge in blocking Bars, offer insights into the motivation for moving toward this new trend of repository-based attacks.

##### 4.1 Bar-based Malicious Web Infrastructure

**Landscape.** As mentioned earlier, BarFinder reported 730 suspicious repositories from 6885 cloud buckets over 20 cloud platforms. Among them, we utilized 694 confirmed Bars (through AV/blacklist scanning or manual validation, see Section 3.3) for the measurement study. These Bars were found to directly serve 156,608 domains

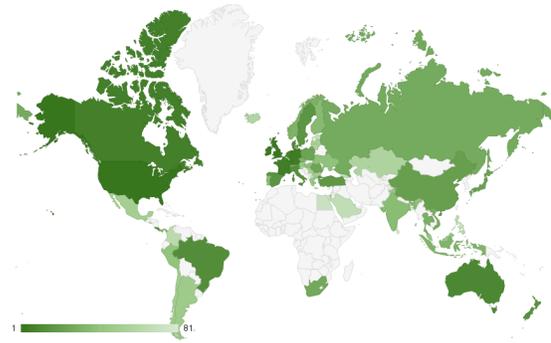


Figure 6: Impact of Bars' front-end websites around the globe.

(i.e., front-end websites), through which they are further attached to 6,513,519 redirection paths involving 166,772 domains. Figure 5 illustrates the number of Bars we found on different cloud platforms. Among them, Amazon S3 is the most popular one in our dataset, hosting the most Bars (45%), which is followed by CloudFront (Amazon's CDN) 25.1% and Akamaihd 9.3%. Note that of these 20 clouds, seven of them provide free storage services (e.g., 15GB free space on Google Drive, 5GB for Amazon S3), and therefore easily become the ideal platforms for low-budget miscreants to distribute their illicit content. Also, eleven of them support HTTPS, on which malicious activities are difficult to catch by existing signature-based intrusion detection systems like snort and Shadow[?][?]. Interestingly, on some of the most prominent platforms, the miscreants are found to take advantage of the cloud providers' reputations to make their Phishing campaigns look more credible: for example, we found that the adversary continuously spoofed Gmail's login page on Google Drive, and the software download page for Amazon FireTV in an Amazon S3 bucket.

Figure 6 shows the distribution of Bars' frontend websites across 81 countries, as determined by the geolocations of the sites. The number of Bars' frontend sites in each country is ranked and described with different levels of darkness in the figure. We observe that most of these frontends stay in United States (14%), followed by Germany (7%) and United Kingdom (5%).

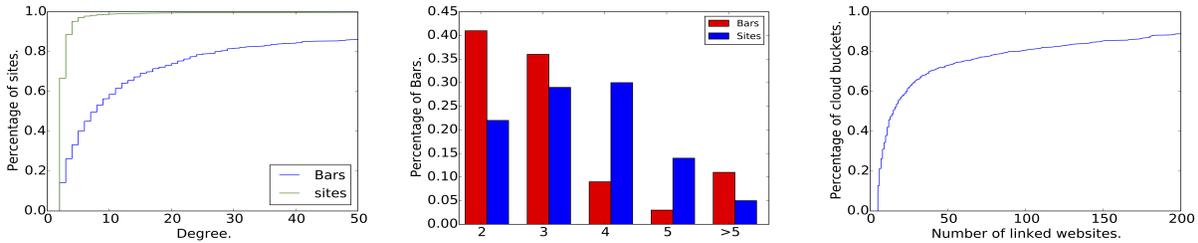
**Role in attack infrastructures.** Actually, most nodes on a malicious infrastructure are the malicious websites with newly registered domains and those that are compromised. To better understand the critical roles of Bars, we compared those nodes with the bad cloud buckets. Specifically, we first identified both types of nodes from the redirection paths and then analyzed the number of unique paths each member in either category is associated with and the position of the member on the path. Figure 7(a) presents the cumulative distribution of the paths going through a Bar and that of a compromised or malicious site. As seen in the figure, compared with other nodes on the infrastructure, Bars clearly sit on much more paths (47.4 on average vs. 8.6), indicating their importance.

Further, Figure 7(b) shows the histogram of position distributions (again, Bars vs. bad sites). The observation is that more Bars (41%, 11%) show up at the beginnings and the ends of the paths than bad websites (22%, 5%), which demonstrates that they often act as first-hop redirectors or attack-payload repositories. For example, in our three-month-long monitoring of the campaign based on the Spyware distribution Bar [akamaihd.net\\_rvar-a](http://akamaihd.net_rvar-a), we found that besides the Bar, 320 newly-registered websites participated in the attack; here the Bar acted very much like a dispatcher: providing JavaScript that identified the victim's geolocation and then using an iframe to redirect her to a selected bad site.

**Content sharing.** Our research reveals that Bars have been ex-

Table 5: Top 10 most popular Bars.

| Rank | Cloud bucket  | # of front-end sites | Avg path len | Popularity |
|------|---|----------------------|--------------|------------|
| 1    | s3.amazonaws.com_content.sitezoogole.com              | 4,429                | 2.9          | 2.8%       |
| 2    | cloudfront.net_d3n8a8pro7vhmx                         | 1,829                | 3.3          | 1.4%       |
| 3    | s3.amazonaws.com_assets.ngin.com                      | 1,643                | 3.2          | 1.2%       |
| 4    | s3.amazonaws.com_publisher_configurations.shareaholic | 1,434                | 2.7          | 0.9%       |
| 5    | cloudfront.net_d2e48lfsb5exy                          | 1,340                | 4.0          | 0.9%       |
| 6    | cloudfront.net_d1t3gia0in9tdj                         | 1,297                | 3.2          | 0.9%       |
| 7    | cloudfront.net_d2i2wahzwrmln5                         | 1,249                | 2.5          | 0.8%       |
| 8    | cloudfront.net_d202m5krfqbp5                          | 1,062                | 2.8          | 0.8%       |
| 9    | s3.amazonaws.com_files.enjin.com                      | 1,020                | 7.1          | 0.7%       |
| 10   | akamaihd.net_cdncache3-a                              | 976                  | 6.4          | 0.6%       |



(a) Cumulative distribution of degrees per site. (b) Percentage of Bars in each position of redirection path (Ignoring those of in-degrees per Bar, traces with length of 2). (c) Cumulative distribution of number of linked websites.

Figure 7: Bars play critical roles in attack infrastructures.

tensively shared among malicious or compromised websites, also across different positions on malicious redirection chains. Figure 7(c) illustrates the cumulative distribution of Bars’ in-degrees in their individual redirection graphs: that is, the number of the sites utilizing these Bars. On average, each Bar shows up on 252 sites and 12% of them are used by more than 200 websites. Table 5 lists the 10 most popular Bars we found. Among them, eight, including [s3.amazonaws.com\\_content.sitezoogole.com](https://s3.amazonaws.com_content.sitezoogole.com), [s3.amazonaws.com\\_publisher\\_configurations.shareaholic](https://s3.amazonaws.com_publisher_configurations.shareaholic), etc., host services for website generation, blackhat SEO or Spam. Particularly, [akamaihd.net\\_cdncache3-a](https://akamaihd.net_cdncache3-a) turns out to be a distributor of Adware, whose scripts are loaded into the victim’s browser to redirect it to other sites for downloading different Adware. Also, we found that another Bar [s3.amazonaws.com\\_files.enjin.com](https://s3.amazonaws.com_files.enjin.com) hosts exploits utilized by 1,020 bad sites. Finding Bars can help to effectively detect more sites with malicious contents.

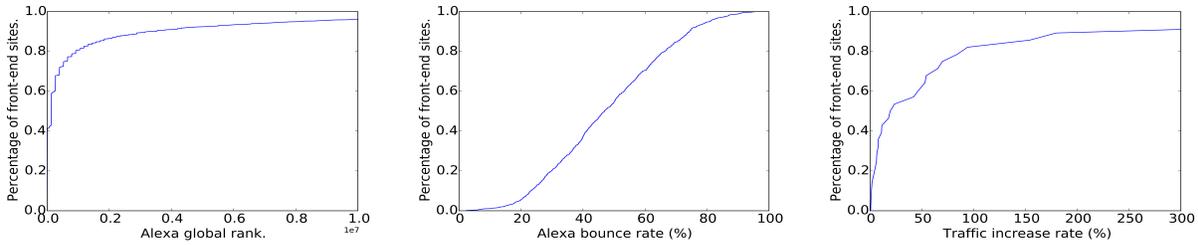
Another interesting observation is that malicious content is also extensively shared across different Bars. To understand such content reuse, we grouped the malicious programs retrieved from different Bars based on the similarity of their code in terms of edit distance. Specifically, we removed the spaces within the programs and ran the Python library *scipy.cluster.hierarchy.linkage* [?] to hierarchically cluster them (now in the form of strings) according to their Jaro scores [?]. In this way, we were able to discover three types of content sharing: intra-bucket sharing, cross-bucket sharing, and cross-platform sharing. Specifically, within the Amazon bucket [akamaihd.net\\_asrv-a](https://akamaihd.net_asrv-a), we found that many of its cloud URLs are in the form of [http://asrv-a.akamaihd.net/sd/\[num\]/\[num\].js](http://asrv-a.akamaihd.net/sd/[num]/[num].js). The JavaScript code turns out to be all identical, except that each script redirects the visitor to a different website. The similar code also appears in another Amazon bucket [akamaihd.net\\_cdncache-a](https://akamaihd.net_cdncache-a). As another example, we discovered the same malicious JavaScript ([JS.ExploitBlacole.zm](https://JS.ExploitBlacole.zm)) from the Bars on CloudFront and Qiniudn respectively, even under the same path (i.e., [media/system/js/modal.js](https://media/system/js/modal.js)). Moreover, we found that attackers used sub-domain generation algorithm to automatically generate sub-domain for Bars, then further

reused the same malicious contents for these Bars. Specifically, we found that 28 content sharing Bars on Akamaihd have the same format in their names. Attackers utilized a word bank based sub-domain generation algorithms [?], which concatenates fixed terms and a series of domain names (remove dot), then truncates the string if its length is over 13, e.g., [apismarterpoweru-a](https://apismarterpoweru-a) (truncated from [smarterpowerunite.com](https://smarterpowerunite.com)). The common patterns of Bars indicate the potential of developing an accurate detection procedure.

**Correlation.** We further studied the relationships between different Bars, fetched by the same websites. From our dataset, 11,442 (3.5%) websites are found to access at least two Bars. Among them, 8,283 were served as front-end websites, and 3,159 other sites on redirection chains. Also, 60.9% of these sites link to the repositories on the same cloud platforms and 39.1% use those on different platforms. In some cases, two buckets are often used together. For example, we found that a click-hijacking program was separated into the code part and the configuration part: the former is kept on CloudFront while the latter is on Akamaihd; the two buckets always show up together on redirection chains. Such a separation seems to be done deliberately, in an attempt to evade detection. Also we saw that Bars carrying the same attack vectors are often used together, which are apparently deliberately put there to serve parties of the same interests: as another example, a compromised website was observed to access four different Bars on different cloud platforms, redirecting its visitors to different places for downloading Adware to the visitor’s system. Our findings show that Bars are widely deployed in attacks and serve in a complex infrastructure.

## 4.2 Bucket Pollution

**Polluted repositories.** To find polluted buckets, we searched the Alexa top 20K websites for the Bars in our dataset and 276 Bars were found. When a legitimate site links to a Bar, the reason might be either the website or the repository is hacked. Differentiating these two situations with certainty is hard, and in some cases, it may not be possible. All we could do is to get an idea about the prevalence of such bucket pollution, based on the intuition that if a website is less



(a) Cumulative distribution of Alexa global ranks per Bars' front-end sites. (b) Cumulative distribution of Alexa bounce rate per Bar's front-end sites. (c) Cumulative distribution of traffic increase rate per Bar's front-end sites.

Figure 8: Alexa global rank, bounce rate and traffic increase rate of Bar's front-end websites.

```
GET /?delimiter=/ HTTP/1.1
Host: (bucket-name).s3.amazonaws.com
Accept-Encoding: identity
content-length: 0
Authorization: AWS (access key):(secret key)
```

Figure 9: Constructed request header.

vulnerable, then it is less likely to be compromised. To this end, we ran WhatWeb, a powerful web vulnerability scanner, on these sites and found 134 Bar's front-end websites contain various flaws, such as using CMS in vulnerable version (e.g. wordpress 3.9), vulnerable plugins (e.g., JCE Extension 2.0.10) and vulnerable software (e.g., Apache 2.2). The remaining 142 Bar's front-end websites look pretty solid in web protection and therefore it is likely that the Bars they include were polluted. This set of potentially compromised buckets takes 19% of all the Bars flagged by BarFinder. These buckets, together with the additional 30 randomly sampled from the set, went through a manual analysis, which shows that indeed they were legitimate buckets contaminated with malicious content.

**Misconfiguration and impact.** It is even more challenging to determine how these buckets were compromised, which could be caused by exploiting either the cloud platform vulnerabilities or the bucket misconfigurations. Without an extensive test on the cloud platform and the repositories, which requires at least direct access to them, a comprehensive study on the issue is impossible. Nevertheless, we were able to identify a misconfiguration problem widely existing in popular buckets. *This flaw has never been reported before but was likely known to the underground community and has already been utilized to exploit these repositories. We reported the flaws to the vendors and they confirmed our finding.*

Specifically, on Amazon S3, one can configure the access policies for her bucket to defines which AWS accounts or groups are granted access and the type of access (i.e., list, upload/modify, delete and download): this can be done through specifying access control list on the AWS Management Console. Once this happens, the cloud verifies the content of the `authorization` field within the client's HTTP request header before the requested access is allowed to go through. However, we found that by default, the policy is not in place, and in this case, the cloud only checks *whether the authorization key (i.e., access key and secret key) belongs to an S3 user, not the authorized party for this specific bucket*: in other words, anyone, as long as she is a legitimate user of the S3, has the right to upload/modify, delete and list the resources in the bucket and download the content. Note that this does not mean that the bucket can be directly touched through the browser, since it does not put anything into the authorization field. However, the

adversary can easily build his own HTTP header, filling in his own S3 key, as illustrated in Figure 9, to gain access to the misconfigured repository. In our research, we verified that all such operations can be performed on any repositories with the configuration flaw, which suggests that site operators need to take more caution when setting the configuration rules.

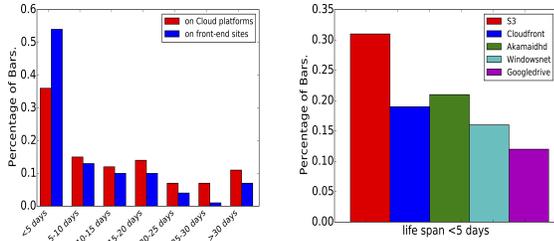
To understand the impact of this problem, we developed a simple web testing tool, which checked a bucket's configuration using our own S3 key. By scanning all 6,885 repositories (including both Bars and legitimate buckets), we discovered that 472 are vulnerable, which were associated with 1,306 front-end websites. The Alexa global ranks and the bounce rates of their front-end websites are illustrated in Figure 8(a) and Figure 8(b). 63% of them have bounce rates from 20% to 60%; 9 sites are ranked within Alexa top 5000 (e.g., [groupon.com](http://groupon.com), [space.com](http://space.com)).

Focusing on the 104 bad buckets with the flaws, we further manually sampled 50 and confirmed that these buckets were indeed legitimate, including high-profile ones like [s3.amazonaws.com\\_groupon](http://s3.amazonaws.com/groupon). Further, looking into the these buckets' file uploading time (retrieved from the buckets through the flaw), we found that in some cases, the attack has been there for six years. Particularly the Amazon bucket [s3.amazonaws.com\\_groupon](http://s3.amazonaws.com/groupon), Groupon's official bucket, was apparently compromised five times between 2012 and 2015 (see Section 4.4 for details), according to the changes to the bucket we observed from the bucket historical dataset we collected from [archive.org](http://archive.org). We also estimated the volume of traffic to those Bar-related sites using a PassiveDNS dataset [?], which contains DNS lookups recorded by the Security Information Exchange. Figure 8(c) illustrates the traffic of the websites during the time period when their buckets were compromised, which was increased significantly compared with what those sites received before their compromise, indicating that they likely received a lot of visits. This provides evidence that the impact of such compromised buckets is indeed significant.

### 4.3 Lifetime and Evasion

In the presence of the severe threat from Bars, we found that cloud providers' responses, however, are far from adequate. This is highlighted by the relatively long lifetimes of malicious repositories we observed.

**Lifetime.** To understand the duration of Bars' impacts, we continuously crawled the front-end bad sites every five days to check whether they were still using the same set of Bars, and also malicious cloud URLs to find out whether the repositories were still alive. Figure 10(a) illustrates the distributions of such bad repositories' life spans within those front-end sites and on cloud platforms. As can be seen in the figure, on average, the references of these Bars on the websites were removed much faster than their cloud URLs and



(a) Distributions of Bars' life spans on front-end sites and moved within 5 days in top 5 cloud platforms. (b) Percentage of Bars' life spans on front-end sites and moved within 5 days in top 5 cloud platforms with most Bars.

Figure 10: Lifetime of Bars.

ultimately their accounts on the cloud platforms. Apparently, the cloud providers are less aggressive, relative to the website owners, in addressing Bar-related infections. In Figure 10(b), we further compare Bars' life spans on different platforms: interestingly, with more bad buckets on its servers, Amazon AWS acted more promptly than other clouds; Google, however, moved much slower: for example, on Google Drive, a repository hosting malware-serving pages, [googledrive.com\\_0B8D1eUrPT\\_z3OVpBTvJ3LUg2UEk](https://googledrive.com_0B8D1eUrPT_z3OVpBTvJ3LUg2UEk), stayed there for over 150 days, longer than the average duration of other exploit servers (non-cloud) reported by the prior work [?][?] (2.5 hours). The observation indicates that cloud providers have noticed such problem, but a likely lack of effective methods to identify and clean Bars.

**Evasion.** Such a long lifetime could be related to a spectrum of evading techniques the adversary deploys to protect his cloud assets, which are described as follows:

- **Content separation.** Apparently, the adversary tends to break his attack assets into pieces and store them at different places. As mentioned earlier, we found that malware's code and configuration files were placed in different buckets. Also, we discovered in this study that there are 32 Bars that host nothing but images used in various attacks, Phishing and Fake AV campaigns in particular. Since the images themselves do not contain any malicious code, these repositories typically stay on the clouds for a long time, >30 days on average.
- **Content change.** Another interesting observation is that the malicious content within Bars changes over time, in an attempt to avoid being linked to blacklisted malicious websites. Specifically, looking into the history of the content (from *archive.org*) retrieved from the Bar through the same cloud URL, we found that part of the content (e.g., the destination of a redirection) changes frequently, moving quickly away from known malicious sites.
- **Redirect cloaking.** Like malicious or compromised websites, Bars are also found to leverage cloaking techniques (rendering different content based on the visitor's cookie, IP, user agent, etc.) to avoid detection. However, different from websites, cloud hosting services typically do not support server-side scripting. As a result, Bars have to run the cloaking code on the client (browser) side, which makes the evasion less stealthy. To make up for this weakness, the adversary was observed to place redirection websites in front of Bars, running cloaking techniques there to hide the repositories.
- **Obfuscation.** We found that the attack payloads in the repositories were often obfuscated. Various kinds of obfuscation techniques were found from simply Base64 encoding to online obfuscation tools (e.g., [api.myobfuscate.com](https://api.myobfuscate.com)). Actually, even the links to refer

Table 6: Comparison of Bars' lifetime under different evasion techniques.

| Evasion technique  | # of Bars | # of front-end sites | Avg. life span |
|--------------------|-----------|----------------------|----------------|
| Content separation | 10        | 743                  | 25-30 days     |
| Content change     | 10        | 1045                 | > 30 days      |
| Redirect cloaking  | 10        | 1220                 | 10-15 days     |
| Obfuscation        | 10        | 1032                 | 10-15 days     |
| None               | 10        | 984                  | 5-10 days      |

to Bars within front-end websites were obfuscated in some cases, apparently, for the purpose of protecting the repositories.

Further, our study shows that these techniques were also utilized together to make identification of Bars even harder. Specifically, we manually choose 10 Bars with each evasion technique (40 in total), combined with 10 Bars without evasion technique, and then compare their life spans. It is clear that evasion techniques do allow Bars to hide longer, as illustrated in Table 6.

## 4.4 Case Studies

In this section, we discuss two prominent examples.

**PUP campaign.** Our study reveals a malicious web campaign dubbed *Potentially Unwanted Programs* (PUP) distribution: the attack redirects the victim to an attack page, which shows her fake system diagnosis results or patch requirements through the images fetched from a Bar, in an attempt to cheat the victim into downloading "unwanted programs" such as Spyware, Adware or a virus. This campaign was first discovered in our dataset. Altogether, at least 11 Bars from 3 different cloud platforms and 772 websites (not hosted on the cloud) were involved in.

Through analyzing the redirection traces of the campaign, we found that two Akamai Bars, [akamai.net\\_cdncache3-a](https://akamai.net_cdncache3-a) and [akamaihd\\_asrv-a](https://akamaihd_asrv-a), frequently inject scripts into compromised websites, which serve as first-hop redirectors to move a visitor down the redirection chain before hitting malicious landing pages (that serve malicious content). Interestingly, all the follow-up redirectors are compromised or malicious websites that are not hosted on the cloud. The scripts in the Bars were found to change over time, redirecting the visitor to different next-hop sites (also redirectors). On average, the life span of such sites is only 120 hours, but the Bar was still alive when we submitted this paper. Such redirections end at at least 216 malicious landing sites, which all retrieve deceptive images from an Amazon S3 bucket [s3.amazonaws.com\\_cicloudfront](https://s3.amazonaws.com_cicloudfront) (a Bar never reported before and is still alive). An example is a system update warning, as shown in Figure 1. From the repository, we collected 134 images, including those for free software installation, updates on all mainstream OSes, browsers and some popular applications. If she clicks and downloads the program promoted on the site, the code will be fetched from multiple Bars, such as [s3.amazonaws.com\\_wbt\\_media](https://s3.amazonaws.com_wbt_media) where the PUP puts a Bitcoin miner on the victim's system, and [cloudfront.net\\_d12mrm7igk59vq](https://cloudfront.net_d12mrm7igk59vq), whose program modifies Chrome's security setting.

**Groupon Bar.** We discovered that a misconfigured Amazon S3 bucket [s3.amazonaws.com\\_groupon](https://s3.amazonaws.com_groupon) belongs to Groupon (Alexa global rank 265), a global e-commerce marketplace serving 48.1 million customers worldwide. The bucket was used as the resource repository for Groupon's official website (i.e., [groupon.com](https://groupon.com)) as well as its marketing sites (12 websites observed in our dataset). When tracking its historical content from *archive.org*, we were surprised to see that the Groupon S3 bucket has been compromised at least eight times in the past five years (e.g., 2015/08/06, 2014/12/18,

2014/06/25, 2014/01/27, 2014/02/26, 2013/06/23, 2011/11/08, 2010/09/28). These attacks caused different types of malicious payloads to be uploaded to their repository, including Adware, Trojan, virus and others. Even though the bucket owner changed the access control policy in 2012 to prevent the unauthorized party from directly listing the bucket content through browser, it remained accessible by our tool mentioned in Section 4.2, which constructs an *Authorization* field in HTTP header, and unauthorized listing, upload and even modification can still occur.

## 5. DISCUSSION

**Limitations of BarFinder.** As mentioned earlier, Bar detection is hard, since cloud repositories cannot be directly accessed by the parties outside the cloud. Therefore, the goal of BarFinder is to leverage the sites served by Bars to find suspicious repositories. For this purpose, we chose to utilize the *collective* features of these sites, such as their topological relations, content shared across sites, etc. This strategy could make the approach more robust, as the collective features are more difficult to evade compared with those from individual sites. On the other hand, it requires that the party running the system first makes efforts to gather the sites using cloud buckets, the more the better. Further, there are repositories that only serve a small set of front-end sites: e.g., we found that among the Alexa top 3K sites, 67 sites are connecting to the cloud buckets only used by themselves. Those “self-serving” buckets are rather popular in reputable websites such as [appsspot.com\\_android-site](http://appsspot.com_android-site) only used by [android.com](http://android.com), [s3.amazonaws.com\\_ttv-backgroundart](http://s3.amazonaws.com_ttv-backgroundart) only used by [twitch.tv](http://twitch.tv), etc. This fact makes the bad apples among them hard to catch by BarFinder simply because not enough sites using them are out there to allow us to differentiate these two types of repositories. Detection techniques covering this type of Bars need to be developed in the follow-up research.

**Other defenses against Bars.** Besides the detection effort made by the third party, as BarFinder does, more can be done to mitigate the threats posed by Bars, from the ends of the website owner, the bucket owner and the service provider. The website owner could perform integrity checks on the resources her website retrieves from the bucket, making sure that it is not compromised. The cloud bucket owner should carefully configure her cloud bucket to avoid the issue we found and other misconfiguration flaws. In this case, an automatic configuration checker could be helpful. Most importantly, the cloud provider does have the responsibility to move more aggressively on detecting and removal of Bars from their systems. This, however, is non-trivial, given the privacy concern and the fact that some Bars can only be considered to be malicious by looking at the malicious activities they are involved in, such as those hosting Phishing pictures. Further research is needed to better understand what the provider can do to address the issue.

**Ethical issues.** Most findings of the paper were made through analyzing the data crawled from the public domain. Regarding the study on the misconfiguration problem we found, our scanner was designed to minimize the privacy impacts on vulnerable repositories: specifically, it only tried out the functionality like file listing, uploading and downloading. The impact of such operations are very much in line with those of running online web testing tools (e.g., Sucuri [?]) on others’ websites. Most importantly, we did this with the full intention to protect such repositories from future exploits, and also carefully avoided changing any existing content there and deleted from our system all the files downloaded. Further, we have already contacted the major vendors such as Groupon and the cloud providers like Amazon about those security breaches, and will continue to notify others and help them fix the configuration

problem. So far, Groupon has acknowledged the importance of our findings and expressed gratitude for our help.

## 6. RELATED WORK

**Bad site detection.** Malicious web activities have been extensively studied [?][?][?][?]. Most related to our work here is the use of HTML content and redirection paths to detect malicious or compromised websites. Examples for the content-based detection include a DOM-based clustering systems for monitoring Scam websites [?], classification of websites for predicting whether some of them will turn bad based on the features extracted from HTML sources, and a monitoring mechanism [?] (called Delta) to keep track of the changes made to the content of a website for detecting script-injection campaigns. For those using malicious redirection paths, prominent prior approaches use short redirection sequences to capture bad sites [?], unique properties of malicious infrastructure (its density) for detecting drive-by downloads [?] or malware distribution [?] and a trace-back mechanism that goes through the redirection paths [?] for labeling malware download sites. Compared with those prior studies, which all rely on the properties of the targets they try to capture, BarFinder utilizes the features found from the front-end websites using cloud buckets, as those repositories may not be directly accessible. Also, our approach leverages a set of unique collective features, based on the connected components of a graph, which, to our knowledge, has never been used before.

**Cloud security.** Previous studies on security and privacy issues in cloud storage primarily focus on the confidentiality of the data stored in the cloud or the attacks targeting the cloud computing infrastructure. Examples include the study on co-locating attack virtual machines (VM) with the target one on Amazon EC2 [?], which enables a cache-based side-channel attack to infer sensitive user information from the target [?], and the work on controlled-channel attacks on multi-user cloud hosting environment [?], which allows an untrusted VM to extract sensitive information from protected applications. More recently, attention has moved to abuse of cloud-based services for fraudulent activities. For example, prior research [?] analyzed Dropbox client software and discovered that it can be exploited to hide files with unlimited storage capacity. Additionally, [?] studied the use of the Amazon EC2 to host malicious domains acting as command and control centers, exploit servers by downloading malware samples and executing them in sandbox environments to analyze their interactions with the cloud. [?] studied the long-tail SEO spam on cloud platforms and measured its effectiveness. Our study differs from these works by proposing BarFinder to identify malicious cloud repositories and provide an in-depth analysis of the use of cloud repositories in malicious campaigns and how they correlate with the websites they serve.

In another study, researchers [?] inspected the fraudulent traffic to cloud-hosted pages for the purpose of squandering the user’s resources and raising her cloud-usage cost. They also proposed detection methods based on the consistency of the requests. Unlike these works, our research investigated the abuse of cloud bucket as a malicious service, an emerging new cloud-based security threat that has never been studied before.

## 7. CONCLUSION

The emergence of using cloud repositories as a malicious service presents a new challenge to web security. This new threat, however, has not been extensively studied and little is known about its scope and magnitude and the techniques the adversary employs. In this paper, we report the first systematic study on malicious and compromised cloud repositories and the illicit online activities built

around them. We collected a small set of seeding Bars and identified a set of collective features from the websites connecting to them. These features describe the effort made by the adversary to protect Bars and utilize them to quickly build up a large campaign. Using these features, we developed a new scanner that detected over 600 Bars on top-of-the-line cloud platforms, including Google, Amazon, and others. Over these Bars, we performed a large-scale measurement study that led to surprising findings of such attacks. Examples include the central roles those buckets play at each stage of a web attack (redirection, displaying Phishing content, exploits, attack payload delivery, etc.), the strategy to separate malware code and configuration files to avoid detection, and a configuration flaw never reported before that was likely exploited to compromise many cloud buckets. Our findings made an important step toward better understanding and effective mitigating of this new security threat.

## 8. ACKNOWLEDGMENT

This work was supported in part by the National Science Foundation (grants CNS-1223477, 1223495, 1527141 and 1618493); We thank our anonymous reviewers for their useful comments.

## 9. REFERENCES

- [1] Buckets. [https://cloud.google.com/storage/docs/json\\_api/v1/buckets](https://cloud.google.com/storage/docs/json_api/v1/buckets). [Online].
- [2] Servnet. <https://servnetshstndci.onion>. [Online].
- [3] A. authors. Details omitted for double-blind reviewing.
- [4] C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [5] K. Borgolte, C. Kruegel, and G. Vigna. Delta: automatic identification of unknown web-based infection campaigns. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 109–120. ACM, 2013.
- [6] BuiltWith. Builtwith. <http://builtwith.com/>, 2015. [Online].
- [7] Clean-MX. Clean mx realtime database. <http://support.clean-mx.de/clean-mx/viruses.php>, 2015. [Online].
- [8] W. Cohen, P. Ravikumar, and S. Fienberg. A comparison of string metrics for matching names and records. In *Proceedings of Kdd workshop on data cleaning and object consolidation*, 2003.
- [9] C. Crawl. Common crawl. <https://commoncrawl.org/>, 2015. [Online].
- [10] damballa. Dgas in the hands of cyber-criminals:examining the state of the art malware evasion techniques. [https://www.damballa.com/downloads/r\\_pubs/WP\\_DGAs-in-the-Hands-of-Cyber-Criminals.pdf](https://www.damballa.com/downloads/r_pubs/WP_DGAs-in-the-Hands-of-Cyber-Criminals.pdf), 2015. [Online].
- [11] M. F. Der, L. K. Saul, S. Savage, and G. M. Voelker. Knock it off: profiling the online storefronts of counterfeit merchandise. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1759–1768. ACM, 2014.
- [12] DNSDB. Passivedns. <https://www.dnsdb.info/>, 2015. [Online].
- [13] Google. Google hosted libraries. <https://developers.google.com/speed/libraries/?csw=1>, 2015. [Online].
- [14] Google. Publish website content. <https://developers.google.com/drive/web/publish-site>, 2015. [Online].
- [15] C. Grier, L. Ballard, J. Caballero, N. Chachra, C. J. Dietrich, K. Levchenko, P. Mavrommatis, D. McCoy, A. Nappa, A. Pitsillidis, et al. Manufacturing compromise: the emergence of exploit-as-a-service. In *Proceedings of the 2012 ACM conference on Computer and communications security*, pages 821–832. ACM, 2012.
- [16] X. Han, N. Kheir, and D. Balzarotti. The role of cloud services in malicious software: Trends and insights. In *DIMVA 2015, 12th Conference on Detection of Intrusions and Malware & Vulnerability Assessment, July 9-10, 2015, Milan, Italy*, Milan, ITALY, 07 2015.
- [17] J. Idziorek, M. Tannian, and D. Jacobson. Detecting Fraudulent Use of Cloud Resources. In *Proc. 3rd ACM Workshop on Cloud Computing Security Workshop*, Chicago, IL, Oct. 2011.
- [18] L. Invernizzi, P. M. Comparetti, S. Benvenuti, C. Kruegel, M. Cova, and G. Vigna. Evilseed: A guided approach to finding malicious web pages. In *Security and Privacy (SP), 2012 IEEE Symposium on*, pages 428–442. IEEE, 2012.
- [19] L. Invernizzi, S. Miskovic, R. Torres, S. Saha, S. Lee, M. Mellia, C. Kruegel, and G. Vigna. Nazca: Detecting malware distribution in large-scale networks. In *Proceedings of the Network and Distributed System Security Symposium (NDSS)*, 2014.
- [20] Z. Li, S. Alrwais, X. Wang, and E. Alowaisheq. Hunting the red fox online: Understanding and detection of mass redirect-script injections. In *Security and Privacy (SP), 2014 IEEE Symposium on*, pages 3–18. IEEE, 2014.

- [21] 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 *Security and Privacy (SP), 2013 IEEE Symposium on*, pages 112–126. IEEE, 2013.
- [22] X. Liao, C. Liu, D. Mccoy, E. Shi, and R. Beyah. Characterizing long-tail seo spam on cloud web hosting services. In *Proceedings of the International World Wide Web Conference*, 2016.
- [23] 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.
- [24] M. Mulazzani, S. Schrittwieser, M. Leithner, and M. Huber. Dark Clouds on the Horizon: Using Cloud Storage as Attack Vector and Online Slack Space. In *Proc. 20th USENIX Security Symposium*, San Francisco, CA, Aug. 2011.
- [25] T. Nelms, R. Perdisci, M. Antonakakis, and M. Ahamad. Webwitness: Investigating, categorizing, and mitigating malware download paths. In *24th USENIX Security Symposium (USENIX Security 15)*, pages 1025–1040, Washington, D.C., Aug. 2015. USENIX Association.
- [26] T. Ristenpart, E. Tromer, H. Shacham, and S. Savage. Hey, you, get off of my cloud: exploring information leakage in third-party compute clouds. In *Proceedings of the 16th ACM conference on Computer and communications security*, pages 199–212. ACM, 2009.
- [27] Scipy. [scipy.cluster.hierarchy.linkage](http://docs.scipy.org/doc/scipy/reference/generated/scipy.cluster.hierarchy.linkage.html). <http://docs.scipy.org/doc/scipy/reference/generated/scipy.cluster.hierarchy.linkage.html>, 2015. [Online].
- [28] Sklearn. [sklearn.svm.svc](http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html). <http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>, 2015. [Online].
- [29] Snort. Snort ssl and tls. <http://manual.snort.org/node147.html>, 2015. [Online].
- [30] solutionary. Threat-intelligence. [https://www.solutionary.com/\\_assets/pdf/research/ser-t-q4-2013-threat-intelligence.pdf](https://www.solutionary.com/_assets/pdf/research/ser-t-q4-2013-threat-intelligence.pdf), 2015. [Online].
- [31] G. Stringhini, C. Kruegel, and G. Vigna. Shady paths: Leveraging surfing crowds to detect malicious web pages. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 133–144. ACM, 2013.
- [32] Sucuri. Sucuri. <https://sucuri.net/>, 2015. [Online].
- [33] Symantec. The future of ids. <http://www.symantec.com/connect/articles/future-ids>, 2015. [Online].
- [34] VirusTotal. Virustotal. <https://www.virustotal.com/>, 2015. [Online].
- [35] WhatWeb. Whatweb. <http://www.morningstarsecurity.com/research/whatweb>, 2015. [Online].
- [36] Y. Xu, W. Cui, and M. Peinado. Controlled-channel attacks: Deterministic side channels for untrusted operating systems. In *Proceedings of the 36th IEEE Symposium on Security and Privacy (Oakland)*. IEEE ã\$ Society of Electrical and Electronics Engineers, May 2015.
- [37] Y. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart. Cross-Tenant Side-Channel Attacks in PaaS Clouds. In *Proc. 21st Conference on Computer and Communications Security (CCS)*, Scottsdale, AZ, Oct. 2014.

## Appendices

### A List of cloud hosting platforms

Table 7: List of cloud hosting platforms.

| Cloud Platform | Domain            |
|----------------|-------------------|
| heroku         | herokuapp.com     |
| amazon S3      | s3.amazonaws.com  |
| cloudfront     | cloudfront.net    |
| windowsnet     | windows.net       |
| azure          | azurewebsites.net |
| google         | googledrive.com   |
| appspot        | appspot.com       |
| msecdn         | msecdn.net        |
| bitbucket      | bitbucket.org     |
| github         | github.io         |
| sina           | sinaapp.com       |
| olymp          | olymp.ee          |
| rackcdn        | rackcdn.com       |
| baiduyun       | duapp.com         |
| qiniu          | qiniucdn.com      |
| akamaihd       | akamaihd.net      |
| yahoo          | hostingprod.com   |
| sogo           | sogoucdn.com      |
| go2cloud       | go2cloud.org      |
| aliyun         | aliyun.com        |