# Method of Generating a Blacklist for Mobile Devices by Searching Malicious Websites

Takashi Ishihara, Masaya Sato, Toshihiro Yamauchi
Graduate School of Natural Science and Technology,
Okayama University, Okayama, Japan
Email: {sato, yamauchi}@cs.okayama-u.ac.jp

*Abstract*—As mobile devices have become more popular, malware and attacks directed at them have significantly increased. One of the methods to attack mobile devices is redirecting a user to unwanted websites by unwanted page transition. One of the countermeasures against such attacks is to generate a blacklist of URLs and hostnames, which can prevent access to malicious websites. To generate a blacklist, first, malicious websites are collected in the web space. Then, URLs and hostnames of the malicious websites are added to the blacklist. However, URLs of the malicious websites are typically changed frequently; thus, it is necessary to keep track of the malicious websites and update the blacklist in a timely manner. In this study, we proposed a method to generate blacklists for mobile devices by searching malicious websites. The method collects many HTML files from the web space using a crawler and searches for HTML files that are highly likely to be malicious using keywords extracted from the known malicious websites to discover the new ones. Thus, new malicious websites can be added to the blacklist in a timely manner. Using the proposed method, we discovered malicious websites that were not detected by Google Safe Browsing. Moreover, the blacklist generated using the method had a high detection rate for certain malicious websites. This paper reports the design process and the results of the evaluation of the new method.

*Index Terms*—Malicious Websites, Blacklist, Web-based Attack, Android.

## I. INTRODUCTION

Mobile devices such as smartphones and tablets have become highly popular all around the world, and the number of mobile device users has increased by 124 million since 2019 and reached 67% of the world population [1] according to a survey published in January 2020. Mobile applications (apps) are increasingly used in everyday activities, and social media constitutes half of the time spent using mobile devices [1]. About 50% of webpage requests come from mobile devices [2]. Especially smartphones equipped with Google Android OS, which is a Linux-based open source operating system for smartphones, are widespread, and the market share of Android phones in 2020 will exceed 70% [3], [4].

As the number of users of mobile devices increases, mobile devices are more frequently targeted in cyberattacks. According to the reports of the damage investigation for mobile malware targeting Android devices, the number of users affected by mobile malware is considerably increasing [5]. Web-based attacks in mobile devices have also been reported [6].

One of the methods to attack mobile devices is redirecting a user to unwanted websites via an unwanted page transition. In this attack, when the user accesses a malicious website that occurs the transition (landing site), the transition to the website occurs automatically or when the user taps on the screen. In addition, it redirects the user to the target website (unwanted website) after passing through multiple websites (intermediate site). Social media platforms are increasingly used by attackers [7] who lead the users to malicious websites via social media posts or messages [8]. In some cases, illegal advertisements direct users to malicious websites [9]. Malicious websites display fake alerts such as virus infection messages and try to make users install suspicious apps for virus removal [10]. In addition, some phishing websites trick the users into entering their login information [8].

There are anti-phishing tools that are designed for the attacks targeting desktop devices. However, the majority of these tools cannot effectively address the phishing attacks on mobile devices [11]. Therefore, it is necessary to take countermeasures against the malicious attacks directed at mobile devices.

One of these countermeasures is to use a blacklist of URLs and hostnames. This measure can be expected to prevent access to malicious websites that are the origin of the redirect, or prevent access to malicious websites that are the redirect destination. However, an attacker may change the IP address or the domain name of the malicious website within a short period of time, making it difficult to take countermeasures using the blacklists. Therefore, it is necessary to search for the malicious websites in a timely manner and update the blacklist frequently.

In this study, we propose a method for searching malicious websites and generating a blacklist for mobile devices. We describe the implementation and evaluation of the proposed method and report the problems encountered during the design of the method.

In summary, our study makes the following contributions:

- The proposed method collects the web content by using crawlers and searches for the web content that is likely to be malicious by using keywords extracted from the existing malicious websites. In addition, it discovers a malicious website via manual analysis. Therefore, there is no need to prepare a data set or a detection method using machine learning tools.
- The URL of an unknown malicious website posted by an attacker can be discovered in a timely fashion by collecting the URLs from Twitter's Streaming API. Most
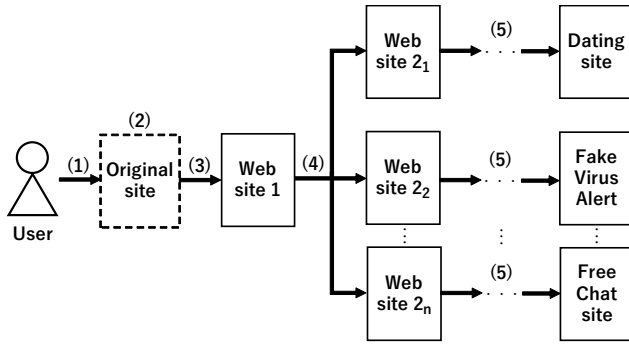
Fig. 1. Flow of an attack that redirects a user to unwanted websites [12]

malicious websites discovered by the proposed method were not detected by Google Safe Browsing.
- The proposed method analyzes the collected malicious websites using an Android device and generates a blacklist for mobile devices. Although Google Safe Browsing could not detect the attacks that redirected the users to unwanted websites, the blacklist generated by the proposed method had a high detection rate for certain malicious websites.

## II. ATTACKS OF REDIRECTING A USER TO UNWANTED WEBSITES

On mobile devices, attacks redirect the users to unwanted websites by unwanted page transition. Figure 1, adapted from [12], shows the flow diagram of an attack that redirects a user to unwanted websites. The steps of the attack are described as follows:

(1) Visiting the landing site (Original site)
(2) Tapping anywhere on the screen
The user taps anywhere on the landing site.
(3) Redirecting to intermediate site 1 (Website 1)
The landing site detects the user's tapping and redirects the user to intermediate site 1. Additionally, the user is forced to move to intermediate site 1 regardless of where they tap on the landing site.
(4) Redirecting to intermediate site 2 (Website 2)
Intermediate site 1 redirects the user to intermediate site 2. There are multiple intermediate sites of redirect destination. This redirection uses a JavaScript code, e.g., "window.location.replace." This JavaScript code can redirect the user to the specified URL without leaving a trace in the browser's history. This makes it impossible for the user to move to the previous website even if they click the back button.
(5) Redirecting to the unwanted website
The user is redirected to the unwanted website. The number of redirects is not the same every time. The user is either redirected to the unwanted website directly or via multiple intermediate sites.

According to this flow diagram, an attack that redirects a user to an unwanted website uses multiple malicious websites

such as the landing site, intermediate site, and unwanted website. In addition, the URL of intermediate site 2 is generated by executing the JavaScript code of intermediate site 1 by either of the following two operations [12].

(1) The JavaScript code of intermediate site 1 creates 10 URLs that are a combination of a specified URL with Base 36 strings created at random. Moreover, the redirecting destination from intermediate site 1 creates a random number and uses it to select the redirecting website from 10 URLs.
(2) The JavaScript code of intermediate site 1 creates one URL, and the number that is a part of this generated URL increases by one each time the user accesses intermediate site 1.

In some cases, redirection to unwanted websites does not occur. On the landing site, attackers set one of the following three conditions to make the redirection occur when the user taps anywhere on the landing site [12].

(1) Web browsers did not store the cookie of the landing site.
(2) The stated time set by attackers for each landing site has passed since the landing site was displayed.
(3) The user has not tapped on the landing site since it was displayed.

The unwanted websites include phishing sites and sites that make the user install apps [10]. On phishing sites, users' personal information is collected by prompting them to enter information such as their address and account number. Other sites make users install suspicious apps under the name of virus removal by displaying fake alerts such as virus infection messages on websites.

## III. METHODS OF GENERATING BLACKLISTS

### A. Purpose

Using blacklists of URLs and hostnames is an effective countermeasure against attacks from malicious websites and for preventing access to a landing site, an intermediate site, or an unwanted website. Therefore, the purpose of this research is to generate a blacklist of malicious websites for mobile devices.

Since malicious websites are continuously constructed, it is necessary to add new malicious websites to the blacklist in a timely manner. Therefore, we propose a method for searching malicious websites and generating a blacklist for mobile devices from these newly discovered websites.

### B. Problem

A blacklist is generated by searching for malicious websites and analyzing the discovered malicious websites. The problems with this method are explained as follows.

(Problem 1) Covering large-scale data on the web
It was reported in 2018 that there are more than 1.6 billion websites on the world wide web [13]. To search for malicious websites, it
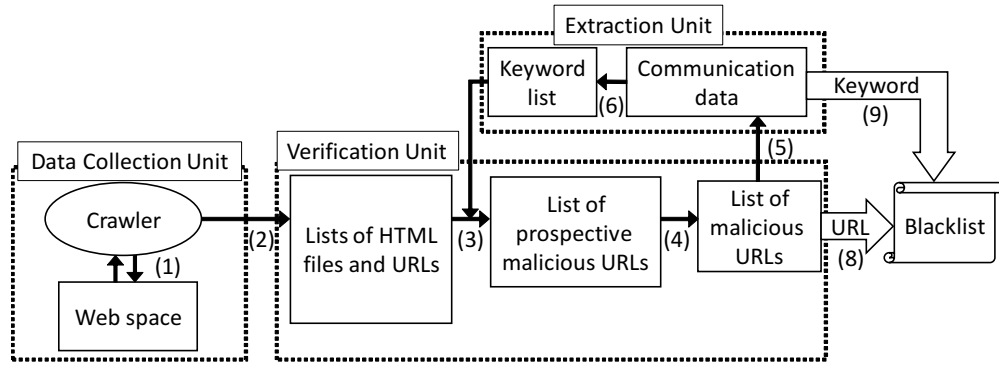
Fig. 2.  Basic design of proposed method

is necessary to manage large-scale data in the vast web space.

(Problem 2)  Timely discovery of malicious websites
Some attackers use a method to change the IP address or the domain name of a malicious website within a short period of time to bypass blacklists. Therefore, it is necessary to discover malicious websites in a timely manner and update the blacklist frequently.

### C. Basic Design

To address problems 1 and 2, the proposed method collects a large amount of web content from the web space using a crawler to discover malicious websites. However, it is difficult to discover malicious websites by verifying all the collected web contents. Therefore, we extract malicious web contents from the collected web contents, and verify and analyze the extracted web contents. This makes it possible to discover malicious websites efficiently.

The proposed method is classified into three parts: a data collection unit, a verification unit, and an extraction unit. The data collection unit collects URLs and corresponding HTML files from the web space. The verification unit searches the HTML files using list of keywords extracted from known malicious websites (keyword list) and finds URLs that are highly likely to be malicious (prospective malicious URLs). We also detect malicious websites in the verification unit. The extraction unit analyzes the communication data of the websites detected to be malicious and extracts the keywords used for keyword searches. The keywords to be extracted are explained in section IV-C. The flow diagram of generating a blacklist in the proposed method is displayed in Figure 2 and explained as follows.

(1) Collect HTML files as web content from the web space using crawler.
(2) Store crawled URLs and collected HTML files.
(3) Search the collected HTML files using keyword list, and if the HTML file contains the keyword, add the URL corresponding to the HTML file to the prospective malicious URL list.

(4) If the URL added to the prospective malicious URL list is a malicious website, add the URL to the malicious URL list.
(5) Acquire communication data to access the malicious website by using an application that collects the URL when accessing the website using Google Chrome.
(6) Expand the keyword list by extracting the keywords from the acquired communication data.
(7) Repeat (3) to (6) using the expanded keyword list until no new malicious URL is found.
(8) Add URL of malicious URL list to the blacklist.
(9) Add extracted keyword to the blacklist.

In the proposed method, the keyword list is expanded each time a keyword is extracted. In addition, the expanded keyword list is used to search the HTML file again. As a result, there is a possibility to discover new malicious websites that were missed by the keyword list before the expansion.

## IV. Implementation

### A. Selecting URLs to Crawl

There are many websites on the world wide web, and the URLs of malicious websites are changed frequently. It is difficult to collect malicious websites by crawling the web space randomly. Therefore, we collect URLs by using Twitter's Streaming API as the URLs to crawl because the attackers have a method of posting the URL of the malicious website on social networking services (SNS) such as Twitter, Facebook, and YouTube. Collecting URLs from Twitter can help discover the malicious websites more efficiently than randomly crawling the web space. In addition, Twitter's Streaming API can acquire the tweets in near real-time, so that the URL of an unknown malicious website posted by an attacker can be discovered in a timely manner.

### B. Detecting Malicious Websites

The malicious website is detected by manual access and checking whether it is redirected to the unwanted websites. Unwanted transition to an unwanted website occurs when the user taps on the screen, as mentioned in Section II. It is necessary to distinguish between the transition by clicking on a legitimate link and the unwanted redirect. When a user

| Target | Keywords to extract |
|---|---|
| HTML files of the landing site | The filename (e.g., example.js) that causes the redirection |
| | FQDN that provides the file that causes the redirection |
| URL of the intermediate site | FQDN |
| URL of the unwanted website | FQDN |

TABLE I
KEYWORDS TO EXTRACT

taps anywhere on the screen and is redirected to the unwanted websites, such as those displaying "Fake Virus Alert," that site is considered a malicious website. Although manual access is time-consuming, this method only checks the prospective malicious URLs; thus, there are not too many websites to check.

### C. Extraction of Keywords

The method of extracting keywords is described in Table I. The HTML file of the website detected to be malicious is the HTML file of the landing site. The landing site is considered to have a file that causes the redirection. Regarding the known malicious websites, there is a strong possibility that common malicious contents are placed in similar domains [14]. Therefore, the filename (e.g., example.js) that causes the redirection and the fully qualified domain name (FQDN) that provides this file are extracted as keywords from the HTML file.

From the communication data, the FQDN is extracted as a keyword from the URLs of the intermediate site and the unwanted website because the landing site may redirect the user to the intermediate site and the unwanted website. The URL of the intermediate site is sometimes created from the specified URL and a randomly created character string [12]. Furthermore, the URL of the unwanted website sometimes includes the user's device information [12]. Therefore, we extract the FQDN from the URLs of the intermediate site and unwanted website.

## V. EVALUATION

### A. Purpose and Evaluation Environment

There are two criteria used for the evaluation of the method:

(Evaluation 1)  Number of malicious websites discovered by the search.
(Evaluation 2)  Detection rate of malicious websites using a blacklist.

For Evaluation 1, we performed steps (1) to (9) in Section III-C between July 23 and December 16, 2019 to determine the number of malicious websites that can be discovered using the method. In the initial keyword list, we set the keywords extracted from an independently discovered malicious website.

For Evaluation 2, we verified whether access to malicious websites could be detected using the blacklist generated by the proposed method in Evaluation 1. For the evaluation, five malicious websites (Site A to Site E) collected using the proposed method from December 20 to December 30, 2019

were used. We accessed the five malicious sites and tapped anywhere on the screen to test the unwanted transition. The malicious website does not always have the same redirect destination for each access; therefore, we accessed it 10 times. In addition, we confirmed that site E transitions to two different intermediate sites and unwanted websites depending on the timing of the tapping on the screen. For this reason, Site E was accessed 10 times each. In addition, for verification, we used the application that notifies the user of access to the malicious website when the character string of the URL bar of Google Chrome is acquired and the acquired character string contains the one registered in the blacklist. The malicious website was accessed using Google Chrome on a real device with OS Android 6.0.

### B. Results and Discussion

*1) Number of Malicious Websites Discovered by Search:* As a result of the search, 200 landing sites were detected from 122,350 crawled websites. These 200 sites included adult sites, video sharing sites, and news sites. In fact, 91 malicious sites had the same URL as a seemingly legitimate world news website. It has been previously reported that legitimate websites are used for attacks that use the links to the actual news sites posted on multiple forums and execute malicious codes [15]. We plan to examine how the attack that redirects users to unwanted websites occurs on legitimate sites in future research.

Google's Safe Browsing technology examines billions of URLs per day looking for unsafe websites and discovers thousands of new unsafe sites every day [16]. Although it is not possible to make a simple comparison with Google's Safe Browsing technology, the proposed method can also efficiently discover malicious websites. Of the 200 landing sites we discovered, 182 landing sites were occurred redirect as of January 7, 2020 could not be detected by Google Safe Browsing, which uses a URL blacklist method. Thus, we need to analyze the reason that about 90% of malicious websites cannot be detected using the existing blacklist. In addition, since the detection rate of the existing blacklist is low, a blacklist that can be updated in a timely manner for mobile devices is urgently needed.

Among 200 landing sites, 54 sites were unique FQDNs. We extracted 111 keywords from the communication data when 54 landing sites with unique FQDNs were accessed. Among the extracted keywords, there were 3 file names and 108 FQDNs. Of the FQDNs extracted as keywords, the number of FQDNs providing the files that causes the redirection was 3, which is equal to number of file names extracted because the file that causes the redirection could not be specified on the landing site due to the obfuscation of the JavaScript code.

*2) Detection Rate of Malicious Websites Using Blacklists:* Table II shows the detection results of malicious websites using a blacklist generated by searching malicious websites from July 23 to December 16, 2019.

To comply with the requirements of ethical research, the matching keywords are represented by labels. For example,

## TABLE II
### DETECTION RESULTS OF MALICIOUS WEBSITES USING A BLACKLIST

| | Use of different intermediate sites or unwanted websites | Detection rate (detection/accesses) | Keyword that matched when detected |
|---|---|---|---|
| Site A | None | 100% (10/10) | FQDN1 of intermediate site<br>FQDN1 of unwanted website |
| Site B | None | 100% (10/10) | FQDN1 of intermediate site<br>FQDN1 of unwanted website |
| Site C | None | 100% (10/10) | FQDN2 of intermediate site<br>FQDN2 of unwanted website |
| Site D | None | 100% (10/10) | FQDN3 of intermediate site |
| Site E | Use | 100% (10/10) | (Transition a)<br>FQDN4 to FQDN6 of intermediate site<br>FQDN3 of unwanted website |
| | | 60% (6/10) | (Transition b)<br>FQDN7 of unwanted website |

FQDN1 of the intermediate site in Table II is one of the FQDNs that are the keywords extracted from the intermediate site.

On sites A to D, the access to the malicious websites was detected by the keyword matched at the time of detection shown in Table II for all 10 accesses.

In Transition a of Site E, access to the malicious websites was detected in all 10 accesses by FQDN4 to FQDN6 extracted from the intermediate site and FQDN3 extracted from the unwanted website. In Transition b, access to the malicious website was detected in 6 out of 10 accesses by FQDN7 extracted from the intermediate site. In Transition b, there were cases that false negatives occurred. In these cases, a FQDN that was not registered in the blacklist was used for the intermediate site. This FQDN was different from the FQDN7 of the intermediate site by only one character.

The keywords matched at the time of detection are the FQDNs of the intermediate site or the FQDNs of the unwanted website. Since the file name that causes the redirection and the FQDNs that provide the file are small in the number of extractions, it is presumed that the keywords did not match at the time of detection.

## VI. RELATED WORK

Sun et al. [17] proposed AutoBLG, which automatically detects a new malicious URL from the web space and generates a blacklist. AutoBLG efficiently discovers malicious URLs by collecting unknown URLs using IP addresses of malicious URLs and reducing the number of URLs to analyze by filtering. AutoBLG focuses on the drive-by download attack as a web-based attack. In contrast, in the proposed method, we focus on the attacks redirecting a user to unwanted websites.

There is another study that proposed kAYO, which detects malicious websites for mobile devices in real-time [18] using supervised machine learning and thus require labeled teacher data used for learning in advance. However, creating such data is very costly. The proposed method uses a relatively small number of keywords extracted from malicious websites to discover new malicious websites.

## VII. CONCLUSION

To detect the cyberattacks that redirect users to unwanted websites, we proposed a method that generates a blacklist for mobile devices. The proposed method efficiently discovers malicious websites by using a large number of HTML files collected from the web space using a crawler, and searching for HTML files that are likely to be malicious using the keywords extracted from the known malicious websites. The blacklist that is generated using the proposed method includes FQDNs of intermediate sites or unwanted websites. This makes it possible to detect access to not only the landing site but also the intermediate sites and unwanted websites.

Using the proposed method, we found 182 landing sites that could not be detected by Google Safe Browsing. We also gained access to malicious websites on mobile devices using the blacklist generated by the new method. Of the five malicious websites used in the evaluation, the blacklist detected all 10 accesses at four websites. This evaluation result shows that the blacklist using keywords has a sufficiently high detection rate for a specific malicious website. However, one of the reasons for the high detection rate is that the malicious websites collected using the proposed method were used for evaluation.

In our future work, we plan to examine how the attack that redirects users to unwanted websites occurs on legitimate sites and evaluate the detection rate of our blacklist for common malicious URL lists.

### REFERENCES

[1] DataReportal, Digital 2020: Global Digital Overview, https://datareportal.com/reports/digital-2020-global-digital-overview Accessed 3 May 2020
[2] StatCounter Global Stats, Desktop vs Mobile vs Tablet Market Share Worldwide, https://gs.statcounter.com/platform-market-share/desktop-mobile-tablet Accessed 3 May 2020

[3] Net Market Share, https://netmarketshare.com Accessed 3 May 2020

[4] StatCounter Global Stats, Mobile Operating System Market Share Worldwide, https://gs.statcounter.com/os-market-share/mobile/worldwide Accessed 3 May 2020

[5] Kaspersky, Mobile malware evolution 2018, https://securelist.com/mobile-malware-evolution-2018/89689/ Accessed 8 Aug 2019

[6] Wandera, Android Malware: 4 Ways Hackers are Infecting Phones with Viruses, https://www.wandera.com/malware-on-android/ Accessed 3 May 2020

[7] Trend Micro, Social media malware on the rise, https://blog.trendmicro.com/social-media-malware-on-the-rise/ Accessed 4 May 2020

[8] CalyptixSecurity, Social Media Threats: Facebook Malware, Twitter Phishing, and More, https://www.calyptix.com/top-threats/social-media-threats-facebook-malware-twitter-phishing/ Accessed 4 May 2020

[9] M. Levinson, Mobile Malware: Beware Drive-by Downloads on Your Smartphone, https://www.cio.com/article/2397969/mobile-malware--beware-drive-by-downloads-on-your-smartphone.html Accessed 4 May 2020

[10] J. Doevan, Android virus. Versions provided. The list of infected apps for 2020, https://www.2-spyware.com/remove-android-virus.html Accessed 4 May 2020

[11] L. Wu, X. Du, and J. Wu, Effective Defense Schemes for Phishing Attacks on Mobile Computing Platforms, *IEEE Transactions on Vehicular Technology*, vol. 65, no. 8, pp. 6678–6691, 2016.

[12] Y. Imamura, R. Orito, K. Chaikaew, C. Manardo, P. Leelaprute, M. Sato, and T. Yamauchi, Threat Analysis of Fake Virus Alerts Using WebView Monitor, In *proceedings of the 2019 7th International Symposium on Computing and Networking (CANDAR'19)*, pp. 28–36, 2019.

[13] Internet Live Stats, Total number of Websites, https://www.internetlivestats.com/total-number-of-websites/ Accessed 8 Aug 2019

[14] L. Invernizzi, P. M. Comparetti, S. Benvenuti, C. Kruegel, M. Cova, and G. Vigna, Evilseed: A Guided Approach to Finding Malicious Web Pages, In *proceedings of the 2012 IEEE Symposium on Security and Privacy*, pp. 428–442, 2012.

[15] Trend Micro, Operation Poisoned News: Hong Kong Users Targeted With Mobile Malware via Local News Links, https://blog.trendmicro.com/trendlabs-security-intelligence/operation-poisoned-news-hong-kong-users-targeted-with-mobile-malware-via-local-news-links/ Accessed 5 May 2020

[16] Google Safe Browsing, Google Transparency Report, https://transparencyreport.google.com/safe-browsing/search?hl=en Accessed 4 May 2020

[17] B. Sun, M. Akiyama, T. Yagi, M. Hatada and T. Mori, Automating URL Blacklist Generation with Similarity Search Approach, *IEICE Transactions on Information and Systems*, vol. E99-D, no. 4, pp. 873–882, 2016.

[18] C. Amrutkar, Y. S. Kim and P. Traynor, Detecting Mobile Malicious Webpages in Real Time, *IEEE Transactions on Mobile Computing*, vol. 16, no. 8, pp. 2184–2197, 2017.