

Protect Sensitive Sites from Phishing Attacks Using Features Extractable from Inaccessible Phishing URLs

Weibo Chu*, Bin B. Zhu†, Feng Xue†, Xiaohong Guan*‡, Zhongmin Cai*

*MOE KLINNS Lab, Xi'an Jiaotong University, Xi'an, China

†Microsoft Research Asia, Beijing, China

‡Center for Intelligent and Networked System and NLIST Lab, Tsinghua University, Beijing, China

Email: (wbchu, xhguan, zmcai)@sei.xjtu.edu.cn, (binzhu, feng.xue)@microsoft.com

Abstract—Phishing is the third cyber-security threat globally and the first cyber-security threat in China. There were 61.69 million phishing victims in China alone from June 2011 to June 2012, with the total annual monetary loss more than 4.64 billion US dollars. These phishing attacks were highly concentrated in targeting at a few major Websites. Many phishing Webpages had a very short life span. In this paper, we assume the Websites to protect against phishing attacks are known, and study the effectiveness of machine learning based phishing detection using only lexical and domain features, which are available even when the phishing Webpages are inaccessible. We propose several novel highly effective features, and use the real phishing attack data against Taobao and Tencent, two main phishing targets in China, in studying the effectiveness of each feature, and each group of features. We then select an optimal set of features in our phishing detector, which has achieved a detection rate better than 98%, with a false positive rate of 0.64% or less. The detector is still effective when the distribution of phishing URLs changes.

I. INTRODUCTION

Online services such as online shopping or online banking have brought us a great convenience yet at the same time new threats. One of these new threats is phishing whereby spoofed emails or instant messages purporting to be from trustworthy sources are used to lure recipients to click the contained URLs that lead to counterfeit websites to trick them into divulging sensitive information such as usernames and passwords, credit card information, social security numbers, etc. Phishing remains to be a serious cyber-security threat. It is the third cyber-security threat globally and the first cyber-security threat in China. For one year from June 2011 to June 2012, there were 61.69 million phishing victims in China alone, with the total annual monetary loss more than 4.64 billion US dollars [1].

To thwart phishing attacks, a great effort has been directed towards detecting phishing. A variation of approaches has been proposed, including blacklisting [5] and whitelisting [6], and anomaly-based detection methods. These methods will be briefly reviewed in Section II. Among anomaly-based

detection methods, a widely used approach is to apply machine learning to a training set consisting of both phishing and benign URLs to build a classification model based on carefully selected discriminative features. Typical discriminative features include lexical features derived from the URL strings, linkage features derived from the relationship between the URL and other Websites, hosting features related to the hosting server of the URL, Webpage features extracted from the Webpage code of the URL, network features derived from accessing the URL. Some features such as Webpage features and network features can be obtained only when the Webpage of the URL is alive.

The existing machine-learning based phishing detectors were typically designed to detect generic phishing attacks that may target any sites and any people. Recent studies [2] indicate recent phishing attacks tended to be “spear-phishing” that targets at specific groups of people. According to Risings report [3], the top four Websites that phishing attacks targeted in the first half year of 2011 in China were Taobao, Tencent, Industrial and Commercial Bank of China (ICBC), and Bank of China (BOC). Taobao is the largest Internet retail website in China, with more than 170 million users. Tencent provides a popular Web portal and the largest instant messenger QQ in China, with more than 600 million active QQ users. ICBC is the largest bank in the world. This phenomenon of highly concentrated phishing targets has also been reported by others. For example, the Anti-Phishing Alliance of China reported that the top four Websites targeted by phishing attacks in the month of April 2012 were Taobao, ICBC, Chinese Central TV, and Tencent. The attacks against these four sites accounts for 93.67% of all the phishing attacks reported to the alliance [4].

In this paper, we investigate the effectiveness of machine-learning based phishing detection when the targeted phishing Websites are known. The actual phishing data targeted at Taobao and Tencent have been used in our studies. This is a position paper for an ongoing project to develop an effective phishing detector to protect the users of the aforementioned major Websites targeted by phishing attacks. The detector can be deployed to these users as a Web browser plugin.

This work was supported by NSFC (60921003, 61175039, 60905018), and the Fundamental Research Funds for Central Universities (xjj20100051, 2012jdhz08).

There are several challenges in our studies. A major challenge is that many phishing Webpages are short-lived, typically less than 20 hours [19], and URLs may change frequently (fast-flux). For example, we received regular (weekly initially and then daily) reports of phishing URLs from Taobao. Upon receiving the report, we immediately access the phishing Webpages but more than 80% of the phishing URLs were inaccessible. As a consequence, the discriminative features obtained from live Websites such as Webpage features and network features used in existing phishing detectors can no longer be used. In this paper, only lexical and domain features are used in our phishing detection. These features are readily available without accessing the Webpage, and thus can be used even if the phishing URLs are no longer accessible. Our studies indicate that our detector is highly effective even with the reduced types of discriminative features, with detection rates better than 98% and false positive rates at 0.64% or less.

This paper has the following major contributions:

1. We have studied phishing detection performance using actual phishing attacks against popular phishing targets in China. The discriminative features used in our detector can be obtained even when a Webpage is inaccessible. As a result, the short-lived phishing URLs have also been included in our studies, and thus the detection performance from our studies is closer to actual performance in real deployment than most previous studies which excluded short-lived URLs.

2. We have studied each discriminative feature's power in detecting phishing attacks using the aforementioned real phishing attacks. This study helps understand the contributions of each feature in the overall detection performance.

3. We have proposed several novel highly effective discriminative features including a similarity measure to the brand names of the sites to be protected against phishing, domain age, and domain confidence level.

The rest of the paper is organized as follows. Section II reviews related work. Section III provides a detailed description of the proposed detector and its discriminative features. The performance evaluation of the detector against more than one year's real-life phishing attacks is reported in section IV. We conclude the paper in Section V.

II. RELATED WORK

Blacklisting [5] uses a blacklist of phishing URLs or domains to block phishing URLs. It incurs no false positive yet is effective only to detect known phishing URLs. A blacklist is generally constructed through time-consuming human feedbacks, and thus ineffective in blocking short-lived phishing Webpages. Whitelisting [6], on the other hand, seeks to identify known good sites by maintaining a whitelist of benign URLs or domains. Any URL not in the whitelist will be blocked. Whitelisting incurs no false negative but may unavoidably result in a high false positive.

The weakness of blacklisting and whitelisting has been addressed by anomaly-based phishing detectors which rely on a classification model based on discriminative rules or features. The classification model can be built with knowledge a priori.

Zhang et al. [10] proposed a system to detect phishing URLs with a weighted sum of 8 features related to Web content, lexical and WHOIS data. Garera et al. [11] used logistic regression over manually selected features to classify phishing URLs. The features include heuristics from a URL such as Google's page rank features. Xiang and Hong [8] proposed a hybrid detection method by discovering inconsistency between a phishing identity and the corresponding legitimate identity. PhishNet [9] provides a prediction method for phishing attacks using known heuristics to identify phishing pages.

The classification model can also be built through machine learning. Fette et al. [12] proposed a system to classify phishing emails. They used a large publicly available corpus of legitimate and phishing emails. Their classifiers examine 10 different features such as the number of URLs in an e-mail, the number of domains and dots in these URLs. Whittaker et al. [14] proposed a phishing webpage classifier to update Google's phishing blacklist automatically. Their detector shares many discriminative features used in [13]. Ludl et al. discussed a system for classifying phishing pages based on Webpage features [15]. Ma et al. published a pair of papers describing another system for identifying malicious URLs by examining lexical features of the URLs and features of the site's hosting information [17][18]. Choi et al. [16] proposed a malicious URL detector that uses a large set of features including lexical, linkage, Webpage, networking, and DNS features. Our detector shares many features with their detector.

Visual similarity has also been exploited. Chen et al. used Contrast Context Histogram (CCH) [21] to describe the images of Webpages and adopts Euclidean distance to find matching between two sites. Fu et al. used Earth Movers Distance (EMD) [22] to measure page similarity. They first convert the involved pages into low resolution images and then use color and coordinate features to represent the image signatures. EMD is employed to calculate the signature distance of the images of the pages. Dunlop [24] experimented with optical character recognition to convert screenshots into text to help detect phishing. Liu et al. [25] used layout and style similarity to evaluate visual similarity, and iTrustPage [26] uses Google search and user opinion to identify visually similar pages.

III. OUR DETECTOR AND ITS DISCRIMINATIVE FEATURES

A. System Overview

Our system consists of two stages, the *learning stage* and the *detection stage*. The system's flowchart is shown in Fig. 1. The *Redirection Parse* module in both stages converts the received URLs into their true URLs. We observed that a significant portion of phishing pages against Taobao were shortened URLs, such as url.cn and goo.cn, or redirected URLs in order to trick recipients. Both the original URLs and their true URLs, if exist, are then passed to the *Feature Extraction* module to extract features for model training and classification.

B. Learning Algorithm

SVM has been widely used as a machine learning method to train a binary classification model with training data. In our

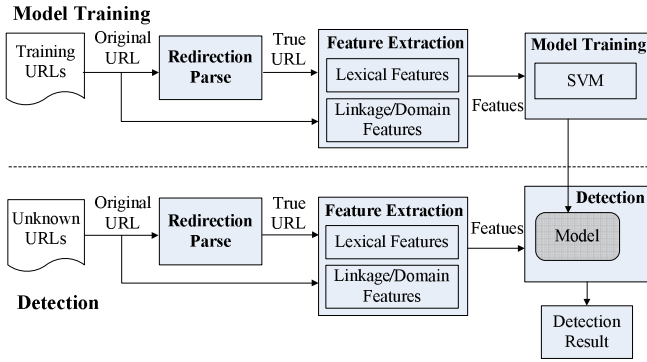


Fig. 1. The flowchart of our detection system

TABLE I
DISCRIMINATIVE FEATURES USED IN OUR DETECTOR

No.	Feature	Category	Type
1	Domain token count	Lexical	Integer
2	Average domain token length	Lexical	Real
3	Longest domain token length	Lexical	Integer
4	Path token count	Lexical	Integer
5	Average path token length	Lexical	Real
6	Longest path token length	Lexical	Integer
7	Domain brand-name distance	Lexical	Integer
8	Path brand-name distance	Lexical	Integer
9	Domain Google links	Domain	Integer
10	Domain Baidu links	Domain	Integer
11	Domain Bing links	Domain	Integer
12	Domain Yahoo! links	Domain	Integer
13	SLD Google links	Domain	Integer
14	SLD Baidu links	Domain	Integer
15	Domain page rank	Domain	Integer
16	Domain alexa rank	Domain	Integer
17	Domain age	Domain	Integer
18	Domain confidence level	Domain	Real

work, SVM was used to train the classification model of our detector. A Gaussian Radial Basis Function (RBF) kernel was used with SVM in our experiments to be reported later.

C. Discriminative Features

The 18 discriminative features listed in Table 1 are used in our phishing detection system. These features can be roughly classified into two groups: lexical and domain features.

1) Lexical features

Phishing URLs often have distinguishable patterns in their URL text. We adopt 8 lexical features (No. 1-8 in Table 1) in our detection method. Among these features, feature No. 1-6 are from the previous work. Feature No. 7 and 8 are novel features never used before.

The first 8 features in Table 1 are lexical features, with the first 6 features from [16] and the last 2 lexical features being novel features never used before. A token is a substring in the URL delimited by ‘:’, ‘/’, ‘?’, ‘=’, ‘-’, ‘_’. Previous detection systems use a binary feature to check whether

a brand name is contained in the URL tokens. A careful examination of phishing attacks against Taobao indicated that a significant portion of phishing URLs contained tokens similar to but different from brand-names. For example, “tac.bao” in “www.tac.bao.com.cn” is different from but similar to brand name “taobao”, one of Taobao’s second-level domain (SLD) names. The binary feature of *brand name presence* [20] would not capture this vital characteristic. Therefore, we propose two new features: *domain brand-name distance* and *path brand-name distance*, which are defined as follows:

Definition 1: Let $B = \{b_1, b_2, \dots, b_n\}$ be the set of brand names of one site or more sites to be protected against phishing. Let s be a string of domain or path that we need to calculate the brand-name distance from B , and $S = \{s_1, s_2, \dots, s_m\}$ be the set of all the substrings derived from s . The brand-name distance between s and b_i is defined as the minimum edit distance (i.e., Levenshtein distance) between all substrings of s and b_i :

$$brand_dist(s, b_i) = \min\{edit_dist(s_j, b_i) | s_j \in S\} \quad (1)$$

The brand-name distance between s and B is defined as the minimum brand-name distance between s and all the brand names in B :

$$brand_dist(s, B) = \min\{brand_dist(s, b_i) | b_i \in B\} \quad (2)$$

For example, if we want to protect Taobao from phishing attacks, $B = \{\text{“taobao”}, \text{“alibaba”}, \text{“alipay”}\}$ contains three SLD names used by Taobao. For URL “www.tao.bac.com.cn”, the *domain brand-name distance* between this URL and B (i.e., Taobao) is 2, which is the edit distance between substring “tao.bac” and brand name “taobao”. The *path brand-name distance* is calculated in the same way.

A phishing URL targeting at Taobao tends to contain a substring similar to one of Taobao’s brand names in either the domain or the path, and thus has a domain or path brand-name distance smaller than that of benign URLs. In our detector, legitimate sites containing substrings similar to the brand-names of B are collected and placed in the whitelist, and thus their brand-name distances are not calculated. The brand-name distance features are a superset of the *brand name presence* feature in [16] and other papers.

2) Domain features

The remaining features in Table 1 can be roughly classified as domain features. They are used to capture information of a site such as its link popularities, domain reliabilities, domain age, etc. Phishing sites tend to have a small value of link popularities, whereas most benign sites, especially those popular, tend to have a large value of link popularity. In our detector, four search engines, *Google*, *Bing*, *Baidu*, and *Yahoo!*, are used to calculate the link popularity of a site and its SLD. Link popularity was borrowed from [16]. In addition, we borrowed page rank from [14] in our detector. *Google’s page rank* and *Alexa rank* of a site were used in our method since they are much harder to forge or manipulate than the above link popularities (e.g., through “link farming” [23]). These discriminative features comprise of features No. 9-16.

Our detector has also adopted two new features, *domain age* and *domain confidence level*. They are designed to capture the characteristics that phishing URLs tend to use domains with a short life than the domains of benign URLs. In order to calculate the *domain confidence level* of a URL, we maintain two lists: a list of benign URLs and a list of phishing URLs. The *domain confidence level* is defined as follows:

Definition 2: Let d be the domain to be checked and $SLD(d)$ be the second level domain of d . Let x be the number of benign URLs hosted by $SLD(d)$ in our benign URL list, and y be the number of phishing URLs hosted by $SLD(d)$ in phishing URL list. The *domain confidence level* of d is defined as follows:

$$domain_conf_level(d) = \left[\frac{x + A}{x + y + 2A} - 0.5 \right] \times \frac{3}{5} + 0.5 \quad (3)$$

Here A is a parameter to avoid oversensitivity to small x and y . We set $A = 1000$ in our experimental studies. The range of the *domain confidence level* calculated with Eq. (3) is (0.2, 0.8), with 1.0 assigned to the domain confidence level for a URL in the whitelist and 0 for a URL in the blacklist.

IV. PERFORMANCE EVALUATIONS

A. Dataset

Phishing URLs: we used 17423 distinct Taobao-phishing URLs received from Taobao in evaluating the performance of our detector. These phishing URLs were reported to and confirmed by Taobao from Jan. 2011 to April 2012.

Benign URLs: We collected 28722 benign URLs from Yahoo!'s directory (<http://random.yahoo.com/bin/ryl>) and also by crawling some well-known Chinese navigation sites.

B. Detection Performance

The Taobao-phishing dataset and the benign dataset described above were used with 5-fold cross-validation in evaluating the performance of our detector using the following metrics: *accuracy* (ACC) which is the ratio of true results (both true positives and true negatives) over all the samples in the datasets; *false positive rate* (FP) which is the proportion of benign URLs that are mistakenly identified as phishing URLs; and *false negative rate* (FN) which is the proportion of phishing URLs that are missed by the detector. Libsvm was used as the SVM implementation in our experiments.

1) Evaluation of Each Feature Group

We first studied the effectiveness of each feature group by performing detection using only the features in each of the two groups: lexical group and domain group. Table II shows the performance results for each feature group. We can see clearly from the table that both feature groups had a good detection performance, with accuracy better than 95%. The lexical feature group had a little worse performance than the domain feature group, 95.88% vs. 98.14% for the accuracy. For the setting of the experiments, the lexical feature group leaned towards a small FP at 0.82%, resulting in a much higher FN than that of the domain feature group, 9.38% vs. 1.37%. The domain feature group, on the other hand, leaned towards

TABLE II
PERFORMANCE OF EACH FEATURE GROUP

Category	ACC	FP	FN
Lexical Features	95.88%	0.82%	9.38%
Domain Features	98.14%	2.56%	1.37%

a smaller FN at 1.37%, resulting in a much larger FP than that of the lexical feature group, 2.65% vs. 0.82%.

2) Evaluation of Individual Features

We then performed detection using each individual feature alone to study its effectiveness and contribution to the detection of phishing URLs. Table III shows the experimental results of the performance of each individual feature. From the table, we can draw the following conclusions:

a) Among all the 18 discriminative features, *domain brand-name distance* is the most effective feature, indicating that Taobao-phishing URLs tended to contain a string similar to Taobao's brand names to trick users. *Domain page rank* and *domain age* are two next most effective features.

b) *Domain token count* is the second most effective features in the lexical feature group, next to *domain brand-name distance* which is the most effective feature among all the features in both feature groups.

c) Among the four search engines, the link popularity provided by *Bing* provides the least effective feature in phishing detection. It is also surprising that *Google* provides a similar performance, much worse than that provided by *Baidu*. We conducted an investigation and found that *Google* reported only a partial list of link popularity information, confirmed by Google's official website.

d) *Domain confidence level* is also a very effective feature in detecting Taobao-phishing URLs. In fact, we found that Taobao phishers tend to use the same domains with the same SLD and TLD (Top-Level Domain) to launch many phishing pages. As a result, new Taobao-phishing URLs tend to share domains with old phishing URLs, leading to effective detection by *domain confidence level*.

3) Performance of Our Phishing Detector

With the effectiveness data of individual discriminative features, we applied the plus-m-minus-r algorithm to select the most effective features to be used in our detector. Due to space limitations, we omit the detailed detection performance data in this feature selection process. Our results show that, with the plus-2-minus-1 algorithm, the 7th feature (i.e., *domain brand-name distance*) contributes the most to phishing detection, whereas the 6th feature (i.e., *longest path token length*) contributes the least. The combination that achieves the best performance comprises of the following features (listed by their indexes in Table 1): {7, 18, 8, 15, 17, 12, 14, 11, 10, 9, 2, 16, 3, 4, 5}. Using this optimized set of features, our detector produced the following results: ACC = 99.35%, FP = 0.45%, and FN=1.01%.

4) Performance Impact of A Changing Distribution

To study our detector's performance under realistic application scenario where new URLs may have a different

TABLE III
PERFORMANCE OF INDIVIDUAL FEATURES

No.	Feature	ACC
1	Domain token count	86.34%
2	Average domain token length	63.31%
3	Longest domain token length	61.89%
4	Path token count	76.54%
5	Average path token length	75.52%
6	Longest path token length	63.72%
7	Domain brand-name distance	88.44%
8	Path brand-name distance	75.93%
9	Domain Google links	69.12%
10	Domain Baidu links	84.29%
11	Domain Bing links	66.39%
12	Domain Yahoo! links	72.90%
13	SLD Google links	69.10%
14	SLD Baidu links	78.44%
15	Domain page rank	87.28%
16	Domain alexa rank	71.56%
17	Domain age	87.19%
18	Domain confidence level	84.67%

distribution from that of the training data, we applied the detection model obtained in the previous subsection to the 599 “new” Taobao-phishing URLs we received in May 2012. Our detector achieved the following performance: $ACC=99.22\%$ and $FN=8.51\%$, with $FP=0.45\%$ unchanged.

5) Performance for Phishing URLs Targeted at Tencent

We also applied our system to detect phishing attacks targeting at Tencent. The phishing dataset comprised of 34657 phishing URLs we received from Tencent, and the same dataset of benign URLs were used. Our detector has achieved the following performance: $ACC=98.72\%$, $FP=0.64\%$, and $FN=1.82\%$, which is close but a little worse than the performance in detecting phishing URLs targeted at Taobao. The result indicates the robustness of our detector in detecting phishing URLs targeted at different Websites.

V. CONCLUSION

In this paper we investigated the effectiveness of machine learning based phishing detection with known protected Websites. Only lexical and domain features were used since many phishing URLs had a short life span, and these features were typically still available even when phishing Webpages were inaccessible. We proposed several novel highly effective features. We studied effectiveness of each feature and selected an optimal set of features in our detector, which achieved a detection rate better than of 98%, with a false positive rate of 0.64% or below. The detection rate with changed distribution of phishing URLs was still above 91%.

REFERENCES

[1] Trusted E-Commerce Promotion Center of China E-Commerce Association, et al., 2012 Report of Trust Verification Development for China's Websites, July, 2012, accessible from http://ectrust.knet.cn/column_2/201207/W020120704645636974021.pdf.

[2] Anti-Phishing Working Group, Phishing activity trends report, 2008, http://www.antiphishing.org/reports/apwg_report_Q4_2009.pdf, 2009.

[3] Rising, The first semi-annual Report on Internet Security, 2011, accessible from <http://www.rising.com.cn/2011/report/report2011.doc>.

[4] Anti-Phishing Alliance of China, April 2012 Report of Phishing Websites, <http://www.apac.org.cn/gzdt/201205/P020120518602784673833.pdf>.

[5] S. Sheng, B. Wardman, G. Warner, et al., “An empirical analysis of phishing blacklists,” In proceedings of 6th Conference on Email and AntiSpam (CEAS 2009), Mountain View, CA, USA, July 2009.

[6] N. Chou, R. Ledesma, Y. Teraguchi, et al., “Client-side defense against web-based identify theft,” In Proceedings of the 11th Annual Network and Distributed System Security Symposium (NDSS04), San Diego, 2004.

[7] A. Ntoulas, M. Najork, M. Manasse, et al., “Detecting spam web pages through content analysis,” In proceedings of the 15th International Conference on World Wide Web, Edinburgh, 2006.

[8] G. Xiang and J. I. Hong, “A hybrid phish detection approach by identity discovery and keywords retrieval,” In proceedings of the 18th International Conference on World Wide Web, Madrid, Spain, 2009.

[9] P. Prakash, M. Kumar, R. R. Kompella, et al., “PhishNet: Predictive blacklisting to detect phishing attacks,” In Proceedings of the 29th conference on Information communications (INFOCOM), 2010.

[10] Y. Zhang, J. Hong, and L. Cranor, “Cantina: A content based approach to detecting phishing web sites,” In proceedings of the 16th International conference on World Wide Web, New York, NY, USA, 2007.

[11] S. Garera, N. Provos, M. Chew, et al., “A framework for detection and measurement of phishing attacks,” In proceedings of the 2007 ACM workshop on Recurring Malcode, VA, USA, 2007.

[12] I. Fette, N. Sadeh, and A. Tomasic, “Learning to detect phishing emails,” In proceedings of the 16th International Conference on World Wide Web, Banff, Alberta, Canada, 2007.

[13] A. Ramachandran and N. Feamster, “Understanding the network-level behaviors of spammers,” ACM SIGCOMM Computer Communication Review-Proceedings of the 2006 conference on Applications, technologies, architectures, and protocols for computer communications, New York, NY, USA, Vol. 35, Issue 4, 2006.

[14] C. Whittaker, B. Ryner, and M. Nazif, “Large-scale automatic classification of phishing pages,” In Proceedings of the 17th Annual Network and Distributed System Security Symposium (NDSS10), 2010.

[15] C. Ludl, S. McAllister, E. Kirda, et al., “On the effectiveness of techniques to detect phishing sites,” In Proceedings of the International conference on Detection of Intrusions and Malware & Vulnerability Assessment (DIMVA), 2007.

[16] H. Choi, Bin B. Zhu, and H. Lee, “Detecting malicious web links and identifying their attack types,” USENIX International Conference on Web Application Development (WebApps), Portland, USA, 2011.

[17] J. Ma, L. K. Saul, S. Savage, et al., “Beyond blacklists: learning to detect malicious web sites from suspicious URLs,” In proceedings of the International Conference on Knowledge Discovery and Data Mining, Paris, France, 2009.

[18] J. Ma, L. K. Saul, S. Savage, et al., “Identifying suspicious urls: An application of large-scale online learning,” In proceedings of the International Conference on Machine Learning, Montreal, Canada, 2009.

[19] T. Moore and R. Clayton, “Examining the impact of website take-down on phishing,” In proceedings of Anti-Phishing Working Group eCrime Researchers Summit (APWG eCrime), ACM, 2007, pp. 1-13.

[20] D. Kevin McGrath and M. Gupta, “Behind phishing: An examination of phisher mod operandi,” In Proceedings of the USENIX workshop on Large-Scale Exploits and Emergent Threats, San Francisco, USA, 2008.

[21] K. T. Chen, J. Y. Chen, C. R. Huang, et al., “Fighting phishing with discriminative keypoints features,” IEEE Internet Computing, vol. 13, no. 3, 2009, pp. 56-63.

[22] A. Y. Fu, L. Wenyin, and X. Deng, “Detecting phishing web pages with visual similarity assessment based on earth mover's distance (emd),” IEEE Trans. Dependable Secur. Comput., vol. 3, no. 4, 2006.

[23] Z. Gyongyi and H. Gracia-Molina, “Web spam taxonomy,” 2005.

[24] M. Dunlop, S. Groat, and D. Shelly, “Using visual website similarity for phishing detection and reporting,” In Proceedings of the 5th International conference on Internet Monitoring and Protection, 2010.

[25] W. Liu, X. Deng, G. Huang, et al., “An antiphishing strategy based on visual similarity assessment,” IEEE Internet Computing, vol. 10, no. 2, 2006.

[26] T. Ronda, S. Saroiu, and A. Wolman, “Itrustpage: a user-assisted anti-phishing tool,” In proceedings of the 3rd ACM SIGOPS/EuroSys European Conference on Computer Systems, Glasgow, Scotland, 2008.