

A User Study of the Effectiveness of Typosquatting Techniques

Jeffrey Spaulding
University at Buffalo & CUBRC, Inc
Buffalo, NY, USA
jeffrey.spaulding@ubrc.org

Ah Reum Kang
University at Buffalo
Buffalo, NY, USA
ahreumka@buffalo.edu

Shambhu Upadhyaya
University at Buffalo
Buffalo, NY, USA
shambhu@buffalo.edu

Aziz Mohaisen
University at Buffalo
Buffalo, NY, USA
mohaisen@buffalo.edu

Abstract—The nefarious practice of cyber *typosquatting* involves deliberately registering Internet domain names containing typographical errors that primarily target popular domain names in an effort to steal their traffic for monetary gain. Typosquatting has existed for well over two decades and continues to be a credible threat to this day.

In this work, we discuss the results of a user study that exposes subjects to several uniform resource locators (URLs) in an attempt to determine the effectiveness of several typosquatting techniques that are prevalent in the wild. We also attempt to determine if security education and awareness of cybercrimes such as typosquatting will affect the behavior of Internet users.

Keywords. Domain Names, Typosquatting, Defenses.

I. INTRODUCTION

An alarming amount of domain names in the Domain Name System (DNS) are deliberately registered with typographical variations that target popular domain names [1]. *Typosquatting*, as this practice became known as, exploits common typographical errors made by users that manually type URLs into web browsers in an attempt to steal traffic or redirect users to unintended destinations. These so-called “typosquatters” employ several techniques (*e.g.*, adding or deleting characters) when typosquatting domain names in order to sufficiently capture enough traffic for monetary or personal gain.

In this work, we present the design and evaluation of a user study for gauging the effectiveness of several typosquatting techniques that are used in the wild. More specifically, we make the following contributions:

- To validate typosquatting techniques identified in prior studies by examining current trends and data sources.
- How security education and awareness of cybercrimes, particularly typosquatting, will affect the behavior of Internet users.
- How to leverage the existing countermeasures and cognitive traits of Internet users to strengthen the defense against typosquatted domains.
- Publicly releasing our data so others can verify and build upon our research.

II. BACKGROUND AND PRELIMINARIES

Over the years, studies have been conducted to understand typosquatting models, including various features of their domain names [2]. In the following, we review these various models and features prevalent in typosquatted domain names.

A. Typo-Generation Models

One of the first and widely cited approaches in the area of typo domain name generation was introduced in 2006 by Wang

et al. [3], where the following five typo-generation models were commonly used in the wild:

- 1) **Missing-dot typos:** the dot following “www” is removed, *e.g.*, `wwwSouthwest.com`.
- 2) **Character-omission typos:** one character is omitted, *e.g.*, `Diney.com` (a typo of the *Disney* brand).
- 3) **Character-permutation typos:** two consecutive characters are swapped, *e.g.*, `NYTiems.com`.
- 4) **Character-substitution typos:** characters are replaced by their adjacent ones on a specific keyboard layout, *e.g.*, `DidneyWorld.com` (“s” → “d”).
- 5) **Character-duplication typos:** characters are mistakenly typed twice, *e.g.*, `Googlle.com`.

Later studies, such as Banerjee *et al.* [4], looked at exhaustively generating typo domains using other methods:

- 6) **N-mod-inplace:** substitutes N characters in the original domain name with all possible alphabet letters.
- 7) **N-mod-inflate:** increases the length of a domain name (or URL) by N characters.
- 8) **N-mod-deflate:** removes N characters from the original domain name (or URL).

B. Features of Typosquatted Domains

Domain Name Length. Early observations showed that most typosquatted domain names had less than 10 characters [4]. However, it was later shown in [5] that no matter the length, typo domains within the Damerau-Levenshtein [6], [7] distance of one or adjacent-keyboard distance of one from popular domains were overwhelmingly typosquatted.

Domain Name Popularity. While Banerjee *et al.* [4] initially suggested that typosquatting decreases significantly with declining domain name popularity, newer studies by Szurdi *et al.* [8] and Agten *et al.* [9] concluded that 95% of typo domains target the “long tail” of the popularity distribution.

Effect of the Top-Level Domain (TLD). Since `.com` is the dominant TLD of all registered domain names, most studies confirm that `.com` domain names have a high chance of being typosquatted—either by modifying the second-level domain (SLD) portions (*e.g.* `googlle.com`) or creating a malicious counterpart in another separate TLD (*e.g.* `Netflix.om`).

III. STUDY: IDENTIFYING TYPO DOMAINS

Our user study presented subjects with a list of actual URLs with a subset of them deliberately “typosquatted”. The subjects were simply asked to select “Yes” or “No” if the given URL

appears to be a typosquatted domain name. To assess how prior knowledge and awareness of security concepts affect a user’s behavior, the user study encompassed three separate phases which incrementally introduced subjects to all of the typosquatting techniques discussed in §II-A.

A. Preliminary Experimental Results

A total of 34 participants completed all three phases of the survey over a one-week period, receiving their score (out of 200) for the number of correct responses after each phase.

Scores and Completion Time. With each phase, the average number of correct responses improved and the average response time decreased slightly. For Phase 1, the scores ranged from 78 to 186 correct responses with a mean, standard deviation and variance of 142.2, 23.6 and 557.1, respectively. In Phase 2, the minimum score increased to give us a range from 110 to 188 (Mean=147.1, s.d.=18.6, variance=345.7). Phase 3’s minimum score increased slightly to range of 117 to 183 (Mean=149.9, s.d.=15, variance=225).

Age. The ages of the participants ranged from 22 to 39 (Mean=25, s.d.=4.1, variance=16.5). Interestingly, younger participants generally scored higher than older participants across all phases of the study. However, surprisingly, while the younger participants scored higher, they also spent more time per question on average compared to their older counterparts.

Education Level. Of the participants, there was only 1 High School Graduate and one who reported some College Education. For the rest, 17 participants (50%) had a Bachelors degree, 13 (38.2%) had a Masters degree, and 2 held Ph.D. degrees (5.88%). Participants holding higher degrees of education actually *scored worse* than those with less education.

Familiarity of Security Concepts. On a scale of 1-5, only 1 participant chose “2”, 15 participants (44.1%) chose “3”, 14 participants (41.1%) chose “4”, and the remaining 4 participants (11.8%) chose “5”. The final results confirm that one’s familiarity with security concepts coincides with how well they performed as scores generally increased.

Domain Name Features. As expected, participants were more successful in correctly identifying typosquatted domain names that targeted popular domains. As for which typo model was the most “effective”, Figure 1 shows that participants were very likely to identify a typosquatted domain name that used *Model 1* (Missing-Dot Typo) as opposed to *Model 2* (Character-omission Typo) and *Model 6* (1-mod-inplace) which caused the most confusion.

Given our sample size of 200 domain names, 167 (83.5%) contained *all* alphabetic characters while 33 (16.5%) contained alpha-numeric characters. Naturally, participants were more likely to identify a domain name that contained all alphabetic characters as opposed to alpha-numeric characters.

Furthermore, we grouped our sample domain names by their TLD (as listed by the Internet Assigned Numbers Authority (IANA) [10]) and found that 119 (59.5%) fall into the “historic” TLD group (e.g., .com, .net), 75 (37.5%) fall into the “country-code” TLD group, and 6 (3%) fall into the “generic”

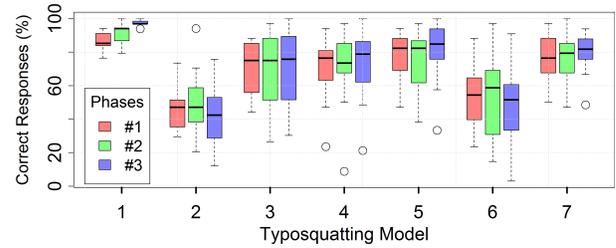


Fig. 1. Responses By Typo Model listed in §II-A.

TLD group. Not surprisingly, participants performed the best when presented with a domain name from a “historic” TLD.

IV. CONCLUSION AND IMPLICATIONS

This study has allowed us to gain valuable insight into the effectiveness of various typosquatting techniques and how security education affects a user’s behavior. Our results confirm that participants generally performed better and faster at identifying typosquatted domain names *after* being educated about typosquatting models between each phase of the study.

Our results also show a trend where older participants spent less time per survey question on average than their younger counterparts. This trend could explain why the younger participants scored better, as the older participants appeared less patient and tended to perform worse at typo identification.

Finally, our results indicated that users tend to correctly identify typosquatted domains that utilize models which add characters (e.g., duplicate or random). As Figure 1 depicts, the most effective typosquatting techniques involve permutations and substitutions. Studies in Cognitive Science, such as the work of Grainger and Whitney [11], highlight the “jumbled word effect” which demonstrates that the human brain can easily read words whose inner letters have been re-arranged.

An ongoing research we are pursuing is how to utilize the findings in this study to devise defenses for typosquatting that take users behavior and cognitive ability into account.

ACKNOWLEDGMENT

This work is supported by NSF grant CNS-1643207.

REFERENCES

- [1] A. R. Kang, J. Spaulding, and A. Mohaisen, “Domain Name System Security and Privacy: Old Problems and New Challenges,” <http://arxiv.org/pdf/1606.07080v1.pdf>, 2016.
- [2] J. Spaulding, S. Upadhyaya, and A. Mohaisen, “The Landscape of Domain Name Typosquatting: Techniques and Countermeasures,” in *Availability, Reliability and Security (ARES)*, 2016.
- [3] Y. Wang, D. Beck, and J. Wang, “Strider typo-patrol: discovery and analysis of systematic typo-squatting,” *USENIX SRUTI*, 2006.
- [4] A. Banerjee, D. Barman, M. Faloutsos, and L. N. Bhuyan, “Cyber-Fraud is One Typo Away,” in *IEEE INFOCOM*, 2008.
- [5] T. Moore and B. Edelman, “Measuring the perpetrators and funders of typosquatting,” in *FC*, 2010.
- [6] F. J. Damerau, “A technique for computer detection and correction of spelling errors,” *CACM*, 1964.
- [7] V. I. Levenshtein, “Binary codes capable of correcting deletions, insertions, and reversals,” *Soviet Physics Doklady*, vol. 10, p. 707, 1966.
- [8] J. Szurdi, B. Kocso, G. Cseh, M. Felegyhazi, and C. Kanich, “The Long Tail of Typosquatting Domain Names,” in *USENIX Security*, 2014.
- [9] P. Agten, W. Joosen, F. Piessens, and N. Nikiforakis, “Seven Months’ Worth of Mistakes: A Longitudinal Study of Typosquatting Abuse,” in *NDSS*, 2015.
- [10] —, “Root Zone Database,” <http://bit.ly/1TBSeck>, 2015.
- [11] J. Grainger and C. Whitney, “Does the huamn mnid raed wrods as a wlohe?” *Trends in cognitive sciences*, vol. 8, no. 2, pp. 58–59, 2004.