# APDs: Framework for Surveillance of Phishing words in Instant Messaging Systems Using Data Mining and Ontology

Mohd. S. Qaseem, Mohd. Nayeemuddin, M. M. Ali, Owais A.W. Siddiqui, and Md. Abdul Rafae

*Abstract* — Instant Messaging Systems (IMS) generically cannot detect many deceptive phishing attacks; hence they are vulnerable for cyber frauds. To overcome, we propose an Active-Phishing Detection System (APDs), developed using Data Mining (Rule-based) and Ontology. APDs monitor the user's psychology and predict type of the detected phishing activity with an alert to achieve zero-minute phishing attacks.

### I. INTRODUCTION

The APDs, dynamically predicts any potential deceptive phishing attacks, when instant messages are exchanged between users of an IMS. Currently, IMS lacks a stronger mechanism to deal with phishing at content-level. Few researchers proposed various methods to detect phishing attacks [1] [2]. But, lack the Rule-based method without Ontology [5], and trapping the users into deceptive phishing attacks [3]. The emails of Google are categorized into Primary, Social, Updates, and Forums, ignoring the issues of phishing attacks. Recently, leaked news of PRISM-NSA as one of the largest surveillance programs monitoring the Networks exploiting against cyber international law<sup>1</sup>.

#### **II. PROPOSED ACTIVE-PHISHING DETECTION SYSTEM**

The operational Framework of APDs is shown in Fig. 1. The APDs algorithm initiates the steps to capture the phishing words that are exchanged between the users and then stores them into database for identifying phishing words using pre-defined phishing rules of Table I. The APD algorithm is shown in Fig. 2. In APDs, the Monitoring system program identifies the culprit details of Phisher and report to the victim client. Steps of algorithm are illustrated as follows:



Fig. 1. Proposed Framework (APDs) to detect phishing messages from Instant Messaging Systems (IMS).

Mohammed. S. Qaseem is with Dept. of CSE, Nizam Institute of Engg. & Tech., Hyderabad, India (e-mail: ms\_qaseem@yahoo.com).

Mohammed Nayeemuddin is with Dept. of Informatics, Nizam College, Osmania University (e-mail: nayeemmca3@yahoo.co.in).

M.M. Ali, Owais A.W. Siddiqui & Md. Abdul Rafae are with Dept of CSE, MJCET, Hyd., India (e-mail:owais.aws@gmail.com, and



Fig. 2. Schematic cum algorithmic representation for proposed Framework named as APD algorithm.

- 1. In this step, the phishing words are identified by using *GSHL* and *Tree Alignment* Algorithms as discussed in [4]. These messages are stored in Text database (TDB), where unnecessary words are filtered using Ontology Based Information Extraction technique (OBIE) (stemming, N-gram technique, ignore words) [6] [7].
- 2. The frequently recurring words are extracted from the TDB dynamically using Association rule mining technique [8] and SSPWDB (pre-defined rules) guided with Ontology database (ODB), later these words are pushed to TPDB. The *metadata* is a gist of information related to instant messages.
- 3. Once the Phishing words are detected, the message is considered as suspicious, as given in Table I (*rule 1*). The *KDB* maintains the detected stem words along with the domain (i.e. type of Phishing activity).
- 4. Profile details are traced from *EDB*, which are provided during the creation of an email id, with the aid of Relational Wrapper Algorithm [9].
- 5. The email- id through which the phishing words are sent is tracked using *metadata*, and the victim is alerted.

Phishing Rule 1 Phishing words / scenario **Threat Type** Financial Account no, bankers, credit card details, dob pain name initials, fav mobile Fame and Password, aliases, hobbies, fav game, notoriety lesignation, prof goal, strengths Identity Kids name, pets name, fav food, fav teacher, access middle name, fav team, school name, pen name, pet name, alias. Deceitful Phony email, phony txtmsg, screen name, elicitation pet-name, phony url, multiple email acc, manage details, phony social networks Acc\_creation Emp id, all caps, special characters, tips alpha-numeric, common pwd Suspicious page URL or link contains- or URL related @, contains more than 5 dots, false login page, host-name portion is an IP address, organization's name (e.g. Ebay) in a URL path but not in the domain name, page is not accessible, inconsistencies are found LIDI ' Phishing Rule 2 (Threshold value) Check the Threshold value for the stem words using OBIE Framework Phishing Rule 3 (multi lingual & undetected words) Multilingual and undetected words to be checked and updated automatically PREDEFINED PHISHING RULES ERROR DETECTION ROOT WORD(S) EXTRACTED (DOMAIN) ONTOLOGY (TREE ALIGNMENT ONTOLOGY BASED EXTRACTED INFORMATION IFORMATION EXTRACTION MODULE(TPDW) Υ PREPROCESSOR (NLP TOOL ERROR DETECTION & TEXT INPUT (TDB)

Table I. Pre-defined Phishing rules

Fig. 3. Proposed OBIE Framework to detect phishing words from Instant Messaging Systems (IMS).

CORRECTION (DICTIONARY)

**III. EXPERIMENTAL RESULTS** 



Fig. 4. Shows Taxonomic Structure of words mapped from Text Pattern Database (TPDB) with pre-defined axioms (Set of Suspicious Phishing Words Data Base (SSPWDB))

The given Text Input (TDB) is converted to pure textual format by the preprocessor component using NLP tools. The information extraction techniques filter unnecessary words from unstructured text. The stem words that are found and stored in TPDB, are guided by ontology internally by the Ontology Generator component using TreeAlignment algorithm. During this process it makes use of semantic lexicon (WordNet) and Error detection and Correction, for stem words that are extracted and builds a Tree with empty root node initially. The input given from pre-defined phishing rules (SSPWDB) to Ontology Editor is again mapped with the Ontology Generator by identifying the exact Domain for the empty root node. Finally the root word(s) along with stem words are stored in knowledge database (KDB). The experimental results obtained for predicting the type of phishing activity detected from IMS are shown in Fig. 4, that shows Taxonomic Structure of words mapped from Text Pattern Database (TPDB) with pre-defined phishing axioms (SSPWDB). The experimental results obtained for predicting the type of phishing activity detected from IMS using OBIE framework of Fig. 3, is shown in Fig. 4.

## IV. ISSUES AND CHALLENGES

- 1. Deceptive phishing messages are sent in any format other than textual (Images, Audio, Video), then they are not detected as shown in Fig. 5.
- 2. Rules lack multilingual support for deceptive phishing detection.
- 3. Issue with the interpretation of a message written in multiple languages



Fig. 5. Deceptive phishing attack using images to avoid detection in IMS

## REFERENCES

- Mahmoud Khonji, Youssef Iraqi, and A. Jones, "Phishing detection: A [1] Literature Survey," IEEE Vol 15, 2013.
- [2] Mohd Mahmood Ali, and Lakshmi Rajamani, "Deceptive phishing detection system: From audio and text messages in instant messengers using data mining approach," IEEE, 2012.
- [Online] APWG www.antiphishing.org, accessed on 2014. [3]
- Mohd. Mahmood Ali, and L. Rajamani, "Framework for surveillance of [4] instant messages," IJITST, Inderscience publisher, 2013.
- [5] [Online] Ontology Portal - www.ontologyportal.org, accessed on 2014
- [6] Jer Lang Hong, "Data Extraction for Deep Web using WordNet," IEEE Transaction on Systems, Man and Cybernetics, 2011.
- [7] C.D. Manning, P. Raghavan, Hinrich Schutze, Introduction to Information Retrieval, 2004.
- R. Srikant, and R. Agarwal, "Mining quantitative association rules in [8] large relational tables," In Proceedings of the ACM - Special Interest Group on Management of Data (ACM SIGMOD), 1996, pp.1-12.
- [9] Sunitha Ramanujam, and et al., "A Relational Wrapper for RDF Reification," E. Ferrari et al. (Eds.): TM 2009, IFIP AICT 300, pp.196-214, IFIP International Federation for Information Processing 2009.