

Phishing email detection using Improved Recurrent Convolution Neural Network

¹SURESH KURA, ²G.NAGGAPA, ³ Dr.ANANTHA RAMAN.GR

¹M-TECH, DEPT OF CSE, MALLAREDDY INSTITUTE OF ENGINEERING AND TECHNOLOGY, DHULAPALLY MEDCHAL, SECUNDERABAD, TELANGANA, INDIA, 500014

²ASSISTANT PROFESSOR, MALLAREDDY INSTITUTE OF ENGINEERING AND TECHNOLOGY, DHULAPALLY MEDCHAL, SECUNDERABAD, TELANGANA, INDIA, 500014

³PROFESSOR, MALLAREDDY INSTITUTE OF ENGINEERING AND TECHNOLOGY, DHULAPALLY MEDCHAL, SECUNDERABAD, TELANGANA, INDIA, 500014

Abstract: Phishing is one of the serious security problems faced by cyber-world and leads to a lot of financial losses for both individuals and organizations. Phishing is an illegal act done by the attackers to retrieve confidential and personal information of the web users by betraying them. So, a more effective era of phishing detection is needed to reduce the risk of phishing emails. In this paper, we first look at the structure of email. So based on an advanced iteration, We proposed a enhanced recurrent convolutional neural network (ERCNN) model with multi-level vector and interest mechanism a new phishing emails detection form called ERCNN, which is used to generate emails in the email header, Email text, person score, and phrase level at the same time. To evaluate efficacy of ERCNN, we use an unbalanced dataset with realistic phishing and valid emails. The Experimental results show that the total accuracy of ERCNN is 99.848%. Meanwhile, an error the fine rate (FPR) is 0.043%. FPR high and low ensures that the filter remains above the spoofed emails with a high chance and filter legitimate emails as low as possible. This promising result is advanced Current detection methods and checking the effectiveness of THEMIS in detecting phishing emails.

Keywords: Phishing emails, phishing detection, Convolutional neural network (CNN), recurrent convolutional neural networks.

I. INTRODUCTION

Phishing is a social engineering attack that aims to exploit a vulnerability identified in device strategies as caused by machine clients. For example, the device might be technically comfortable enough against password theft, but users might mistakenly

filter their passwords if an attacker requests them to replace them via an HTTP hyperlink.) Given that in the long run the general protection of the device.

In addition, attackers can use technical vulnerabilities (e.g, cache poisoning) to build more compelling forms of social engineering messages (i.e. using valid, if fake) domains that can be more persuasive than using exceptional domains). This makes phishing attacks a multi-layered nuisance, and effective mitigation may require tackling technical and human problems.

Since phishing attacks aim to exploit specific vulnerabilities in people (i.e. customers who stop using machines), it is very difficult to mitigate them. For example, as evaluated in [1], customers who have stopped doing so have not detected 29% of phishing attacks, even if they are trained with the user's global focus program. On the other hand, phishing detection strategies are evaluated against massive phishing attacks, which makes their overall performance almost unknown when it comes to focused patterns of phishing attacks. These bottlenecks in phishing mitigation strategies have virtually resulted in a breach of protection against different groups alongside the leading data security operators [2], [3].

Phishers use Special techniques to lie to users and steal their monetary and non-public actions. Fake emails and creating fake websites are the most used strategies by them. E-mails are sent by attackers to clients who request their personal records by changing some components of the message and causing the person to accept it from an approved source. Online phishing is a method by which attackers create a fake or fake internet site that looks identical to the only site that tricks users into stealing

their data. Fake emails seem incredibly regularly, or even the website where the internet consumer is required to enter private information is very similar to legitimate mail. Phishing messages are posted by email, SMS, instant messengers, social networking sites, VoIP etc.

II. LITERATURE SURVEY

There are various techniques has been implemented in the recent years. Some of the works discussed below.

Choon Lin et al. [5] with a focus on discovering phishing webpages, phishing is a definition of online fraud because it misused users online for more than a decade. The authors suggested a phishing method that relies entirely on the approach to URL tags and likewise symptoms, which derive key phrases from the query site. With the help of identity keywords such as search phrases, the search engine is called to discover the destination domain call, which can be used to clarify the legitimacy of the question's web page. Like the authors, experiments were performed on more than 1,000 datasets, resulting in 99.20% true positive results and 92.20% true negative. The results confirm that the proposed anti-phishing method can correctly detect phishing webpages without using traditional language-based keyword extraction techniques.

Zhijun Yan et al. [6] were proposed a "new model for discovering a Chinese e-commerce site for phishing," which combines URL properties and net website capabilities. Some of the unique features of Chinese e-commerce sites are integrated, and the Sequential minimal Optimization Algorithm (SMO) is used to identify e-commerce sites for phishing. Meanwhile, the authors choose a genetic set of rules to improve the detection model. The effects of the evaluation of the proposed set of rules indicate that the overall performance of the SMO algorithm is more reliable than the reference model, and that GA has significantly expanded detection accuracy.

Shivam Aggarwal et al. [7] Phishing provides false social engineering strategies that extract accurate statistics from harmless victims. In this document, its planner proposes to locate phishing email messages that do not include a hyperlink, but that it completely

depends on the interests of the affected person, who wants to respond with personal statistics. We use fashion capabilities in all phishing emails, including not mentioning the victim's name in the email, a financial incentive, and a word that persuades the recipient to return. This text test may also be related to the email header evaluation so that final combined predictions can be made for each invoice. We have shown that this technology is much better than previous phishing detection e-mail strategies, because it includes e-mail messages without hyperlinks, while the current methods mainly depend on the concept of hyperlinks.

Mohammad et al. [8] they have proposed "Intelligent rule-based Phishing Websites Classification" Phishing defined as a function of following the website of a famous company that intends to collect the user's personal information, such as username, password and social security. Phishing websites include tips on their content as well as browser-based security ads. Several explanations have been proposed to detect phishing. Although there is no magic bullet that can explain this threat altogether. A promising technique can be used to predict phishing attacks based on data mining. In particular, "incorporating rating rules", since anti-phishing solutions aim to accurately predict the type of website, and these exactly match the rating data mine. In this approach, we have highlighted the essential features that distinguish phishing websites from the legitimate, and determine the rule-based data mining techniques for phishing website prediction.

Islam et al. [9] Phishing attacks continue to pose serious risks for consumers and businesses as well as threatening global security and the economy. Therefore, developing countermeasures against such attacks is an important step towards defending critical infrastructures such as banking. Although different types of classification algorithms for filtering phishing have been proposed in the literature, the scale and sophistication of phishing attacks have continued to increase steadily. In this paper, we propose a new approach called multi-tier classification model for phishing email filtering. We also propose an innovative method for extracting the features of phishing email based on weighting of

message content and message header and select the features according to priority ranking. We will also examine the impact of rescheduling the classifier algorithms in a multi-tier classification process to find out the optimum scheduling. A detailed empirical performance and analysis of the proposed algorithm is present.

Parmar et al [10] Faced with severe risks and a growing number of trade regulations, organizations are constantly challenged to ensure that security and compliance is adequate across the entire IT infrastructure. Although scams and tricks are nothing new, their pace and logistics have increased with the increasing dependence on the Internet, email and social networks worldwide. In particular, the proliferation of workplace email has not only helped companies succeed, but also opened the door to security risks

Hsu et al. [11] Phishing sites generate billions of dollars in profits by stealing non-public identities and private information. In this file, by examining URL structures and threading rankings, a URL Category method is proposed to prioritize suspicious URLs across a web page. Given the fact that the average time for phishing sites is low, it is very important that the proposed technology interact on time with phishing URLs while the URLs are valid. Since the approved approach does not consist of any crawling of the Internet or analysis of content materials, you can create the required signatures from the phishing URLs in real time. Moreover, the proposed approach uses only a few computer sources which, along with a small additional computer, can be integrated into any existing real-time URL evaluation device.

Yue et al. [12] Given many anti-phishing mechanisms currently Help users verify if a website is genuine. However, usable studies show that a preventative approach fails to effectively suppress lonely phishing attacks and prevents Internet users from displaying their credentials on phishing sites. In this document, a new method is proposed to protect human users from "bite-bashing" phishing attacks, rather than preventing them from "bite the hook". We developed BogusBiter, a specialized client-side anti-phishing tool that opens up a relatively large number of fake documents on the phishing site. Bogus Better

hides the victim's original credentials under false credentials and also allows a legitimate website to identify the stolen credentials in a timely manner. Taking advantage of the power of the client's automatic phishing detection technique, BogusBiter complements current phishing prevention methods. We implemented BogusBiter as an extension to the Firefox 2 Web browser, and evaluated its efficacy through real experiments on both phishing and legitimate Web sites

Sheng et al [13] studied the effectiveness of phishing black-lists. We used 191 fresh phish that were less than 30 minute sold to conduct two tests on eight anti-phishing tool bars. We found that 63% of the phishing campaigns in our dataset lasted less than two hours. Blacklists were ineffective when protecting users initially, as most of them caught less than 20% of phish at hour. We also found that blacklists were updated at different speeds, and varied in coverage, as 47% - 83% of phish appeared on blacklists 12 hours from the initial test. We found that two tools using heuristics to complement blacklists caught significantly more phish initially than those using only blacklists. However, it took a longtime for phish detected by heuristics to appear on blacklists. Finally, we tested the toolbars on a set of 13,458 legitimate URLs for false positives, and did not and any instance of mislabeling for either blacklists or heuristics. We present these findings and discuss ways in which anti-phishing tool scan be improved.

Basnet et al [14] Phishing is a form of identity theft that occurs when a malicious Web site impersonates a legitimate one in order to acquire sensitive information such as passwords, account details, or credit card numbers. Though there are several anti-phishing software and techniques for detecting potential phishing attempts in emails and detecting phishing contents on websites, phishers come up with new and hybrid techniques to circumvent the available software and techniques.

III. PROPOSED METHODOLOGY

The emails in this document are divided into categories, valid emails and phishing emails. Also detecting phishing emails is a trouble for binary rating. We improve inconvenience and cut email into

two components, the head and the frame. Most current documents are more practical in the world stage conveyor. However, the char-stage vector can better focus on spelling mistakes, personal spelling habits, big phrases and small words. All this is difficult to obtain by moving the stage of the word myself. For email body, email frame content can be very individual, for example, spelling behavior and big and small words, and some spelling errors are likely to occur. Therefore, we can take advantage of

SYSTEM ARCHITECTUE

character email writing features to distinguish phishing emails from legitimate emails. Also, the values in the email header are not all phrases, and there are also some static character sequences with delicate phishings, including the site name. A continuous person sequence with a unique blend is easily reached with a char conveyor. Therefore, similar to the use of the phase phrase vector, the char character vector is used.

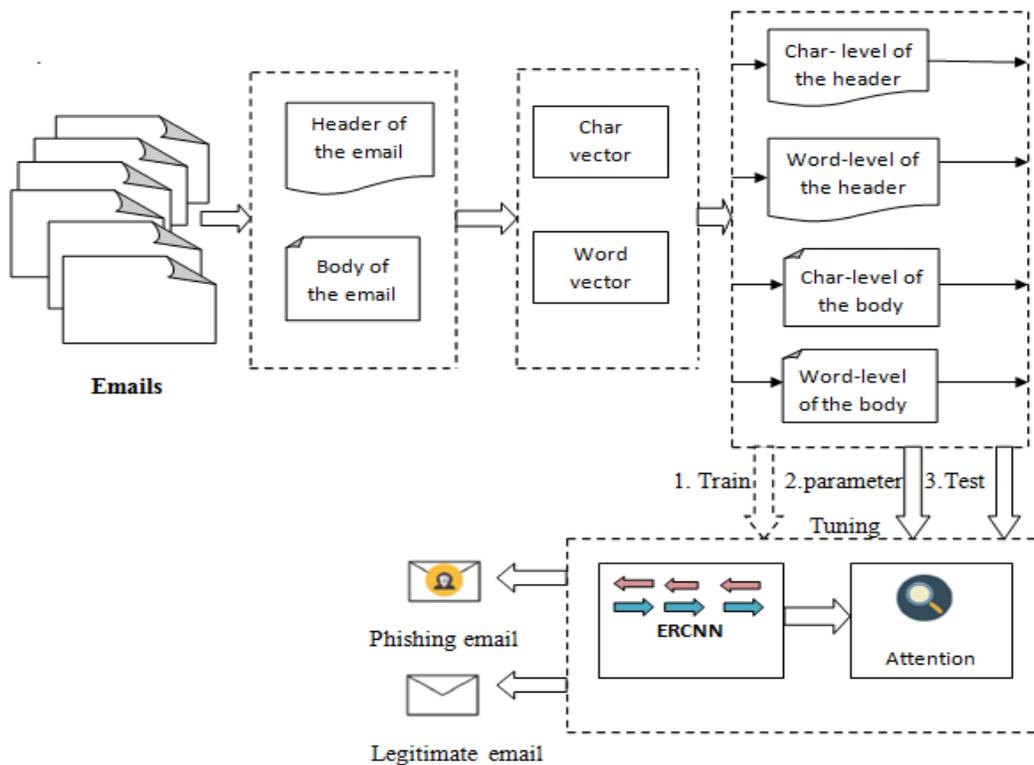


Fig.1 Proposed system architecture

In this paper, letter grading and sentence representation for e-mail is obtained using Word2Vec. It is a popular version suggested by Mikalov, which is used to create camel studs in text facts. Reproduces the linguistic context of words with the help of training the surface layer structure. Word2vec input is a large group, and the resulting outputs are vectors of some distinct dimensions. Each specific word (or person) in the set has a corresponding vector, which reflects the context of the phrase (or person). This makes mastering the

expression of phrases (or letters) much faster than the previous methods.

In summary, to represent the email, we can describe it from multiple levels of the char-level of the email header, char-level of the email body, word-level of the email header and word-level of the email body to better reflect all information contained in the email. The char-level vector embedding model and the word-level vector embedding model are obtained through Word2Vec tool training. Enter the email into these two models and the results are the vector

sequences of the char-level email header, the word-level email header, the char-level email body and the word-level email body.

The attention mechanism in deep learning is substantially similar to the selective visual attention mechanism of human beings, and the core goal is to select information more critical to the goal of the current task from much information. The enhanced RCNN and attention mechanism are combined to form an improved ERCNN-Attention model.

IV. RESULTS AND DISCUSSION

The proposed ERCNN performance evaluated using a well-known matrix, which includes accuracy, precision, recall. These measurements are calculated using

True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN). The following equations can be used to achieve this

$$Accuracy = \frac{TP+TN}{TP+FP+TN+FN} \times 100$$

Precision: Precision measure shows the amount predicted nodules that are actually related to the cancer.

$$Precision = \frac{TP}{TP+FP}$$

Recall: The recall is called the ability to classify positive class patterns. The following equations can be used to achieve this

$$Recall = \frac{TP}{TP+FN}$$

Table.1 Comparison between various methods

	Accuracy	Precision	Recall
LSTM	0.974	0.934	0.841
CNN	0.966	0.855	0.848
ERCNN	0.998	0.996	0.99

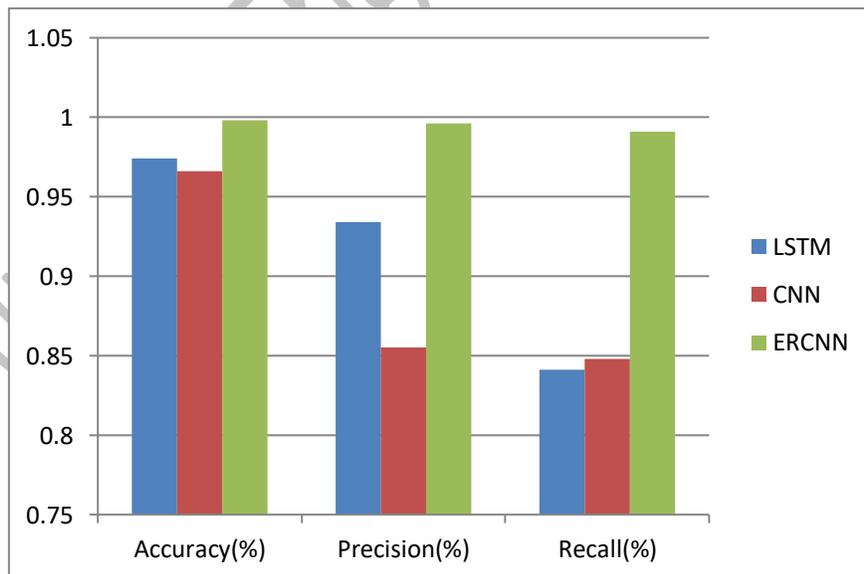


Fig.2 The summary of test results

To perform an objective evaluation, we implemented the methods as the proposed baseline ERCNN. As shown in Figure 2, the ERCNN version accuracy is 99.8%, precision is 99.6%, and recall is 99.0%. All offers are superior to CNN and LSTM strategies

V. CONCLUSION

In this paper, we use a new deep learning paradigm called the enhanced recurrent convolutional neural network (ERCNN) to detect phishing emails. The model uses an advanced RCNN network to generate the email header and text at the person and word level email. Therefore, noise in the model is added to a minimum. At launch, we used the head-and-body attention mechanism, which makes the model pay more attention to the additional valuable facts between them. We use an unbalanced dataset for the actual international scenario of behavioral testing and copy exam. The ERCNN model obtains a promising final result. Several experiments were completed to clarify the blessings of the proposed ERCNN model. For self-employment, we may be aware of how we can improve our version of phishing emails and not use an email header and the best email body.

REFERENCES

- [1] S. Sheng, M. Holbrook, 2010, "Who falls for phish?: a demographic analysis of phishing susceptibility and effectiveness of interventions", pp. 373–382.
- [2] B.Krebs,2011, " HBGray Federal hacked by anonymous", accessed December 2011.
- [3] B. Schneier, 2011, "Lockheed Martin hack linked to RSA's SecurID breach," <http://>, accessed December 2011.
- [4].Kaspersky Lab, 2013, "Spam in January 2012 love, politics and sport,".
- [5] Choon Lin Tan, Kang Leng Chiew, San Nah Sze , 2017, "Phishing Webpage Detection Using Weighted

URL Tokens for Identity Keywords Retrieval", pp 133-139.

[6] Zhijun Yan, Su Liu, Hangzhou Yang, 2016, "A Genetic Algorithm Based Model for Chinese Phishing E-commerce Websites Detection in HCI in Business".

[7] Shivam Aggarwal, S D Sudarsan, Vishal Kumar, 2015," Identification and Detection of Phishing Emails Using Natural Language Processing Techniques".

[8] R. M. Mohammad, F. Thabtah, and L. McCluskey, "Intelligent rule-based Phishing Websites Classification," 2014.

[9] Islam, R., & Abawajy, J. (2013). A multi-tier phishing detection and filtering approach. *Journal of Network and Computer Applications*, 36(1), 324–335.

[10] B. Parmar, "Protecting against spear-phishing," *Computer Fraud & Security*, vol. 2012, no. 1, pp. 8–11, 2012.

[11] C. H. Hsu, S. Pu, 2011, "Identify fixed-path phishing attack by STC," pp. 172–175.

[12]Yue, C., and Wang, G H. Bogusbite, "A transparent protection against phishing attacks", *ACM Transactions on Internet Technology (TOIT)* Vol.10 n.2, pp.1-31, 2010.

[13] Sheng G, Cranor, L, Hong, J, Zhang, C, "An empirical analysis of phishing blacklists". In: *Proc. 6th Conf. on Email and Anti-Spam (2009)*.

[14] Basnet, R., Mulkamala, S., Sung, A.: Detection of phishing attacks: A machine learning approach. In: *Soft Computing Applications in Industry*, pp. 373–383 (2008).