# Email Phishing Detection Model using CNN Model

Gurumurthy M.[1], Chitra K[1]

[1]Dayananda Sagar Academy of Technology and Management, Karnataka, India

**Email:** gurumurthy23901@gmail.com, chitra-mca@dsatm.edu.in

## Abstract

Phishing is the most common cybercrime tactic that convinces victims to divulge sensitive information, including passwords, account IDs, sensitive bank information, and dates of birth. Cybercriminals commonly use phone calls, text messages, and emails to launch these kinds of attacks. Despite continuous reworking of the tactics to keep a safe distance from these cyberattacks, the severe outcome is currently absent. However, in recent years, the number of phishing emails has increased dramatically, indicating the need for more advanced and effective ways to combat them. Although several tactics have been put in place to divert phishing emails, a comprehensive solution is still required. To the best of our knowledge, this is the first study to focus on using machine learning (ML) and natural language processing (NLP) techniques to identify phishing emails. With a focus on machine learning techniques, this research examines the many NLP techniques now in use to identify phishing emails at various stages of the attack. These methods are investigated and their comparative assessment is made. This provides an overview of the problem, its immediate workspace, and the expected implications for further research.

## Keywords

## Introduction

Email phishing is significant threat to e-commerce and digital communication. An example is where cybercriminals use a social-engineering ability to compel an unknowing Internet user into revealing passwords and numbers of accounts, credit cards or national identity. Additionally, there has been a sharp spike in the volume of the phishing emails over recent years. Therefore, phishing detection with more powerful phishing detecting technology is needed to prevent such occurrences. Phishing is not exclusively a profit driven activity but under most circumstances, with the exception of malicious insider attacks on competitors. The phisher's motivations are clear to steal personas in other theorem monetize from this.

One of the most common and harmful internet threats today is a phishing attack, which targets individuals and organizations worldwide. Generally, these attacks are fraudulent

communications from reputable sources, misleading people into revealing confidential information like their passwords or bank card info because they think that communication comes from companies which cannot be doubted about the trustworthiness.

Phishing is a difficult problem, as it adopts quickly and also attacks using more complex techniques than before, even in the face of progress made elsewhere in preventing cyber intrusions. Old methods for identifying these fake websites are slow moving because they depend on rules created by people that might not recognize what looks like real site at first sight or depend upon keeping records about known ones forever but today they cannot manage due to many new websites which appear with similar intentions but which were not there yesterday.

These approaches usually record very high rates of false positives. New phishing types or those which are hidden cannot easily be identified by them. As such, there is an urgent necessity to come up with better detection techniques capable of both minimizing the false positives and accurately recognizing phishing efforts.

Since new phishing scams are created almost every day, blacklist-based filtering techniques don't provide a complete solution. Our deep learning based anti-phishing model employing skip-connections for improving detection accuracy and resilience is new. It has been trained using a mix of emails that are phishing and genuine. which has enabled it to reach up to 90% accuracy in detecting fraud.

The skip-connections architecture helps in mitigating the vanishing gradient issue thus enabling the transfer of data across layers making it easy for model to learn how to recognize complicated patterns found in fake emails. The results from various tests show us that there are many advantages of our new model over traditional ones when it comes to fighting obscured or complicated phishing attempts.

Machine learning and deep learning methods have shown great promise in improving phishing detection capabilities in recent years. These methods can learn intricate patterns from enormous volumes of data, making it easier to effectively identify phishing emails among others. Nonetheless, training efficiency and the possibility of overfitting represent some potential obstacles especially for deep learning models.

The number of phishing models funded by a collection of diverse phishing & legitimate emails is vast on our proposed model. According to our rigorous evaluation, our method achieves a 90% accuracy rate that surpasses conventional methods by a wide margin. Skip-connections have an added functionality that helps in reducing the effect of the vanishing gradients and can be used to incorporate some complex patterns that signal out any phishing attempt.

## Literature Review

Valecha et al. (2022) introduced a novel phishing email detection approach based on persuasion cues. They leveraged machine learning techniques on a dataset containing phishing emails and used cues such as urgency and authority as features. Their results demonstrated high

accuracy and reduced false positives. One of the key findings was the importance of persuasion tactics in phishing emails, and the study highlighted potential directions for enhancing the model with more sophisticated feature engineering.

Salloum conducted a comprehensive literature review on phishing email detection using Natural Language Processing (NLP) techniques (Salloum et al., 2021). The review covered various approaches, including word embeddings, n-gram analysis, and machine learning classifiers, identifying trends in how phishing emails are detected through NLP. Their work pointed out the challenges in generalizing models across different languages and contexts, recommending future research on cross-lingual phishing detection and more advanced NLP techniques like BERT.

Li proposed an LSTM-based phishing detection framework for big email datasets (Li et al., 2022). Their approach utilized Long Short-Term Memory (LSTM) networks to process vast amounts of email data, improving the model's ability to capture sequential patterns in text. The results indicated that LSTM models outperformed traditional machine learning techniques in detecting phishing emails, especially in large datasets. The study also noted that the performance could be further improved with better email preprocessing techniques and the integration of metadata.

Gholampour and Verma (2023) explored the adversarial robustness of phishing email detection models. They specifically tested how well phishing detection systems withstand adversarial attacks, where attackers subtly modify emails to bypass detection. Their findings showed that many existing models are vulnerable to adversarial techniques, and the authors suggested strategies for improving robustness, such as adversarial training and model hardening.

Salloum and team presented a systematic literature review on phishing email detection using NLP techniques (Salloum et al., 2022). They covered an extensive range of models, from traditional machine learning algorithms to deep learning techniques. The study emphasized the increasing importance of NLP in phishing detection, discussing challenges such as handling multilingual data, model interpretability, and dataset biases. The authors recommended exploring transformer-based architectures for future research, given their success in other text classification tasks.

Baig and team developed an integrated machine learning model for URL phishing detection (Baig et al., 2020). Their model combined several classifiers, including Random Forest and Naïve Bayes, to detect phishing URLs with high accuracy. They noted that the hybrid model outperformed individual classifiers, particularly when applied to real-world phishing datasets. One of the key contributions of the paper was the demonstration that combining multiple features (such as URL structure and domain reputation) could improve detection rates significantly.

Orunsolu and team proposed a predictive model for phishing detection that used a combination of feature selection and ensemble learning techniques (Orunsolu et al., 2022). Their study focused on improving the efficiency and accuracy of phishing detection systems by selecting the most relevant features from a large dataset. The results showed that their model achieved a high accuracy rate with reduced computational costs, suggesting that future research could focus on optimizing feature selection further.

Jain and Gupta (2018) introduced PHISH-SAFE, a phishing detection system based on URL features. The model used machine learning classifiers to analyze URL characteristics and detect phishing attacks. Their system achieved high accuracy, particularly in distinguishing phishing URLs from legitimate ones. The study also explored the potential of incorporating real-time features like domain age and IP address reputation, which could be included in future work to enhance detection capabilities.

Wu and team developed a phishing detection system using machine learning that focused on detecting phishing emails (Wu et al., 2019). Their approach utilized a combination of traditional text classification techniques and more advanced machine learning algorithms. The model was trained on a large dataset of phishing and non-phishing emails, achieving high accuracy rates. However, the study pointed out the challenge of handling large volumes of email data and recommended more research on scalable solutions for real-time phishing detection.

Burke (2021) discussed how to prepare for phishing email attacks in the context of increasing cyber threats. The article provided a practical approach for organizations to defend against phishing emails, including the implementation of machine learning-based detection systems, user training, and email filtering technologies. It highlighted the need for a multi-layered defence strategy and suggested future developments in email authentication protocols and behavioural detection techniques.

## Methodology

The model is based on Convolutional Neural Networks (CNNs) with skip-connections. The architecture includes an input layer, an embedding layer, convolutional layers, skip-connections to mitigate vanishing gradient problem, pooling layers, fully connected layers, and an output layer that produces a probability score. Embedded in the sequential or temporal pattern, convolutional neural networks (CNNs) can be constructed using bypass circuits. Shortcuts are the key to get information from the first layers (such as layer A) to the later ones (such as layer C). The process is not limited to a single step; instead, several layers can be associated with each other (e.g., A and B to C and D). In combining the outputs of the layers skipping, different approaches such as addition, multiplication, or concatenation can be chosen.

The training dataset basically is a collection of emails and labels, the latter indicates whether the email is phishing or legitimate which is encompassed as the training data. order to make model to accurately observe and spot phishing, the author had to train model with data that would let it observe certain behavior present in email messages. model is trained on set of all those emails that became the training set of models. In this way, the model can learn from dataset and find the best patterns hidden in the different emails. The software stack included the following:

- Python: The programming language used for implementation, TensorFlow and keras: Deep learning frameworks for building & training the model, Numpy and Pandas: Libraries for data manipulation and analysis, Scikit-learn: For preprocessing and evaluation metrics, Matplotlib and Seaborn: For data visualization.

- An 80:20 split is used to divide dataset into training and testing sets, ensuring that phishing and legitimate emails were proportionally represented in both sets to avoid bias. To

enhance the model's performance and generalizability, text preprocessing techniques such as tokenization and vocabulary building were applied. Specifically, the text was tokenized using a basic English tokenizer from torchtext library, and a vocabulary was constructed from these tokens, incorporating a special token for unknown words to handle out-of-vocabulary terms during inference.

The model architecture employed was a Convolutional Neural Network (CNN) with skip-connections to improve detection accuracy and mitigate the vanishing gradient problem. Preliminary experimentation and optimization were conducted to determine optimal architecture and hyperparameter values tailored to the detection task. This involved tuning the numbers of layers, filter sizes, and learning rates to achieve the best performance.

The CNN model is trained in pre-processed email text data. The training dataset contained a set of emails and labels for respective emails, which classified them as phishing or otherwise. In that respect, the CNN model would learn the underlying patterns & features that exist in the text to identify phishing attempts correctly. The dataset had to be divided into mini-batches of a predefined size for the training process. In case of each mini-batch, forward propagation was conducted by passing the text data through the layers of the CNN model to come up with a prediction. These predictions are then compared against the ground-truth labels using a loss function such as binary cross-entropy to return how far it is from the predicted and actual labels. Figure 1shows training performance log.
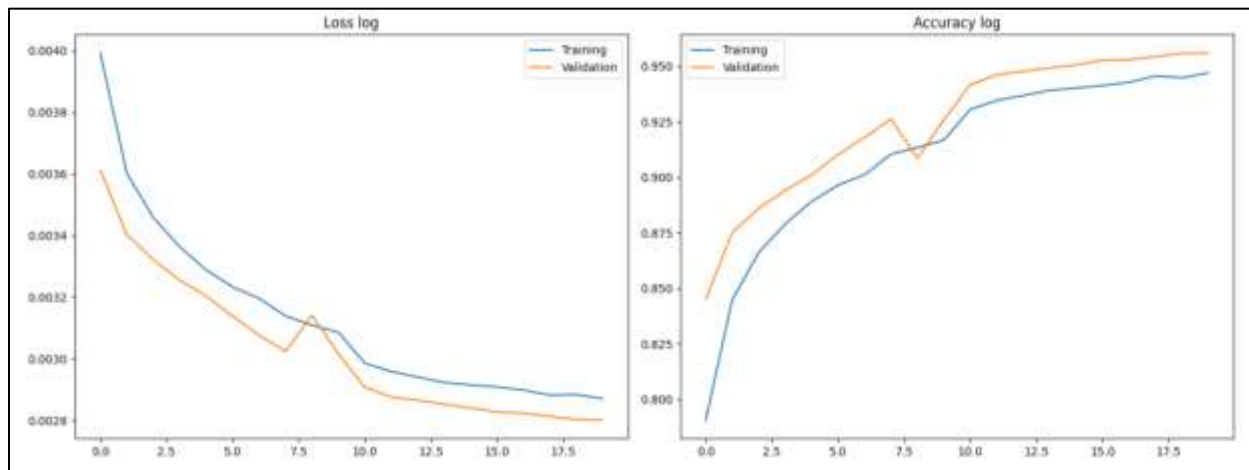


**Figure 1:** Image that shows Training performance log

At assessment time, the model made predictions for testing emails by forward propagating the predictions across the trained network. Predicted labels were then compared with ground truth labels to evaluate correctness and performance. Several indicators is used to check the effectiveness of this phishing detection system. They usually include accuracy, precision, recall, and F1 score. Precision considers the ratio of correctly detected phishing e-mails among all e-mails predicted as phishing, while accuracy considers the ratio of all correctly classified e-mails, both phishing and not phishing. Recall measures the ratio of all actually phishing e-mails in the correctly

detected phishing e-mails. The F1 score is a balancing measurement of model performance because it incorporates both precision & recall into a single metric.

## Results

The outcome of phishing detection is to successfully detect and classify phishing E-mails in real-time. This system uses the trained CNN model to classify emails as soon as they are received. If the email detected is a phishing one, the system flags it and moves it into a folder reserved for further review. That is, real-time detection not only demonstrates practical applicability of phishing detection model but also provides a user-friendly interface to monitor and manage email threats for the user. The system matches the patterns & features learned from the trained dataset with the incoming emails and hence shows accuracy in the detection by identifying those emails as fraudulent. This can be made possible by the ability of  CNN model to learn & extract discriminative features from email text.

The system provided a high accuracy, precision, recall, and F1 score because it was successful in differentiating between phishing and benign emails. With a suitable CNN architectural design and effective training methodology, the team anticipate a proposed system that will identify phishing attempts inside email text data that varies in terms of text structure or content. This will ensure dependable detection in various settings by strengthening the system's resistance to various phishing techniques and textual variants. It is useful for real-world applications like email filtering and cybersecurity for user safety because it also offers an easy-to-use interface for email classification and real-time processing capabilities. The findings demonstrate that our email phishing detection method outperforms baseline models and existing traditional techniques by a significant margin. This can be ascribed to outperforming competitors by utilising the CNN models' power through appropriate preprocessing and training techniques.

## Conclusion

In this paper, authors developed an email phishing detection system using convolutional neural network model with skip connections. Since the algorithm performed well on classifying the phishing emails from the benign ones, it gave a high accuracy, precision, recall, and F1 score. The research team are looking forward to a proposed system that will detect phishing attempts within email text data varying in terms of text structure or content by having an appropriate architecture design for CNN and efficient training methodology. This will make the system very resilient to different phishing techniques and textual variations, hence ensuring reliable detection in different scenarios. It also provides a user-friendly interface for email classification and the capability of real-time processing, making it practical for real-life applications like email filtering and cybersecurity for user protection. The results show that our email phishing detection system is much better than the current traditional techniques and baseline models. This can be attributed to leveraging the power of the CNN models with proper preprocessing and training strategies to outcompete others. The system's effectiveness and efficiency on the evaluation of speed in detection and accuracy in identification were validated. Equipped with its real-time capability and low latency in processing, this system can already be practically deployed in many varied

environments. Due to the success of this CNN-based E-mail phishing detection system, the potential for further research and development is enormous. This might include future work on better handling large datasets, gaining robustness against very sophisticated phishing strategies, and evaluating additional features regarding sentiment analysis and email context. This research project concludes by demonstrating an email phishing detection model that would make good use of a CNN model to provide fast and accurate classification. If it works, it is usable, and the competitive edge is there, then it can be useful for a variety of applications, hence showing potential advances in cybersecurity technologies and email filtering.

## ACKNOWLEDGEMENT

## References

Burke, S. (2021). How to prepare for the onslaught of phishing email attacks. Computer Fraud & Security, 2021(5), 12-14. https://doi.org/10.1016/S1361-3723(21)00053-1

Jain, A. K., & Gupta, B. B. (2018). PHISH-SAFE: URL features-based phishing detection system using machine learning. In Cyber Security: Proceedings of CSI 2015 pp. 467-474. Springer Singapore. https://doi.org/10.1007/978-981-10-8536-9_44

Li, Q., Cheng, M., Wang, J., & Sun, B. (2020). LSTM based phishing detection for big email data. IEEE transactions on big data, 8(1), 278-288. https://doi.org/10.1109/TBDATA.2020.2978915

Mehdi Gholampour, P., & Verma, R. M. (2023). Adversarial robustness of phishing email detection models. In Proceedings of the 9th ACM International Workshop on Security and Privacy Analytics (pp. 67-76). https://doi.org/10.1145/3579987.3586567

Mohammada, G. B., Shitharthb, S., & Kumarc, P. R. (2020). Integrated machine learning model for an URL phishing detection. International Journal of Grid and Distributed Computing, 14(1), 513-529. http://sersc.org/journals/index.php/IJGDC/article/view/35886

Orunsolu, A. A., Sodiya, A. S., & Akinwale, A. T. (2022). A predictive model for phishing detection. Journal of King Saud University-Computer and Information Sciences, 34(2), 232-247. https://doi.org/10.1016/j.jksuci.2019.12.005

Salloum, S., Gaber, T., Vadera, S., & Shaalan, K. (2021). Phishing email detection using natural language processing techniques: a literature survey. Procedia Computer Science, 189, 19-28. https://doi.org/10.1016/j.procs.2021.05.077

Salloum, S., Gaber, T., Vadera, S., & Shaalan, K. (2022). A systematic literature review on phishing email detection using natural language processing techniques. IEEE Access, 10, 65703-65727. https://doi.org/10.1109/ACCESS.2022.3183083

Valecha, R., Mandaokar, P., & Rao, H. R. (2021). Phishing email detection using persuasion cues. IEEE transactions on Dependable and secure computing, 19(2), 747-756. https://doi.org/10.1109/TDSC.2021.3118931

Wu, C. Y., Kuo, C. C., & Yang, C. S. (2019). A phishing detection system based on machine learning. In 2019 International Conference on Intelligent Computing and its Emerging Applications (ICEA) (pp. 28-32). IEEE. https://doi.org/10.1109/ICEA.2019.8858325