

A MACHINE LEARNING AND NATURAL LANGUAGE PROCESSING-BASED SMISHING DETECTION MODEL FOR MOBILE MONEY TRANSACTIONS

*Aaron Zimba**, *Katongo O. Phiri*, *Chimanga Kashale*, *Mwiza Norina Phiri*

School of Computing, Technology, & Applied Sciences, ZCAS University,
Zambia

* Corresponding Author, e-mail: gvsfif@gmail.com

Abstract: As mobile services proliferate to include financial transactions, the threat of phishing attacks targeting users has equally been escalating. Attackers have been using different kinds of phishing techniques, especially in third world where mobile services are prevalent. As such, this paper presents a Smishing (SMS phishing) Detection model leveraging Natural Language Processing (NLP) and Machine Learning (ML) techniques. It aims to detect smishing threats in real-time by the integrating NLP with ML. The developed model harnesses NLP algorithms to analyse text-based messages, scrutinizing linguistic patterns and contextual clues indicative of smishing attempts. Through ML algorithms, the model learns to distinguish between legitimate (Non-Smishing) and fraudulent messages (Smishing), adapting dynamically to evolving smishing tactics. The model's efficacy is evaluated through comprehensive testing, demonstrating promising results of precision, recall, and accuracy with F-1 measure at 0.902 and AUC at 0.95. The Model stands as a proactive defence mechanism against smishing in mobile money environments, contributing to enhanced user security and trust in financial transactions.

Key words: machine learning, natural language processing, smishing, mobile money, phishing.

1. INTRODUCTION

The proliferation of smartphones has become a global phenomenon, notably in Zambia, where the liberalization of the telecommunications sector has catalysed widespread adoption [1]. This surge in smartphone usage has led to a paradigm shift in communication dynamics, with Short Message Service (SMS) and instant messaging notifications becoming the predominant means of notifications for financial transactions. Mobile money technology has emerged as a transformative force in financial transactions through mobile phones. Leveraging Short Message Service (SMS) technology, mobile money allows users to conduct various financial and banking

operations, including buying airtime, paying utility bills, managing savings, and engaging in mobile banking [2].

However, the increasing adoption of mobile money services has also given rise to new challenges, notably the surge in smishing attacks. Smishing, a portmanteau of "SMS" and "phishing," involves cybercriminals sending deceptive text messages to individuals, posing as reputable entities with the intent of acquiring sensitive information for financial exploitation. This form of social engineering [5] leverages human behavioural vulnerabilities to manipulate users into compliance, and attackers often exploit SMS as a trusted source during the exchange of confidential information. Despite the escalating threat of smishing attacks in the mobile money landscape, existing detection methods predominantly focus on email phishing and web-based attacks [6]. The unique behavioural characteristics of smishing in SMS messages within the context of mobile money services have been largely neglected. This research addresses this gap by proposing a natural language processing and machine-learning-based [4] detection model specifically tailored to classify Bemba and English smishing text messages targeting mobile money users in Zambia.

The research focuses on detecting smishing attacks in English and Bemba language datasets, limited by the survey results and targeting mobile money services in Zambia, specifically MTN Mobile Money and Airtel Money. The proposed detection model aims for broad usability across various mobile devices and operating systems, emphasizing real-time identification and classification of smishing attacks to enhance mobile money transaction security. By analysing text content within SMS messages, the model ensures accessibility and cost-effectiveness for a diverse audience of mobile money users, prioritizing practicality over prevention frameworks.

This research proposes a tailored smishing detection model for mobile money services, using natural language processing and machine learning. By enhancing transaction security, it aims to protect users' financial interests and personal information, making significant contributions to cybersecurity and linguistics. This research also advances understanding in language-specific smishing detection, offering insights for future research in these fields.

The remainder of the paper is organised as follows: Section 2 presents the related literature whilst Section 3 presents the Methodology and developed model. Experiments and data are brought forth in Section 4 while the results and the discussions thereof are presented in Section 5. The conclusion is drawn in Section 6.

2. RELATED LITERATURE AND CONCEPTS

The predominant usage of the SMS service in mobile money transaction has necessitated its widespread adoption even in the third world. Zambia, for example, has experienced significant growth, as evidenced by the volume of transactions reported by the Bank of Zambia. The statistics reveal an 89.7% increase in Mobile Money transactions from 2012 to 2022 [3], far surpassing the growth in Electronic Funds Transfer (EFT) as depicted in Figure 1.

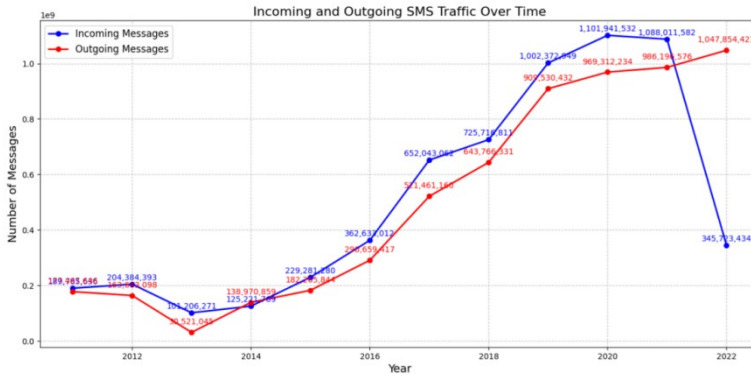


Fig. 1. SMS Traffic Visualization from 2011 to 2022

This has attracted cybercriminals who employ a subtype of phishing in social engineering through the use of SMS, hence Smishing. Over the years, the domain of Smishing detection has predominantly relied on a combination of methodologies including blacklisting, heuristics, and visual analytics [7], serving as the foundation for identifying and mitigating Smishing threats across diverse communication systems and digital platforms. While these methods have been foundational for mitigating smishing threats, they have limitations. Perpetrators have adeptly exploited vulnerabilities in existing solutions, devising techniques that circumvent conventional security measures.

A rule-based method that applies predefined rules to evaluate each SMS passing through an SMS gateway was proposed in [8]. Despite the attempt to establish rules, the ever-changing nature of attackers, who frequently alter mobile numbers, renders blacklisting and whitelisting techniques ineffective. They introduced a feature-based approach, identifying ten distinct features to differentiate Smishing messages from legitimate ones. Two features encoded as '0' represent legitimate messages, '1' for Smishing messages, while the remaining eight features predominantly indicate Smishing messages. Experimental results demonstrated a classifier with notable performance metrics: a true positive rate of 94.2%, a true negative rate of 99.08%, and an overall detection accuracy of 98.74%.

The 'S-Detector,' a system designed to identify Smishing attacks, comprising components such as an SMS monitor, SMS analyzer, SMS determinant, and Database was proposed in [9]. Employing a Naïve Bayesian Classifier, the authors evaluated both URL and content of text messages, identifying commonly used words in Smishing messages.

In a comprehensive investigation into spear phishing attacks, a natural language processing (NLP) detection algorithm to identify SMS spear phishing attacks was implemented in [10]. Collaborating with 360-mobile-safe, the study utilized a dataset of 31 million real-world spam messages associated with spear phishing. Employing Word2Vec and TFIDF vectorization techniques, the study achieved a remarkable F1-Score of 93.41% with the Logistic Regression and Word2Vec combination.

A Naïve-Bayes algorithm to classify spam communications targeting Kenyan mobile money customers is utilised in [11]. By gathering English-language spam mails and employing the Weka toolbox, the study achieved an accuracy rate of 96.1039%.

A prototype system utilizing the Backpropagation Algorithm and comparison with conventional classifiers such as Random Forest, Decision Tree, and Naïve-Bayes is presented in [12]. The Backpropagation Algorithm demonstrated superiority, achieving an impressive accuracy rate of 97.93% in classifying Smishing messages through a two-phase system: domain checking and SMS classification.

A framework designed to identify Smishing and vishing attacks associated with mobile money transactions, outlining recommended actions for customers when faced with such attacks is presented in [13].

A Smishing management system based on trust management principles, aiming to regulate and screen Smishing attempts by assessing trust relationships between message senders and recipients is introduced in [14].

A classification method for detecting phishing on mobile devices, encompassing various types of mobile device phishing, such as Bluetooth phishing, SMS phishing, voice phishing, and mobile web application phishing is employed in [15]. The researchers introduced and conducted comparisons of technologies designed for detecting mobile device phishing. Table 1 below summarises comparisons with other works.

Table 1. Comparisons with related works

Work	Classifier	Domain	Language	Modelling Approach	Evaluation Metrics
Joo et al. [9]	Naïve Bayesian	Smishing	English	Machine Learning	Accuracy
Liu et al. [10]	Logistic Regression	Smishing	English	Natural Language Processing	Prec, Recall, FN, FP, and F1-Score
Mishra & Soni [12]	Back-propagation	Smishing	English	Deep Learning	Acc., AUC & Exec. Time
Kipkebut et al. [11]	Naïve Bayesian	Smishing	English	Machine Learning	Prec., Recall, Acc., TP, FN, TN & FN-Rates
Mambina [16]	Random Forest	Smishing	Swahili	Machine Learning	Log-Loss, AUC & Exec. time
Proposed Model	RF & Naïve Bayesian	Smishing	Bemba & English	NLP & ML	MCC, Prec., Recall, F1-score

The predominant challenges faced by existing smishing detection approaches lie in the limitations of conventional methods such as blacklisting, heuristics, and visual analytics. Rule-based techniques, exemplified by Jain and Gupta's method [8], encounter shortcomings as attackers frequently change mobile numbers, rendering the rules ineffective. Feature-based approaches, as demonstrated by Jain and Gupta [8], have shown promise but are not foolproof in distinguishing smishing messages from genuine ones. While some studies like Joo et al. [9] and Liu et al. [10] introduce innovative methods employing Naïve Bayesian Classifier and natural language processing (NLP),

respectively, there remains a need for context-specific models tailored to the dynamics of smishing attacks in the context of mobile money services. The proposed smishing detection model aims to address these challenges by leveraging advanced natural language processing and machine learning techniques to enhance the accuracy and adaptability of smishing detection in the specific scenario of mobile money transactions.

3. METHODOLOGY

3.1. Research Framework and Approach

This study employs a hybrid descriptive research design [17] to investigate Smishing behaviours among mobile money users and to develop a detection model against smishing. This design approach integrates Natural Language Processing (NLP) techniques and Machine Learning (ML) models to develop a robust Smishing detection system tailored for mobile money transactions. This methodology relies on the synergy between NLP algorithms for text analysis and a suite of ML algorithms, including Naive Bayes, Random Forest, and Logistic Regression.

We justify this method selection based on the demonstrated effectiveness of these algorithms in handling unstructured text data, prevalent in Smishing attempts. NLP algorithms decode textual nuances, while ML models excel in identifying patterns within the dynamic landscape of mobile money transactions.

3.2. Implementation and Model Enhancement

The chosen methodologies provide a direct pathway for analysing textual content from SMS messages and transaction logs. Through the strategic application of Natural Language Processing (NLP), semantic meanings embedded within the text are extracted, enabling the detection of suspicious Smishing patterns.

The integration of Machine Learning (ML) models within these methodologies plays a pivotal role. Trained to detect patterns and features, these models efficiently categorize incoming messages into two distinct categories: legitimate messages (Non-Smishing) and Smishing attempts. The conceptual framework is depicted in Figure 2

3.3. Smishing Detection Model

The proposed Smishing detection algorithm follows a systematic process of preprocessing SMS data, extracting informative features, training machine learning models, and evaluating their performance to identify smishing messages.

By incorporating NLP techniques and machine learning, the algorithm aims to enhance the accuracy of smishing detection and mitigate the risks associated with SMS phishing attacks. This is shown in Algorithm 1 below. The algorithm begins by tokenizing the SMS messages (in English or Bemba) to break them down into individual words or tokens, followed by normalization to ensure consistency across the data. Feature engineering is then performed to extract relevant features, including computing TF-IDF scores to quantify the importance of each word in each message relative to the entire dataset.

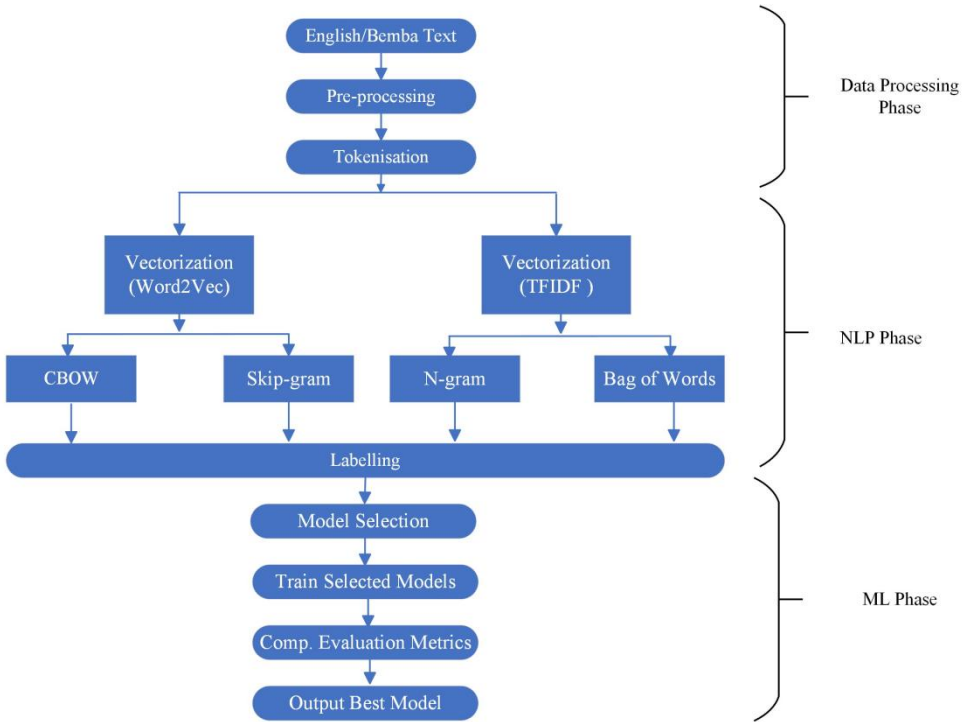


Fig. 2. Smishing Detection Conceptual Framework

Algorithm 1: Smishing detection algorithm

Input: $x_{E-d}^i = \{x_{e-d}^1, x_{e-d}^2, x_{e-d}^3, \dots, x_{e-d}^i\}$, labelled English dataset where $x_{e-d}^i \in \mathbb{R}^n, i = 1, 2, 3, \dots, n$

$x_{B-d}^i = \{x_{b-d}^1, x_{b-d}^2, x_{b-d}^3, \dots, x_{b-d}^i\}$, labelled Bemba dataset where $x_{b-d}^i \in \mathbb{R}^n, i = 1, 2, 3, \dots, n$

Output: TP & FP rates – Smishing detection accuracy

1. Read the labelled data x_{E-d}^i and x_{B-d}^i
 2. Tokenise data elements $T(D)$: $T(x_{E-d}^i/x_{B-d}^i) = \{t_1, t_2, t_3, \dots, t_n\}$
 3. Normalise the original data to $\overline{x_{E-d}^i}$ and $\overline{x_{B-d}^i}$
 4. Feature engineering & compute $TF - IDF = TF(t, d) * IDF(t)$
 5. Train pre-processed data and engineered features with **RF, NB, & LR**
 6. Generate TP & FP and other smishing detection accuracy rates
 7. Return smishing detection accuracy rates
-

Machine learning models Random Forest, Naive Bayes, and Logistic Regression are trained using the pre-processed data and engineered features to learn patterns distinguishing between legitimate and smishing messages.

Finally, the trained models are evaluated on a separate test dataset to generate metrics including True Positive and False Positive rates, providing insights into the performance of the algorithm.

4. EXPERIMENTS AND DATA

This research relies on two labelled datasets comprising confirmed Smishing messages and non-Smishing ones samples shown in Table 2 and Table 3, serving as fundamental components of the study.

Table 2: English Smishing dataset

Label	Text
Smishing	To send that money use this number of <i>MNO</i> xxxxxxxxxx the name will come Edward Sichula. my number is not working in Airtel money. Thanks
Smishing	Please call me now. The money for CDF and youth empowerment is out's
Smishing	Ok use this airtel number to send that money name will come Joyce. My number is not working in mobile money

Table 3: Bemba Smishing dataset

Label	Text
Smishing	Please call me now. mupoke indalama sha C.D.F ishabalanda naba youth. tumeni NRC registration number mupoke indalama shenu.
Smishing	indalama sha cdf na youth impowerment nashifuma tumeni NRC number yenu mupokeko ulupiya. Tumeni phone.
Smishing	Natukwata gold amasaka yabili tuleshitisha tumeni indalama

These datasets play a pivotal role in training and validating models and algorithms aimed at discerning between malicious Smishing attempts and legitimate messages. They represent a critical resource for developing effective strategies to detect Smishing, contributing significantly to the advancement of cybersecurity measures. Table 2 and Table 3 display the first five rows of the Smishing English Dataset and the Smishing Bemba Dataset respectively, offering a glimpse into the nature of the data used for analysis and model development.

In this paper, we present a hybrid algorithm depreciated in Algorithm 1 that incorporates specific algorithms, including Random Forest, Naive Bayes, and Logistic Regression, chosen for their suitability in natural language processing (NLP) and machine learning tasks, particularly in analysing text and classifying data. Random Forest's ensemble learning approach effectively captures complex interactions among features [18], while Naive Bayes excels in simplicity and efficiency, making it adept at handling textual data [19]. Logistic Regression, meanwhile, offers effectiveness in scenarios where the relationship between features and the target variable is linear [20]. These algorithms are leveraged to develop a comprehensive smishing detection system, supported by diverse model characteristics and a range of NLP techniques such as tokenization, stop words removal, stemming, and TF-IDF for feature extraction. Feature selection methods and machine learning algorithms, including Random Forest Classifier with Synthetic Minority Over-sampling Technique (SMOTE), further contribute to addressing imbalanced classes and optimizing performance.

Random Forest: Achieved the highest performance with cross-validation scores between 0.805 and 0.902. It also identified important features for prediction and demonstrated high precision, recall, and F1-score (0.902).

Naive Bayes: Showed a slightly lower performance compared to Random Forest with cross-validation scores ranging from 0.725 to 0.878. The F1-score was 0.852, and precision and recall metrics indicated room for improvement.

Logistic Regression: Performed well with cross-validation scores between 0.800 and 0.927, achieving an F1-score of 0.902 (identical to Random Forest). Its performance mirrored Random Forest's based on the confusion matrix and classification report.

Table 4: Model Performance Metrics

Metric	Random Forest	Naïve Bayes	Logistic Regression
Matthews Correlation Coefficient (MCC)	0.822	0.743	0.822
F1-score (Non-Smishing)	0.9	0.86	0.9
F1-score (Smishing)	0.9	0.84	0.9
Precision (Non-Smishing)	0.83	0.76	0.83
Precision (Smishing)	1	1	1
Recall (Non-Smishing)	1	1	1
Recall (Smishing)	0.82	0.73	0.82

The analysis of model performance suggests Random Forest as the preferred choice for smishing detection. While both Random Forest and Logistic Regression excelled, Random Forest offers several advantages. Its ensemble approach promotes better generalization and robustness, making it suitable for handling the complexities of smishing data. While interpretability might be slightly lower compared to Logistic Regression, insights can still be gained through feature importance analysis. Additionally, Random Forest has the potential for further optimization through hyperparameter tuning, making it a strong overall choice for this application.

Receiver Operating Characteristic (ROC) Curve

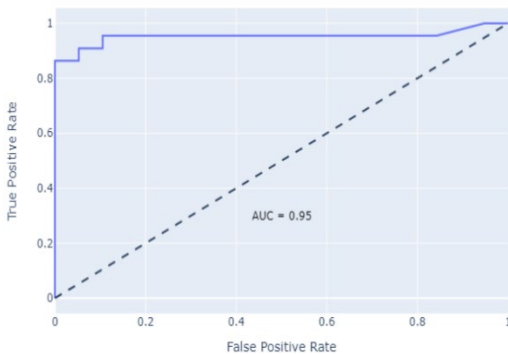


Fig. 5. RF ROC Curve

Receiver Operating Characteristic (ROC) Curve

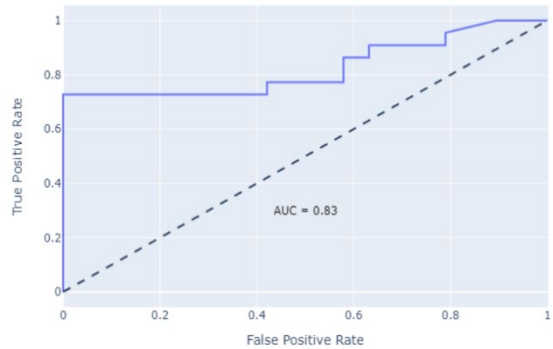


Fig. 6. NB ROC Curve

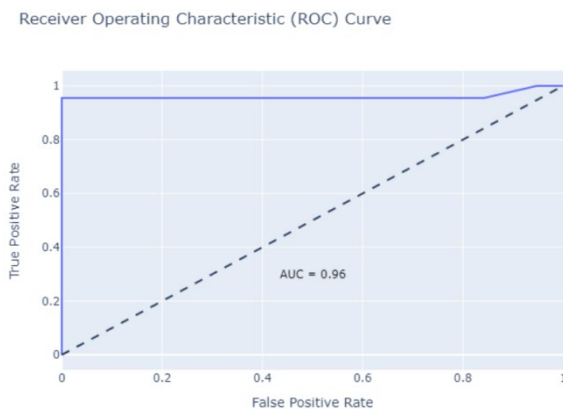


Fig. 7. LR ROC Curve

The ROC curves above reiterate the superior performance of Random Forest. The novelty of this research lies in its focus on detecting smishing attacks within the context of mobile money transactions, particularly targeting SMS messages in both English and Bemba languages in Zambia and other applicable countries. The concrete contribution of this study is the development of a robust, real-time smishing detection model that integrates Natural Language Processing (NLP) and Machine Learning (ML) techniques. This model not only enhances the security of mobile money transactions by accurately identifying and classifying smishing messages but also addresses a critical gap in the existing literature, which predominantly focuses on email and web-based phishing.

6. CONCLUSION

The exponential rise of mobile money transactions globally has introduced unprecedented convenience but also attracted the looming threat of smishing, a form of SMS phishing targeting users through deceptive messages. This study emphasizes the critical need for robust security measures, advocating for the development of effective smishing detection models utilizing Machine Learning (ML) and Natural Language Processing (NLP) techniques to safeguard users and financial transactions. Findings regarding mobile network usage highlight the potential risk exposure of dominant networks like Airtel and MTN, while linguistic analysis sheds light on language preferences and awareness levels among users, underscoring the urgency of educational initiatives to combat smishing threats. Despite challenges such as limited data access and awareness among respondents, the study contributes significantly to the body of knowledge by identifying linguistic markers unique to smishing messages, evaluating model performance, and emphasizing the robust performance of the Random Forest model in accurately identifying smishing attempts. The adoption of Random Forest is justified by its superior performance metrics, including high Matthews Correlation Coefficient (MCC) of 0.822, F1-scores of 0.9 for both Non-Smishing and Smishing classes, precision of 0.83 for Non-Smishing and 1 for Smishing, and recall of 1 for Non-Smishing and 0.82 for Smishing, underscoring its effectiveness in detecting smishing attempts compared to Naïve Bayes and Logistic Regression models. However,

limitations such as restricted data access and the need for further model optimization highlight the evolving nature of ML models and the ongoing quest for improvement in predictive analytics.

REFERENCES

- [1] Zimba A, Mbale TF, Chishimba M, Chibuluma M. Liberalisation of the International Gateway and Internet Development in Zambia: The Genesis, Opportunities, Challenges, and Future Directions. *arXiv preprint arXiv:2102.10629*. 2021 Feb 21.
- [2] Ghosh I, O'Neill J. The unbearable modernity of mobile money. *Computer Supported Cooperative Work (CSCW)*. Vol.29, no. 3, Jun 2020, 227-61.
- [3] Zambia Information and Communications Technology (ZICTA), *2023 Annual Market Report: A Supply Side Assessment of Developments in the ICT Sector*. March 2024.
- [4] Savyanavar AS, Mhala N, Sutar SH. Star-galaxy classification using machine learning algorithms and deep learning. *International Journal on Information Technologies and Security*. Vol.15, no.2, 2023, pp.87-96.
- [5] Romansky R. A survey of informatization and privacy in the digital age and basic principles of the new regulation. *International Journal on Information Technologies and Security*. Vol. 11, no. 1, 2019, pp.95-106.
- [6] Jain AK, Gupta BB. A survey of phishing attack techniques, defence mechanisms and open research challenges. *Enterprise Information Systems*. Vol.16, no. 4, 2022, pp.527-65.
- [7] Timko D, Rahman ML. Commercial anti-smishing tools and their comparative effectiveness against modern threats. *Proceedings of the 16th ACM Conference on Security and Privacy in Wireless and Mobile Networks*, 29 May 2023, pp. 1-12.
- [8] A. K. Jain, B. B. Gupta. Feature based approach for detection of Smishing messages in the mobile environment. *Journal of Information Technology and Research*, vol. 12, no. 2, 2019, pp. 17–35.
- [9] Joo, J.W., Moon, S.Y., Singh, S., Park, J.H.: S-Detector: an enhanced security model for detecting Smishing attack for mobile computing. *Telecommunication Systems*, vo.66, no.1, 2017, pp.29–38
- [10] 21.M. Liu, Y. Zhang, B. Liu, Z. Li, H. Duan, D. Sun. Detecting and characterizing SMS spearphishing attacks. *Annual Computer Security Application Conf.*, 2021, pp. 930–943.
- [11] A. Kipkebut, M. Thiga, E. Okumu. Machine learning SMS spam detection model. Proc. Kabarak Univ. Int. Conf. Comput. Inf. Syst., C. M. Maghanga and M. Thiga, Eds., Nakuru, Kenya: Kabarak Univ., Oct. 2019, pp. 63–70.
- [12] S Mishra, D Soni, (2019) SMS phishing and mitigation approaches. In: Twelfth International Conference on Contemporary Computing (IC3), Noida, India pp. 1–5, doi: 10.1109/IC3.2019.8844920
- [13] B. M. Nturihi, *A mobile money social engineering framework for detecting voice & SMS phishing attacks—A case study of M-Pesa*. Ph.D. dissertation, United States Int. Univ. Africa, Nairobi, Kenya, 2018.

- [14] L. Chen, Z. Yan, W. D. Zhang, R. Kantola. TruSMS: A trustworthy SMS spam control system based on trust management. *Future Generation Computer Systems*, vol. 49, Aug 2015, pp. 77–93.
- [15] Foozy, C. F. M., Ahmad, R., Abdollah, M. F. Phishing detection taxonomy for mobile device. *International Journal of Computer Science Issues*, vol.10, no.1, 2013, pp.338–344
- [16] I. S. Mambina, J. D. Ndibwile, K. F. Michael. Classifying Swahili smishing attacks for mobile money users: A machine-learning approach, *IEEE Access*, vol. 10, 2022, pp. 83061-83074, doi: 10.1109/ACCESS.2022.3196464.
- [17] Baran ML. Mixed methods research for improved scientific study. *IGI Global*; 2016 Mar17.
- [18] Cao Y, Geddes TA, Yang JY, Yang P. Ensemble deep learning in bioinformatics. *Nature Machine Intelligence*. Vol. 9, no. 2, 2020:500-8.
- [19] Dixit A, Mani A, Bansal R. Feature selection for text and image data using differential evolution with SVM and naïve Bayes classifiers. *Engineering Journal*. Vol. 24, no.5, 2020, pp.161-72.
- [20] Ray S. A quick review of machine learning algorithms. *2019 Int'l conference on machine learning, big data, cloud and parallel computing (COMITCon) 2019 Feb 14*, pp. 35-39, IEEE.

Information about the authors:

Aaron Zimba, PhD – is the Dean of the School of Computing, Technology, & Applied Sciences at ZCAS University and obtained his PhD in Network and Information Security at the University of Science and Technology Beijing. His research interests include network & information security, machine learning, and artificial intelligence. He's an IEE member.

Katongo Ongani Phiri – is a lecturer at ZCAS University in the department of Information Technology Systems. He holds a master's degree in IT from the ZCAS university. His research interests span Machine Learning and AI, IoT, Cybersecurity, & Mobile Applications.

Chimanga Kashale – is the current HoD for the Computer Science department at ZCAS university where he is pursuing his PhD degree. He holds a master's degree from the Copperbelt university, and his research interests include IoT and artificial intelligence.

Mwiza Norina Phiri – is the current HoD for Information Technology Systems at ZCAS university where she is doing her PhD studies. She holds a master's degree from the University of Zambia and her research interests include cybersecurity and ICT management.

Manuscript received on 07 May 2024