

Deep Learning-Based Speech and Vision Synthesis to Improve Phishing Attack Detection through a Multi-layer Adaptive Framework

*

1st Tosin Ige
Dept. of Computer Science
The University of Texas at El Paso
Texas, USA
toige@miners.utep.edu

2nd Christopher Kiekintveld
Dept. of Computer Science
The University of Texas at El Paso
Texas, USA
cdkiekintveld@utep.edu

3rd Aritran Piplai
Dept. of Computer Science
The University of Texas at El Paso
Texas, USA
apiplai@utep.edu

Abstract—The ever-evolving ways attacker continues to improve their phishing techniques to bypass existing state-of-the-art phishing detection methods pose a mountain of challenges to researchers in both industry and academia research due to the inability of current approaches to detect complex phishing attack. Thus, current anti-phishing methods remain vulnerable to complex phishing because of the increasingly sophistication tactics adopted by attacker coupled with the rate at which new tactics are being developed to evade detection. In this research, we proposed an adaptable framework that combines Deep learning and Randon Forest to read images, synthesize speech from deep-fake videos, and natural language processing at various predictions layered to significantly increase the performance of machine learning models for phishing attack detection. To validate both the effectiveness and adaptability of our proposed framework in overcoming limitations in current approaches and its ability to detect complex phishing site, we created 4 categories of phishing sites and uploaded them to a secure server with a compromised DNS on a friendly URL; the first was a text-only phishing site, image-only phishing site, video-only phishing site, and a phishing site combining all the features. We use SEO friendly URLs, and hacked legitimate DNS on the text-only phishing site, so that they can evade detection at 1st layer until the 4th layer of the framework where they were detected, we also created phishing sites where text are in image only format, text-only, and video only format using deep-fake video to test the adaptability of our proposed framework to different scenarios of a sophisticated or complex phishing site, our proposed framework successfully overcome limitations in existing approaches, significantly improve phishing attack detection, and successfully detect complex phishing webpages with multi-dimensional deep-fake videos, images, and texts.

Index Terms—Phishing, Random Forest, Deep Learning, Recurrent Neural Network, Long Short-Term Memory, Speech Synthesis, Vision Synthesis, Phishing Detection Framework, Adaptive Framework

I. INTRODUCTION

The insufficiency of traditional phishing detection methods such as user education [27] and rule-based methods [13]

against sophisticated phishing attack techniques has led researchers to exploration of possible AI-based solutions. While several machine learning-based models have been proposed, the fact that attackers use advanced innovative methods that are continuously changing to carry out phishing attacks renders previously proposed machine learning models ineffective against sophisticated attacks [11]. Although tools like PhishTank, and OpenPhish were created for the effective detection of malicious Uniform Resource Locator (URL) the rate at which malicious websites are created coupled with the sophistication of the deception method easily overwhelm the system as phishing sites are being created every 11 seconds according to dataprot 2023 phishing statistic report.

One of the common problems with the current ML model is the quality of the dataset used for the training model which has a significant impact on both the accuracy and overall performance of the model [30]. These data do not reflect the ever-changing strategies through which attackers continue to fool existing machine learning-based models to evade detection. In addition, the balancing problem between human factors and model accuracy causes illegal flagging of newly registered legitimate websites due to weak domain authority. Phishing websites have a short life span as they are quickly taken down before detection and another one is created, It is the rate at which an existing phishing website is taken down after launching a campaign and the immediate creation of a new phishing website [6] to begin another campaign coupled with the ever-changing but sophisticated techniques that makes the problem very potent and significant.

While several models and machine learning-based frameworks have been proposed, the ever-evolving ways attacker increases the sophistication of phishing attack to bypass existing state-of-the-art anti-phishing detection and prevention systems pose a mountain of challenges leading to the relative ineffectiveness of previously proposed models against a more complex phishing attack. Thus, the constant evolvment and innovation in phishing techniques adopted by attackers are

the reason why current detection method remains vulnerable to complex or more sophisticated forms of phishing due to their reliability on [1], [5], [8], blacklists/whitelists [9], natural language processing [15], visual similarity [15], rules [14], [24], remains vulnerable to attack due to the following reasons;

- Having understood how the machine learning-based model works, attackers are now increasingly relying on asymmetrical methods by uploading images and videos to evade detection under various pretexts, and none of the proposed models can single-handedly be effective against such.
- Very small or minute changes to the uniform Resource Locator (URL) of a blacklisted URL will make the blacklist/whitelist phishing detection method fail. Also, the fact that there is no worldwide centralized database for whitelisted or blacklisted URLs makes this method even more vulnerable, and so if company X blacklisted my phishing URL on their internal server, I can try it with company Y and be successful.
- In machine learning phishing detection method that relies on relevant features like URL, webpage content, website traffic, search engine, WHOIS record, and Page Rank have their vulnerabilities because firstly, such classifier will misclassify a phishing URL that is hosted on a hacked or compromised server as benign leading to false negative, secondly using domain age as a feature to train a model will always lead to higher false positive simply because the URL of a newly registered legitimate company website will be misclassified. After all, the domain name was recently registered, the page rank is zero, and with low traffic, and thirdly the fact that parameters for those features are gotten from a third-party website is another concern. What will happen if the third-party website is having a downtime?
- The issue with the visual similarity-based heuristic method which compares both the pre-stored signature such as images, font styles, page layout, screenshot, and so on of the new website with the old website will have general difficulty in detecting anomalies in a newly hosted phishing site.
- The fact that the majority of the existing machine learning models are trained based on textual features such as “#”, “”, Internet Protocol address, URL Length, domain levels, and so on from the Uniform Resource Locator (URL) does not help as any phisher or attacker with little web technologies can develop what we called “friendly URL” depending on the programming language adopted whether JAVA, C#, Python, PHP or framework to avoid all those features. With a friendly URL, such models are bound to misclassify leading to an increment in false negative rate.

For any Machine learning-based phishing detection method to be effective in real-time combat against phishing attacks, it must address each of the stated reasons above for which existing state-of-the-art anti-phishing methods continue to be

vulnerable due to the increasingly sophisticated techniques by which phishing attacks are being carried out. It is worth noting that past research work on phishing attack detection had been largely based on approaches, classification, etc. RASHA ZIENI et al. [35] focus their review on list-based, similarity-based, and machine learning-based categories of approaches for phishing detection to identify pending research gap, Angad et al. [21] focus theirs on the advantages and limitations of existing approaches to phishing detection, while also using discussion of related application scenarios as guidance to propose a new method of anti-phishing detection, Yifei Wang [32] categorizes widely used phishing detection methods into seven categories and summarizes them. All previously proposed models, approaches, and frameworks have common limitation, there limitation was that they are either text-based or URL-based which makes it difficult for them to detect complex phishing attack where the attacker uses deep-fake videos, deep-fake images, textual-based images, or combination of any with traditional textual content.

In this research, we first reviewed some of the most recent works on phishing detection, and state-of-the-art algorithms from the past 5 years to investigate the performance of state-of-the-art machine learning and deep learning classifiers for phishing detection tasks, before proposing a multi-layered adaptive framework that uses computer vision to read images on a phishing webpage, and condense videos from a webpage to audio before synthesizing the speech into a condensed text to increase detection of a phishing attack. We use a combination of random forest algorithm and Long-Short Term Memory (LSTM) at different layers of the framework for effective coordination. The contributions of our research include the proposal of an adaptive multi-layered framework that uses computer vision to read graphic images, synthesize speech from uploaded videos, and natural language processing at various predictions layered to significantly increase the performance of machine learning-based models for phishing detection. Our artifacts which consist of source code, dataset, images, videos, and audio files for this research had been uploaded to a public GitHub repository for reproducibility of our research. Artifact can be found on GitHub at;

Deep Learning-Based Speech and Vision Synthesis to Improve Phishing Attack Detection through a Multi-layer Adaptive Framework and also at Code Ocean computational research platform, with the exception of the internally generated deep-fake video and audio data files for privacy.

II. RELATED WORK

A. Natural Language Processing (NLP)

NLP-based models use existing relationships between sentences, words, or letter parts of a language in a given text dataset. This made us explore the possibility of synthesizing an uploaded video from a phishing webpage to feed our neural network model. NLP architectures use modeling, preprocessing, and feature extraction:

Data preprocessing: It is imperative for text in a given dataset to be preprocessed into a pattern that the model

can easily understand because preprocessing effectively turns every character and word in the dataset into a format that the machine learning classifier can understand to extract useful patterns or learn from them. The fact that algorithms learn from data and the quality of the dataset used in training an ML model directly impacts the performance of that model making AI to be data-centric, and hence, priority is given to data preprocessing during NLP.

NLP stemming and Lemmatization

$$\begin{aligned}
 & a < b \\
 & ab < bc \\
 & abc > bcd \\
 & abcd > bcde
 \end{aligned}
 \quad \text{Lemma ??} \quad (1)$$

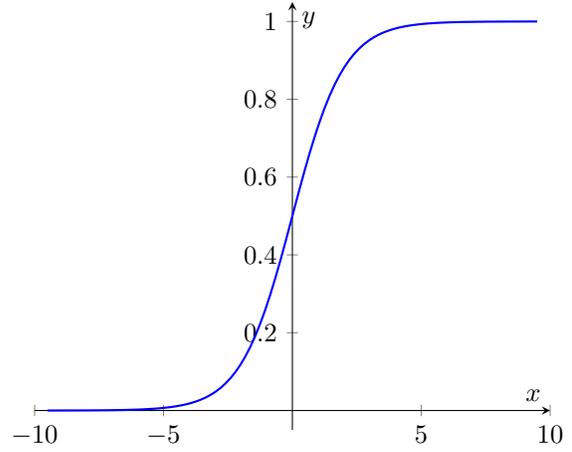
Stemming and lemmatization are the two major data preprocessing tasks for natural language processing. During stemming, there is an end-to-end iteration of each word in the dataset to convert them to their base forms such as the mapping of "university" to "univers", and "calamity" to "calam" while lemmatization uses the word's morphology from vocabulary dictionary to find their corresponding roots.

$$[T_i = \begin{cases} 1, & T \leq 1 \\ 1 + \beta T, & T > 1 \end{cases}, \text{ in which } T = \begin{cases} T_{\text{now}} - T_{\text{last}}, & T_{\text{last}} \neq \text{NULL} \\ T_{\text{now}} - T_{\text{update}}, & T_{\text{last}} = \text{NULL} \end{cases} \quad (2)$$

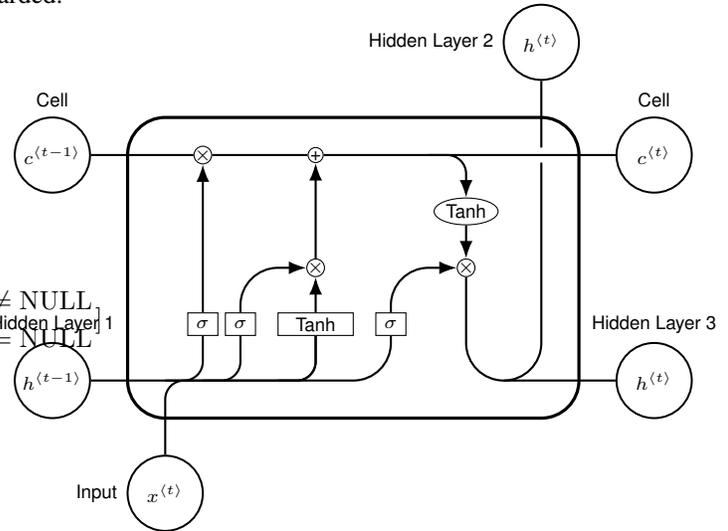
The final preprocessing stage of NLP is sentence segmentation, this process breaks large text into linguistically meaningful sentences where trivial words such as "an," "the," "a," etc that don't add much meaning or information to the text are removed during stop word removal, next we use tokenization to split every text into words and fragments, the result is a combination of word index and tokenized text which could be represented by a numerical token before feeding them to any of the deep learning or machine learning models for prediction.

B. Long Short-Term Memory (LSTM)

For this research, We opted for Long Short-Term Memory a variant of recurrent neural network (RNN) because of its effective solution to vanishing and exploding gradients which are Long-term dependency problems in Recurrent Neural Networks. The most important functioning part of an LSTM network is the cell state which serves as a memory to the network thereby enabling it to remember the past. Hence their suitability for capturing long-term dependencies and sequence prediction problems [3]. LSTM network has an input gate, a forget gate, and an output gate which are sigmoid activation functions with an output value of 0 or 1.



It was easy to use the sigmoid function as a gate because we are only given out positive values that could give a straight answer on whether a particular feature should be kept or discarded.



In an LSTM network, the Input Gate tells what new information is to be stored in the cell state, the forget gate gives clear instructions by telling what information is to be thrown away from the cell state, while the output gate gives activation to the output for more accurate prediction. It is during this activation which occurs after filtering the cell state that the output goes through the activation function where the output portion to be predicted is determined, and this occurs when the current LSTM block goes through softmax layer to predict value for the current block.

To mitigate the effect of phishing attack, several methods, frameworks had been proposed for phishing attack detection but with varying results, these methods are classified based on their different approaches which we classified as Non-Machine Learning, machine learning (Bayesian-based, non-Bayesian-based) and deep learning-based. As attackers continue to navigate potential vulnerabilities to existing phishing detection solution, they are beginning to rely on several images and uploaded videos rather than traditional text to enable them to evade detection, the inability of existing machine learning-based model to detect such phishing site is a peculiar

limitation to existing AI-based solution. Palla Yaswanth and V. Nagaraju [33] used Huang and Premaratne data from Kaggle repository with an equal number of phishing and legitimate datasets for novel network of phishing predictions with an accuracy of 95% for naive Bayes and 94.67% for random forest based on parameter turning. During the comparison of the performances of naive Bayes [10] and random forest for detection of phishing sites in a network, there was no testing of the model against sophisticated form of phishing attack and causes of the 5% failure rate of naive bayes in the research.

Abdul Karim et al. [19] proposed a hybrid model which combines logistic regression, support vector machine, and decision tree in conjunction with soft and hard voting, the proposed hybrid model used Grid Search Hyper-parameter Optimization, cross fold validation, and canopy feature selection method to select relevant features from the dataset. The proposed hybrid model resulted in an accuracy of 98.2% by using the only attribute properties of the uniform resource locator. The sole reliance on the attribute of the URL makes this approach extremely vulnerable to URL manipulations as any attacker with little experience in web technology can use a malicious webpage with a friendly URL to fool the model.

Ishwarya et al. [12] proposed a phishing detection method comprising of Naive Bayes algorithm, SVM, KNN, and random forest including evaluation of the performances of each of the four (4) classifiers in detection of phishing email. The implementation of each classifiers resulted in the highest accuracy of 98.2% for naive Bayes, albeit the use an imbalance dataset comprising 87% ham and 13% spam for the research surely indicate biased in the proposed model, and the problem of Bayesian poisoning was not addressed in the proposed model.

Kamal Omari [26] used the UCI phishing domains dataset to proposed machine learning-based model for the purpose of investigating Logistic Regression (LR), k-Nearest Neighbors (KNN), Support Vector Machine (SVM), Naive Bayes (NB), Decision Tree (DT), Random Forest (RF), and Gradient Boosting for phishing detection task. Hence, we believed that the 98.1% accuracy for phishing detection task obtained from the Naive Bayes classifier by Ishwarya et al. (2023) [12] was due to a massive imbalance in the dataset having 87% ham and 13% spam which was not addressed, also the proposed model doesn't address detection evasion through uploaded video and images on a phishing webpage.

Ann Zeky et al. [20] proposed an extraction-based Naive Bayes model for phishing detection with emphasis on the extraction of relevant features like unusual characters, spelling mistakes, domain names, and URL analysis from unseen web pages for effective classification of a website into malicious and benign. By training the proposed model with a relatively balanced dataset of 7000 records in which 54% are malicious and 46% are benign leading to an accuracy of 99.1%. By using a combination of content extraction and URL analysis, we believed the proposed model would not be vulnerable to malicious URLs in the sense that even if the attacker tried to use a friendly URL to deceive the model, that the model

does not rely on the properties of URL alone but also uses background webpage extraction means the proposed model will still be able to classify webpages correctly, albeit an attacker will still be able to use Bayesian poisoning.

Nishitha et al. [23] compared performances of machine learning algorithms and deep learning for phishing detection classification by implementing KNN, Decision tree, Random Forest, Logistic Regression as machine learning algorithm, convolutional neural network and recurrent neural network as deep learning in which logistic regression and CNN had the best performances with an accuracy of 95% and 96% respectively, albeit the proposed model only uses the URL properties and so couldn't be used for a sophisticated phishing attack that relies on images and video content.

Twana and Murat [22] while assuming the absence of a single solution to detect most phishing attacks and to investigate the impact of feature selection on Naive Bayes model. They [22] developed 6 Naive Bayes-based models in which each model involves a single feature selection technique chosen from individual FS, forward FS, Backward FS, Plus-I takeaway-r FS, AR1, and All. The experiment resulted in the Naive Bayes model with Plus-I takeaway-r feature selection having the best performance with an accuracy of 93.39% while the Naive Bayes classifier with individual feature selection technique has the least performance with an accuracy of 92.05% thereby leading to the conclusion that feature selection has a direct impact on the accuracy of phishing detection.

Jaya T et al. [16] explored the prospect of using unsupervised learning to cluster spam and ham messages in mail using frequency weight-age of words in the message content in more of a natural language processing task and comparing the performances of each of Random Forest, Logistic, Random Tree, Bayes Net, and Naive Bayes algorithms with LSTM Algorithms for phishing detection. The experiment resulted in LSTM which is deep learning based having an encouraging performance, followed by random forest.

One limitation that is peculiar to each and every previously proposed models, frameworks, and approaches is that they can only detect text-based and URL-based phishing webpages and URLs as they are only trained based on text and properties of the Uniform Resource Locator. Current machine learning and deep learning models are not trained to detect more complex and increasingly sophisticated phishing attack which relies heavily on SEO friendly URL, putting text-on images, and Deep-fake AI generated video to evade detection. Hence, there vulnerabilities to complex form of phishing attack.

III. EXPERIMENTAL SETUP

A. Dataset

The complexity of the research means we cannot rely on a single data. So, we use two publicly available datasets. There is no publicly available dataset for video-based, audio-based, and image-based phishing dataset, so we use simulation to internally generate them.

We use the "B" version of Mendeley phishing dataset which was designed as a benchmark dataset for training a machine

TABLE I
NEGLIGIBLE IMPACT OF MAX_DEPTH AND RANDOM_STATE ON
ACCURACY.

MAX_DEPTH	RANDOM_STATE	ACCURACY
30	0	90.1
20	0	90.0
10	0	88.1
40	1	88.2
50	1	89.8
60	1	88.1

learning models for phishing detection. It includes 11430 URLs and 87 extracted features from which models could be trained. Features in the dataset are classified further into three (3) different categories in which 56 extracted features are from the structure and syntax URL, 24 features were extracted from the content of the URL correspondent pages, and the remaining 7 features which are features with the greatest impact on prediction outcome are extracted from external services. 50% of the dataset used are "phishing" while the remaining 50% are from "legitimate" URLs. This balance of the dataset ensures that the prediction result is not unfairly tilted toward or against a particular category.

The second public dataset that we used was the Spam Message Classification dataset from KAGGLE containing 5157 unique records. The remaining datasets in the form of deep-fake videos and images were simulated and internally generated due to unavailability of such datasets in the public repository.

B. Settings

The proposed framework has 4 predictive layers, with each layer suitable for a specialized category of dataset to ensure adaptability. We show the results in several settings.

- **Layer 1 (URL-Based Training):** We did traditional machine learning training on the first layer using the Mendeley phishing dataset. Out of possible 87 features, we use chi2 from sklearn feature selection library to select the best 19 features, having set the hyper-parameter k-value to 7 for optimal result which gave us a combination of the best 19 features. The dataset was split into two such that 80% was used for training, while the remaining 20% was used for validation tests. We choose random forest because of its suitability for URL-based phishing detection relative to other classifiers [4], [29], [31], [25], [28], [18]. During iteration, we set both the depth and random variable to several values for optimal result but only observed a small but negligible change in the variation of the accuracy until 39. with depth ≥ 39 , the accuracy remains constant, at least till when we increase the randomness of the tree to 1 before observing little change. We finally settled on setting the randomness state to 0 so that each tree remains the same each time it is generated.

- **Layer 2 (Image Processing):** This is the layer where the Hypertext Markup Language (HTML) of the actual phishing webpage is secretly web-scraped behind the scenes without any actual navigation for security purposes. The behind-the-scenes mode of web scrapping the HTML content protects the server and the network from potential drive-by attacks that might originate from the phishing site. All the syntax of HTML mark-up language was removed From the extracted HTML by REGEX as we needed only the content within the opening and closing of the body tag which is the section being served by web server to potential victims while on a phishing website, this step securely brings whatever content (textual, videos, or images) that will be served to potential victim into the framework for series of processing, and this effectively ensure that they cannot evade detection.

Next, we wrote an algorithm to iterate through every filtered word in the sentence, returning only the list of words with any image extension. The fact that the webpage was webs crapped means our program automatically returns a list of the full path of those images from the web server where they are hosted. The returned list is further iterated and passed through an Optical Character Recognition (OCR) library which uses computer vision to read content of each images into a raw text message and forward it to the next layer for further processing.

Algorithm 1 LSTM Model Training for Natural Language Processing (NLP) task

```

train_LSTM( $f_i, w_i, o_j$ )
for epochs = 1 to  $N$  do
  while ( $j \leq m$ ) do
    Randomly initialize  $w_i = \{w_1, w_2, \dots, w_n\}$ 
    input  $o_j = \{o_1, o_2, \dots, o_m\}$  in the input layer
    forward propagate ( $f_i \cdot w_i$ ) through layers until getting
    the predicted result  $y$ 
    compute  $e = y - y^2$ 
    back propagate  $e$  from right to left through layers
    update  $w_i$ 
  end while
end for

```

- **Layer 3 (Speech Synthesis):** Having successfully web-scraped the hypertext mark-up language of the potential phishing site behind the scenes, and without any actual navigation for security purposes at the previous layer. Returned content from the previous layer 2 is further iterated through with "for" loop. "for" loop iterates through every filtered word in the sentence, returning only the list of words with any video extension, the return list is automatically the full path of those videos from the server where they are hosted.

Next, we did further iteration through each of the returned video files, and on each iteration step, we used a combination of gttts, pydud, and moviepy for conversion from video file to audio file ".wav" format, after which the actual synthesis of each speech across the "loop" began with natural language processing speech recognition. The final operation output at this layer is a raw text file obtained from synthesizing the

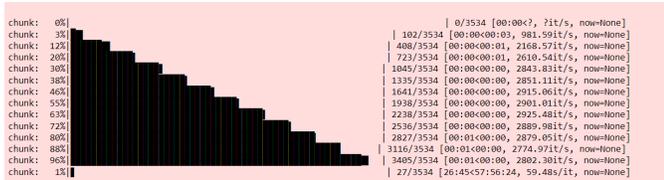


Fig. 1. Step-wise speech synthesis of each audio file during execution of "for" loop in layer 3 to produce text which was later passed on to layer 4. Texts from the phishing sites were processed at Layer 1, images were processed at Layer 2, while Layer 3 processed videos. All text was finally outputted to layer 4 for final prediction using a variant of Recurrent Neural Network in Long Short-Term Memory.

speech. At this stage, we have the images read to text from the previous layer 2 and speeches in the video synthesized to text, next, we combined each of the text from layer 1, layer 2, to layer 3 forwarding them to layer 4 for final prediction.

• **Layer 4 (Speech Synthesis):** We choose LSTM network because of the effective solution it offers to vanishing and exploding gradient which are Long term dependency problem in Recurrent Neural Network, the cell state in LSTM network serves as a memory to the network thereby given it the ability to remember the past. At layer 4, we have all outputted and processed text contents from each of the previous layers, and there is need to capture every long-term dependencies, short term dependencies, and sequences which could be provided by the cell state in LSTM network to ensure a more accurate prediction.

We built an LSTM deep learning-based model, in which 80% of 5572 samples were used as training samples while the remaining 20% was used for validation. The dataset has a maximum of 10,000 features from the word sample, out of which we have 9004 unique words from the dataset. During training and validation, we had wide validation loss leading to low prediction ability but continued to adjust the number of layer, features, epoch, batch size, and activation. We obtained the best result at the following parameter;

- dense layer = 1
- activation layer= sigmoid
- epochs = 10
- batch size=60

feature size = 32 The built LSTM network model is at the 4th layer of the framework where all processed output from each of the previous layers are merged and passed to the newly built LSTM model to make the final prediction.

IV. FRAMEWORK ADAPTABILITY AND PERFORMANCE EVALUATION

To stretch and validate our multilayered adaptive framework for its effectiveness in the detection of phishing sites containing any of (Text, videos, and images), or a combination of any, or all of the 3. It is worth remembering that all existing AI or machine learning-based phishing detection techniques and frameworks can only detect text-based [7], [17], [13], [31] or URL-based [29], [7], [34], [2] phishing sites leading to their



Fig. 2. LSTM network resulting in 0.98 accuracy at optimal parameter



Fig. 3. LSTM network resulting in 0.08 loss at optimal parameter.

- vulnerabilities to;
- phishing sites with friendly URL
- phishing site on hacked legitimate domain name server (DNS)
- Image-only phishing site
- video-only phishing site

or, combination of any of them in any order. To validate both the effectiveness and adaptability of our proposed framework in overcoming such limitations, we created 4 categories of phishing sites and uploaded them to a secure server with a compromised DNS on a friendly URL; the first was a text-only phishing site, image-only phishing site, video-only phishing site, and a phishing site combining all the features. We use friendly URLs, and hacked legitimate DNS on the text-only phishing site, so that they can evade detection at layer 1 until the 4th layer where it will be detected, while we created phishing sites containing each image-only, text-only, and a combination of both to test the adaptability of the framework to different scenarios of a phishing site.

In each scenario, we have 100% accuracy as the framework

successfully adapts to each scenario and detects accordingly, thereby overcoming limitations associated with current approaches to phishing detection methods.

V. CONCLUSION

In this research, we proposed a multi-layer adaptive framework that uses the computer vision capability of Optical Character Recognition (OCR) to read images on live phishing sites to text, and synthesize speech from uploaded deep-fake videos, while using Random Forest, and LSTM network, along with web scrapped text at various predictions layered of the framework to significantly improve the detection rate and performance of AI-based models for phishing detection. Considering the fact that existing AI-based phishing detection techniques, frameworks, and approaches can only detect text-based [7], [17], [13], [31] or URL-based phishing [29], [7], [34], [2] sites which leads to their vulnerability and inability to detect image-based, or video-based phishing sites, the proposed framework is able to overcome limitations in existing approaches, significantly improve phishing attack detection, and successfully detect complex phishing webpages with multi-dimensional deep-fake videos, images, and texts.

VI. LIMITATION AND FUTURE RESEARCH DIRECTION

We used Mendeley and Kaggle phishing datasets which are URL-based and Text-based respectively. image-based and video-based phishing datasets are not publicly available because they are newly adopted forms of phishing websites to evade detection, we simulated them to get the data internally generated for this research especially with regard to deep-fake videos, hence getting publicly available image-based or Video-based phishing datasets will significantly help the research community in this.

The other research direction we will point at is the computational aspect during training. The proposed framework uses Random Forest and LSTM network at Layer 1 and Layer 4 respectively. The fact that Random Forest algorithm creates multiple trees each time to combine individual tree decisions for more accurate prediction leads to an increment in computation time, we have to set the random state to zero while changing the maximum depth for optimal hyper-parameter. Apart from the training computation time, there is also the server response time as the framework web scrapped the content behind the scenes thereby protecting the server against potential drive-by attacks. Hence, reducing the server response and computational to fraction of a second is an area open to future research in this domain.

It is also worth noting that our artifacts which consist of source code, dataset, images, videos, and audio files for this research had been uploaded to a public GitHub repository for reproducibility of our research. Artifact can be found on GitHub at;

Deep Learning-Based Speech and Vision Synthesis to Improve Phishing Attack Detection through a Multi-layer

Adaptive Framework and also at Code Ocean computational research platform, with the exception of the internally generated deep-fake video and audio data files for privacy.

REFERENCES

- [1] Lozan Mohammed Abdulrahman, Sarkar Hasan Ahmed, Zryan Najat Rashid, Yousif Sufyan Jghef, Teba Mohammed Ghazi, and Umed H Jader. Web phishing detection using web crawling, cloud infrastructure and deep learning framework. *Journal of Applied Science and Technology Trends*, 4(01):54–71, 2023.
- [2] Moruf Akin Adebowale, Khin T Lwin, and Mohammed Alamgir Hossain. Intelligent phishing detection scheme using deep learning algorithms. *Journal of Enterprise Information Management*, 36(3):747–766, 2023.
- [3] Sikiru Adewale, Tosin Ige, and Bolanle Hafiz Matti. Encoder-decoder based long short-term memory (lstm) model for video captioning. *arXiv preprint arXiv:2401.02052*, 2023.
- [4] R Alazaidah, A Al-Shaikh, MR AL-Mousa, H Khafajah, G Samara, M Alzyoud, N Al-Shanableh, and S Almatameh. Website phishing detection using machine learning techniques. *Journal of Statistics Applications & Probability*, 13(1):119–129, 2024.
- [5] Ali Aljofey, Qingshan Jiang, Abdur Rasool, Hui Chen, Wenyin Liu, Qiang Qu, and Yang Wang. An effective detection approach for phishing websites using url and html features. *Scientific Reports*, 12(1):8842, 2022.
- [6] Zainab Alkhalil, Chaminda Hewage, Liqaa Nawaf, and Imtiaz Khan. Phishing attacks: A recent comprehensive study and a new anatomy. *Frontiers in Computer Science*, 3:563060, 2021.
- [7] Faisal S Alsubaie, Abdulwahab Ali Almazroi, and Nasir Ayub. Enhancing phishing detection: A novel hybrid deep learning framework for cybercrime forensics. *IEEE Access*, 2024.
- [8] J Anitha and M Kalaiarasu. A new hybrid deep learning-based phishing detection system using mcs-dnn classifier. *Neural Computing and Applications*, pages 1–16, 2022.
- [9] Z Ghaleb Al-Mekhlafi, B Abdulkarem Mohammed, Mohammed Al-Sareem, Faisal Saeed, Tawfik Al-Hadhrani, Mohammad T Alshammari, Abdulrahman Alreshidi, and T Sarheed Alshammari. Phishing websites detection by using optimized stacking ensemble model. *Computer Systems Science and Engineering*, 41(1):109–125, 2022.
- [10] Tosin Ige and Christopher Kiekintveld. Performance comparison and implementation of bayesian variants for network intrusion detection. *arXiv preprint arXiv:2308.11834*, 2023.
- [11] Tosin Ige, William Marfo, Justin Tonkinson, Sikiru Adewale, and Bolanle Hafiz Matti. Adversarial sampling for fairness testing in deep neural network. *arXiv preprint arXiv:2303.02874*, 2023.
- [12] R Ishwarya, S Muthumani, Siva Sharma Karthick PG, and S Suriya. Separation of phishing emails using probabilistic classifiers. In *2023 9th International Conference on Advanced Computing and Communication Systems (ICACCS)*, volume 1, pages 1676–1679. IEEE, 2023.
- [13] Ankit Kumar Jain and BB Gupta. A survey of phishing attack techniques, defence mechanisms and open research challenges. *Enterprise Information Systems*, 16(4):527–565, 2022.
- [14] Ankit Kumar Jain and Brij B Gupta. A novel approach to protect against phishing attacks at client side using auto-updated white-list. *EURASIP Journal on Information Security*, 2016:1–11, 2016.
- [15] Ankit Kumar Jain and Brij B Gupta. A machine learning based approach for phishing detection using hyperlinks information. *Journal of Ambient Intelligence and Humanized Computing*, 10:2015–2028, 2019.
- [16] T Jaya, R Kanyaharini, and Bandi Navaneesh. Appropriate detection of ham and spam emails using machine learning algorithm. In *2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI)*, pages 1–5. IEEE, 2023.
- [17] R Jayaraj, A Pushpalatha, K Sangeetha, T Kamaleshwar, S Udhaya Shree, and Deepa Damodaran. Intrusion detection based on phishing detection with machine learning. *Measurement: Sensors*, 31:101003, 2024.
- [18] Kunj Joshi, Chintan Bhatt, Kaushal Shah, Dwireph Parmar, Juan M Corchado, Alessandro Bruno, and Pier Luigi Mazzeo. Machine-learning techniques for predicting phishing attacks in blockchain networks: A comparative study. *Algorithms*, 16(8):366, 2023.

- [19] Abdul Karim, Mobeen Shahroz, Khabib Mustofa, Samir Brahim Belhaouari, and S Ramana Kumar Joga. Phishing detection system through hybrid machine learning based on url. *IEEE Access*, 11:36805–36822, 2023.
- [20] Ann Zeki Ablahd Magdacy Jerjes, Adnan Yousif Dawod, and Mohammed Fakhrulddin Abdulqader. Detect malicious web pages using naive bayesian algorithm to detect cyber threats. *Wireless Personal Communications*, pages 1–13, 2023.
- [21] Amgad Muneer, Rao Faizan Ali, Abdo Ali Al-Sharai, and Suliman Mohamed Fati. A survey on phishing emails detection techniques. In *2021 International Conference on Innovative Computing (ICIC)*, pages 1–6. IEEE, 2021.
- [22] Twana Mustafa and Murat Karabatak. Feature selection for phishing website by using naive bayes classifier. In *2023 11th International Symposium on Digital Forensics and Security (ISDFS)*, pages 1–4. IEEE, 2023.
- [23] U Nishitha, Revanth Kandimalla, Reddy M Mourya Vardhan, and U Kumaran. Phishing detection using machine learning techniques. In *2023 3rd Asian Conference on Innovation in Technology (ASIANCON)*, pages 1–6. IEEE, 2023.
- [24] Amos Okomayin, Tosin Ige, and Abosede Kolade. Data mining in the context of legality, privacy, and ethics. 2023.
- [25] Dele Olukoya. Heterogeneous ensemble feature selection and multilevel ensemble approach to machine learning phishing attack detection.
- [26] Kamal Omari. Comparative study of machine learning algorithms for phishing website detection. *International Journal of Advanced Computer Science and Applications*, 14(9), 2023.
- [27] Orvila Sarker, Asangi Jayatilaka, Sherif Haggag, Chelsea Liu, and M Ali Babar. A multi-vocal literature review on challenges and critical success factors of phishing education, training and awareness. *Journal of Systems and Software*, 208:111899, 2024.
- [28] Muhammad Waqas Shaikat, Rashid Amin, Muhana Magboul Ali Muslam, Asma Hassan Alshehri, and Jiang Xie. A hybrid approach for alluring ads phishing attack detection using machine learning. *Sensors*, 23(19):8070, 2023.
- [29] Sanjeev Shukla, Manoj Misra, and Gaurav Varshney. Http header based phishing attack detection using machine learning. *Transactions on Emerging Telecommunications Technologies*, page e4872, 2024.
- [30] Muhammed Kürşad Uçar, Majid Nour, Hatem Sindi, Kemal Polat, et al. The effect of training and testing process on machine learning in biomedical datasets. *Mathematical Problems in Engineering*, 2020, 2020.
- [31] RJ van Geest, G Cascavilla, J Hulstijn, and N Zannone. The applicability of a hybrid framework for automated phishing detection. *Computers & Security*, page 103736, 2024.
- [32] Yifei Wang. A survey of phishing detection: from an intelligent countermeasures view. In *2022 IEEE Conference on Telecommunications, Optics and Computer Science (TOCS)*, pages 761–769. IEEE, 2022.
- [33] Palla Yaswanth and V Nagaraju. Prediction of phishing sites in network using naive bayes compared over random forest with improved accuracy. In *2023 Eighth International Conference on Science Technology Engineering and Mathematics (ICONSTEM)*, pages 1–5. IEEE, 2023.
- [34] Erzhou Zhu, Kang Cheng, Zhizheng Zhang, and Huabin Wang. Pdhf: Effective phishing detection model combining optimal artificial and automatic deep features. *Computers & Security*, 136:103561, 2024.
- [35] Rasha Zieni, Luisa Massari, and Maria Carla Calzarossa. Phishing or not phishing? a survey on the detection of phishing websites. *IEEE Access*, 11:18499–18519, 2023.