## DETECTING PHISHING WEBSITES USING DEEP LEARNING TECHNOLOGIES

**Ailneni Banu Harshini**, Department of Computer Science, St.Francis College for Women
Hyderabad – 16. banurao111@gmail.com
**Dr. Sr. Sujatha Yeruva**, Department of Computer Science, St.Francis College for Women
Hyderabad – 16. sr.sujatha@sfc.ac.in

Abstract-The rapidly growing phishing sites are one of the most alarming cybersecurity threats where users are cheated to provide their login credentials, financial information, and personal data. Traditional rule-based and machine learning methods find it difficult to cope with the ever-changing phishing tactics of the cybercriminals. To this end, the present work aims to propose a deep learning-based phishing detection model using a Fully Connected Neural Network (FCNN). The model trains on a set of datasets having various URL-based features for categorizing the classification of websites into legitimate and phishing. The results of the experimental results indicate the proposed model classifies with accuracy of 85% in training set and with 84% in test sets, hence validating its practical suitability.

To achieve better accessibility, the model is integrated into an easy-to-use Graphical User Interface (GUI) using Tkinter, through which r URL analysis and phishing detection take place. It is an efficient and scalable method for detecting phishing attempts, and thus the chance of cyber fraud is reduced to a great extent. This paper shows the effective application of deep learning in the field of cybersecurity and opens ways for future developments in phishing prevention strategies.

*Keywords: Phishing Detection, Deep Learning, Fully Connected Neural Network, Website Security, FCNN, Tkinter GUI*

## I.     INTRODUCTION:

As quickly the internet has expanded, it has made communication smooth, e-commerce efficient, and digital transactions easy. However, this has been accompanied by a grave rise in cyber dangers where phishing attacks are surfacing as one of the most common and crafty ways to commit online fraudulence. Phishing sites pretend to be the real ones that look fully authentic and trick users into providing their login credentials, financial information, and private personal details. According to reports on cybersecurity, millions of phishing websites are developed every year. Their traditional rule-based detection systems become ineffective as they don't have the ability to change based on attack strategies.To counter all these, machine learning (ML) and deep learning (DL) techniques have been highly explored in the domain of phishing detection. ML-based approaches utilize handcrafted features, classifiers such as decision trees, support vector machines, and random forests. However, the problem of feature selection often arises in this approach and it fails to capture the deeper, non-linear relationships within the data. In contrast, deep learning models, especially Fully Connected Neural Networks (FCNNs), have recently demonstrated promising performance in pattern recognition tasks by automatically extracting and learning relevant features from raw data.

This paper shows  a deep learning-based phishing website detection system based on an FCNN model trained on a dataset containing URL-based attributes. The proposed model has achieved 85% training accuracy and 84% testing accuracy, which means that the model can effectively classify phishing and legitimate websites. Furthermore, a Tkinter-based graphical user interface (GUI) has been developed to allow users to input domain names and receive classification results. It also used data preprocessing techniques like feature scaling and selection to improve the model's performance.

## II.     DATASET:

The data set used in this paper is the most important in training and evaluation of this phishing website detection model. Various features extracted from website URLs, indicating whether a given URL points to safe or malicious sources, are available in this data set. Each entry is labeled as "Phishing" (1) or "Legitimate" (0), so the deep learning model learns to distinguish between a legitimate website and fraud.

**Dataset Summary**
- File Name: urldata.csv
- No. of Rows: 10,000 rows
- No. of Columns: 18 columns
- TargetVariable: Label Binary classification: 0 for Safe, 1 for Phishing
- Key Attributes:Various URL-based and domain-based features, which are extracted and analyzed to detect phishing attempts.

## III.     LITERATURE SURVEY:

Security online is facing a persistent threat from phishing attacks with financial loss and identity theft. Researchers in previous work focused on different machine learning-based approaches to detect phishing websites effectively. This literature survey reviews recent works that utilize machine learning methods in the detection of phishing websites. Mohammed Hazim Alkawaz, Stephanie Joanne Steven, and Asif Iqbal Hajamydeen proposed a system based on capturing blacklisted URLs to alert users via pop-ups and emails. Future improvements may include text message notifications for added security[1].A. Mishra and B. B. Gupta presented a hybrid approach to zero-day phishing detection by using similarity analysis of URIs and images with a significantly lower false positive rate than compared to the approaches that exist already. They also define further scope in detecting previously unseen phishing sites[2]. R. Kiruthiga and D. Akila reviewed several approaches for phishing site detection using various machine learning methods, including algorithms such as Naïve Bayes, SVM, Decision Tree, and Random Forest. They recommend updating features with the advancement of phishing techniques for better detection[3].Ankit Kumar Jain and B. B. Gupta suggested a client-side, language-agnostic phishing detection mechanism using hyperlink information that resulted in more than 98.4% accuracy by logistic regression. Their method doesn't rely on third-party dependencies for security purposes[4].Rishikesh Mahajan (2018) attained 97.14% accuracy for phishing website detection with the help of Random Forest algorithm with minimal false positives. The future enhancements are to make use of the hybrid method by integrating both machine learning and blacklist methods to increase the accuracy level[5].Purvi Pujara and M. B. Chaudhari (2018) has also reviewed the methods of detection for phishing websites. They concluded that tree-based classifiers in machine learning are the best for detection[6].David G. In 2016, Dobolyi and Ahmed Abbasi launched an open archive of real phishing websites called PhishMonger, designed to aid in research into mechanisms of phishing as well as the susceptibility of users. It supports the development of anti-phishing tools, usability studies, and phishing trend analysis[7]. Satish.S and Suresh Babu.K (2013) proposed an anti-phishing technique, which applied URL domain identity and scripting mechanisms using a classification algorithm for finding an approximate classification while improving detection speed. Their approach outperforms the existing tools by reducing latency in identifying phishing URLs[8]. Ping Yi proposed a deep learning framework using Deep Belief Networks (DBN) for phishing detection based on original and interaction features, in 2018. The model obtained nearly 90% True Positive Rate (TPR) in an extensive dataset[9]. Marwa Al-Saedi and Nahla Abbas Flayh surveyed techniques of machine learning for phishing such as hybrid methods like Random Forest that achieved accuracy greater than 99%. They propose deep learning methods may be a rich area of investigation to mitigate continuously changing phishing threats[10]. Almaha Abuzuraiq and Mouhammd Alkasassbeh surveyed phishing detection techniques that ranged from content-based, heuristic-based, and fuzzy rule-based methodologies. They conclude that no single approach is perfect, and continuous improvements are needed to counter evolving phishing techniques[11].Y. Sönmez et al. (2018) proposed a phishing detection model using feature extraction and classification techniques, including SVM, NB, and Extreme Learning Machine (ELM). ELM, with six activation functions, got the highest accuracy[12].M. A. El-Rashidy (2021) provided a phishing detection model based on a novel feature selection method about URL, HTML, JavaScript, text, images, and domain names. The method attained accuracy of 96.66% which supersedes other algorithms while the risk to download malware could be reduced[13]. G. Harinahalli Lokesh and G. BoreGowda (2020) discussed phishing attacks on Nigerian internet banking users; trends and countermeasures were identified. They proposed an

adaptive anti-phishing model to assist financial institutions and policymakers in combating the evolving threats[14]. Ume Zara et al. proposed a phishing detection model based on deep learning and ensemble learning, which achieves 99% accuracy using top features selected by information gain, gain ratio, and PCA. Their approach enhances adaptability to evolving phishing techniques[15].Gopal, S. (n.d.) was trained the Machine learning models and deep neural nets to predict the phishing website[16]. T. Peng, I. Harris, and Y. Sawa (2018) designed the phishing detection model based on Natural Language processing method for identifying text malicious content. Their approach, which is validated on a large dataset, focuses on semantic analysis for phishing attacks identification[17]. Sahingoz, BUBEr, and Kugu proposed DEPHIDES, a deep learning-based phishing detection system that uses RNN, BiRNN, CNN, ANN, and attention-based networks. Their findings indicated that deep learning outperformed the traditional methods, and all codes and datasets were shared for transparency[18].Wenyin, Huang, Xiaoyue, Min, and Deng proposed a phishing detection approach based on visual similarity by analyzing block, layout, and style similarities. Their method demonstrated promising results on 328 webpages with plans for further improvements and commercial applications[19].Zhang, Hong, and Cranor developed CANTINA, a content-based phishing detection approach using the TF-IDF algorithm. Their method achieves an accuracy of around 95% with further heuristics to reduce false positives[20].Mohammad, Thabtah, and McCluskey proposed a phishing detection model with structuring neural networks based on the ever-changing web threat. Their model achieved high accuracy, fault tolerance, and high generalization performance across experiments[21].Buber, Diri, and Sahingoz proposed an NLP-based phishing detection system through machine learning along with visual similarity analysis. Their experiments showed that the Random Forest algorithm was able to find phishing attacks with a 97.2% success rate[22]. Abu-Nimeh, Nappa, Wang, and Nair compared the effectiveness of phishing detection mechanisms using machine learning techniques on a dataset of 2,889 emails and 43 features. Their research analyzed LR, CART, BART, SVM, RF, and Neural Networks for predictive accuracy[23].Babagoli, Aghababa, and Solouk proposed a phishing detection approach using a metaheuristic nonlinear regression algorithm with feature selection. The authors' methodology achieved 96.32% accuracy, and their proposed methodology was superior to SVM in classification[24].Butnaru, Mylonas, and Pitropakis used a lightweight phishing detection system based on supervised machine learning with URL-based features. Their model outperformed Google Safe Browsing and remained effective against phishing URLs even a year after training[25].

## IV. METHODOLOGY:

Deep learning's methodology for spotting phishing websites consists of numerous stages, including dataset preprocessing, feature selection, model creation, training, evaluation, and GUI implementation.

**1.Data collection and preprocessing**.

The dataset used in this study consists of 10,000 cases with 18 features collected from website URLs. These elements include domain-related attributes that can help differentiate between phishing and authentic websites. Preprocessing of the dataset includes:

● Remove unnecessary columns, such as domain names.
● Handle any missing values.
● Normalizing numerical features with StandardScaler to achieve uniformity in input data.
● To evaluate model performance, split the dataset between 80% training and 20% testing sets using train_test_split.

**2. Feature Engineering**

Feature selection is an important step, which improves the accuracy of classification. The extracted features are many, based on URL length, existence of HTTPS, a dot count in the domain, and so on. All these attributes are scaled to achieve optimized learning in the model to avoid bias.

**3. Deep Learning Model Development**

A Fully Connected Neural Network (FCNN) is used to determine whether a website is phishing or safe. The model architecture includes an input layer of 128 neurons with ReLU activation.

The model has two hidden layers (64 and 32 neurons) with ReLU activation and dropout layers (30%) to reduce overfitting.

A single neuron output layer with sigmoid activation is used for binary classification.

The Adam optimizer optimizes weight updates efficiently and the binary cross-entropy loss function is used for training.
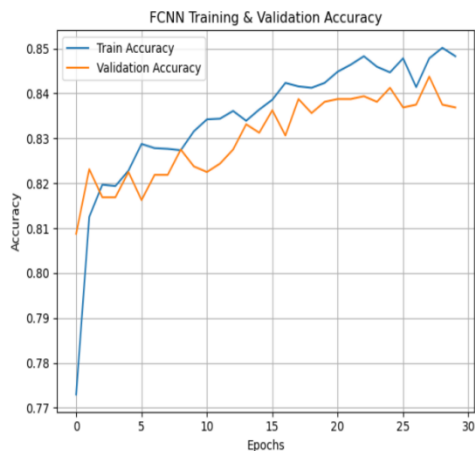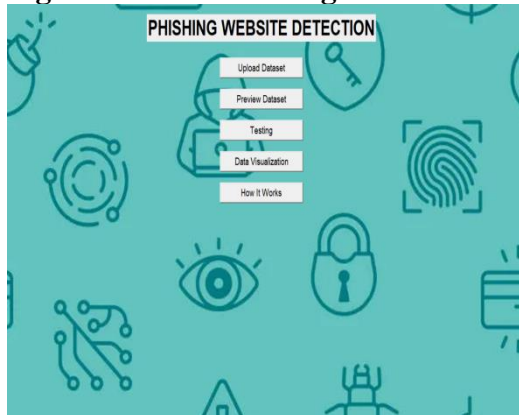


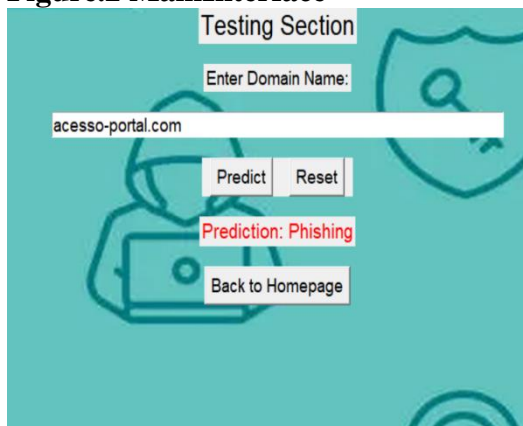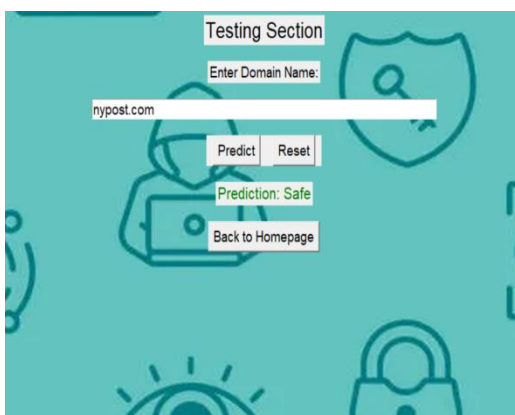**Figure.1 FCNN training and Validation** ...



**Figure.2 MainInterface**

### 4. Model Training and Evaluation

The training dataset is used to train the model across 30 epochs with a batch size of 32. Accuracy measures are used to assess the performance on both training and testing data. The final model's training accuracy is 85%.The testing accuracy was 84%.These findings show that the FCNN model successfully identifies phishing websites and generalizes well.

#### GUI Implementation for User Interaction

To make the detection system user-friendly, a Tkinter-based GUI is developed. This allows users to:

Upload and visualize datasets.
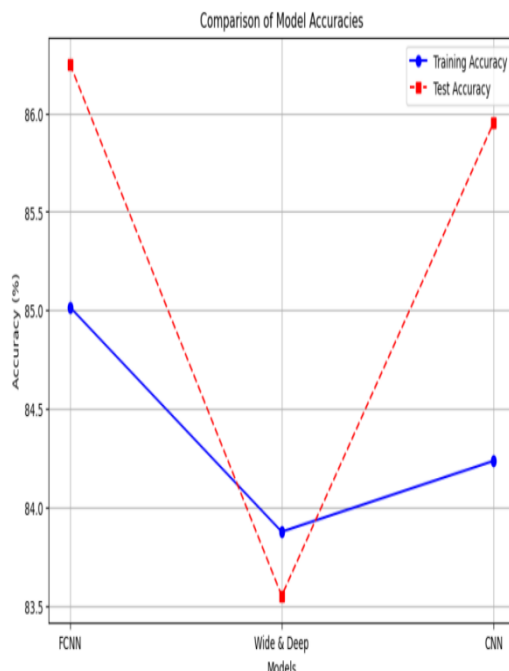
Input a domain name for phishing detection.

Observe data distribution and feature importance through visualization dashboards.

The GUI makes the system more accessible and user-friendly for non-technical users so that it is practically applicable in cybersecurity domains.



**Figure.3 Testing the Domain name**



**Figure.4  Testing the domain name**

### V.EMPRICAL RESULTS AND GRAPHICAL REPRESENTATION

Figure.1 Training curve of fully connected neural network (FCNN) on 30 epochs: training accuracy in blue, validation accuracy in orange line. Training starts with a less higher level and increases constantly. This would represent that the model is indeed learning from the training data. Validation accuracy follows an upward trend in general, which is quite fluctuated because the model generalizes well on other epochs. By the 30th epoch, both curves converge to a high accuracy level of about 85%, suggesting that the FCNN achieves good performance on both training and validation sets, with only a slight gap between them that may be attributable to mild overfitting.Figure.2 shows the main user

interface for the phishing website detection system. It has a central title, "Phishing Website Detection," with a sequence of buttons: Upload Dataset, Preview Dataset, Testing, Data Visualization, and How It Works to take the user through different stages of the system. The background contains security-related icons like locks, fingerprints, and keys to emphasize the theme of cybersecurity and phishing prevention. Figure.3&4 shows the testing interface of the phishing website detection system. In this interface, a user can type in a domain name in the text field and click Predict to classify it as either phishing or legitimate. This interface also contains options to reset the input and return to the Homepage.

Figure.5 shows the comparison of model accuracies for all three deep learning models, FCNN, Wide & Deep, and CNN applied on the project "Safe Surfing: Detecting Phishing Websites Using Deep Learning Technologies". Training Accuracy (Blue Solid Line): The training accuracy by the FCNN model is about 85%, with a small margin above that of the CNN whose training accuracy is at about 84.5%. The Wide & Deep model, however, has the lowest training accuracy at just above 83.5%.Test Accuracy (Red Dashed Line): The test accuracy for FCNN is the highest at about 86.2%, which indicates good generalization. In contrast, the Wide & Deep model has the lowest test accuracy, at about 83.6%, which might indicate underfitting. Test accuracy, in this case, which is over training accuracy by about 86%, the CNN model is generalizing very well. Overall, Fig.5 shows that although the test accuracy of the FCNN model is the highest, the CNN model is also delivering well without much overfitting. Wide & Deep is complex with both low train and low test accuracy that indicates this model may not be the most suitable for this task.



**Fig.5 Comparison of model Accuracies**

## VI. CONCLUSION

In this paper, we analyze the detection of phishing websites with deep learning models. We compared FCNN with Wide & Deep and CNN to distinguish between a phishing website and a legitimate website. The FCNN model gives the highest accuracy .The research reflects the usefulness of deep learning in cybersecurity by providing a solid solution for phishing detection.We have made an interactive interface with the user requirement to feed in website information, and once that is provided, instant prediction is returned whether the site is phishing or actual. This really helps users immensely in avoiding such online threats in order to allow safer browsing experience.

Future work includes optimization of model performance, the incorporation of other features, and integration of real-time URL analysis to achieve a more accurate detection.

## VII. REFERENCES

[1].Alkawaz, M. H., Steven, S. J., & Hajamydeen, A. I. (2020, February 28–29). Detecting phishing website using machine learning. 16th IEEE International Colloquium on Signal Processing & Its Applications (CSPA 2020).

[2].Mishra, A., & Gupta, B. B. (2014). Hybrid solution to detect and filter zero-day phishing attacks. ERCICA.

[3].Kiruthiga, R., & Akila, D. (2019). Phishing websites detection using machine learning. International Journal of Recent Technology and Engineering (IJRTE), 8(2S11).

[4].Jain, A. K., & Gupta, B. B. (2018). A machine learning-based approach for phishing detection using hyperlinks information. Springer Nature.

[5].Mahajan, R. (2018). Phishing website detection using machine learning algorithms.

[6].Pujara, P., & Chaudhari, M. B. (2018). Phishing website detection using machine learning: A review.

[7].Dobolyi, D. G., & Abbasi, A. (2016). PhishMonger: A free and open-source public archive of real-world phishing websites.

[8].Satish, S., & Babu, S. K. (2013). Phishing websites detection based on web source code and URL in the webpage.

[9].Yi, P. (2018). Web phishing detection using a deep learning framework.

[10].Al Saedi, M., & Flayh, N. A. Phishing website detection using machine learning: A review.

[11].Abuzuraiq, A., & Alkasassbeh, M. Review: Phishing detection approaches.

[12].Sönmez, Y., Tuncer, T., Gökal, H., & Avci, E. (2018). Phishing websites features classification based on extreme learning machine. Proceedings of the 6th International Symposium on Digital Forensic & Security (ISDFS 2018).

[13].El-Rashidy, M. A. (2021). A smart model for web phishing detection based on new proposed feature selection technique. Menoufia Journal of Electronic Engineering Research, 30(1), 97–104.

[14].Lokesh, G. H., & BoreGowda, G. (2020). Phishing website detection based on an effective machine learning approach. Journal of Cyber Security Technology.

[15].Zara, U., Ayyub, K., Khan, H. U., Daud, A., Alsahfi, T., & Ahmad, S. G. Phishing website detection using deep learning models.

[16].Gopal, S. Phishing website detection by machine learning techniques – Dataset [CSV file]. GitHub.

[17].Peng, T., Harris, I., & Sawa, Y. (2018). Detecting phishing attacks using natural language processing and machine learning. Proceedings of the IEEE 12th International Conference on Semantic Computing (ICSC), 300–301.

[18].Sahingoz, O. K., Buber, E., & Kugu, E. (2024). DEPHIDES: Deep learning-based phishing detection system. IEEE Access, 12, 8052–8070.

[19].Wenyin, L., Huang, G., Xiaoyue, L., Min, Z., & Deng, X. (2005, May). Detection of phishing webpages based on visual similarity. Proceedings of the 14th International Conference on World Wide Web.

[20].Zhang, Y., Hong, I., & Cranor, F. (2007, May 8–12). CANTINA: A content-based approach to detecting phishing websites. Proceedings of the 16th International Conference on World Wide Web.

[21].Mohammad, R. M., Thabtah, F., & McCluskey, L. (2014, August). Predicting phishing websites based on self-structuring neural networks. Neural Computing & Applications, 25(2), 443–458.

[22].Buber, E., Diri, B., & Sahingoz, O. K. (2018). NLP-based phishing attack detection from URLs. In Intelligent Systems Design and Applications (Vol. 736, pp. 608–618). Springer.

[23].Abu-Nimeh, S., Nappa, D., Wang, X., & Nair, S. (2007). A comparison of machine learning techniques for phishing detection. Proceedings of the Anti-Phishing Working Groups 2nd Annual eCrime Researchers Summit (eCrime 2007), 60–69.

[24].Babagoli, L., Aghababa, M. P., & Solouk, V. (2019). Heuristic nonlinear regression strategy for detecting phishing websites. Soft Computing, 23(12), 4315–4327.

[25].Butnaru, A., Mylonas, A., & Pitropakis, N. (2021, June). Towards lightweight URL-based phishing detection. Future Internet, 13(6), 154.