# Real-Time Social Media Monitoring for Data Leak Detection

Nitesh Singh[1], Vinit Sharma[2], Satyam Singh[3], Raunit Dubey[4] Ranjeet Singh[5]

[1-4]Student, Department of Information Technology, Buddha Institute of Technology, GIDA, Gorakhpur

[5]Assistant. Professor, Department of Computer Science & Engineering, Buddha Institute of Technology, GIDA, Gorakhpur

**Abstract:** In today's digital environment, social media has become a dominant platform for interaction, enabling rapid two-way communication between individuals and organizations. However, its widespread adoption has also introducedevolving cyber risks, including phishing attacks, data breaches, malicious URL dissemination, identity theft, and financial scams. These threats frequently exploit social engineering techniques, making traditional signature-based detection methods slow and ineffective against emerging attack patterns. To address this challenge, this research presents a real-time social media monitoring system based on a hybrid approach integrating machine learning algorithms, natural language processing, and behavioral analysis. The proposed system continuously analyzes posts, comments, and direct messages to identify suspicious indicators and potential malicious activities. Natural language processing techniques extract semantic and contextual features to detect deceptive language and phishing intent, while behavioral analysis captures abnormal interaction patterns associated with compromised or malicious accounts. In addition, advanced feature extraction methods are employed to model linguistic, structural, and behavioral characteristics, which are used to train supervised learning models for accurate threat classification. Designed for real- time operation, the system enables early detection and rapid response to cyber incidents. Experimental evaluation shows superior accuracy, scalability, and response time. The framework offers an intelligent and adaptive security solution for platforms.

**Keywords:** Social media monitoring, Data leak Detection, Cybersecurity, Threat intelligence,Phishing detection, Fake notification, Malicious attachments, Financial scams, NLP.

## Introduction

Social media platforms have evolved from simple communication tools into large-scale digital ecosystems supporting instant messaging, multimedia sharing, financial transactions, and professional as well as business communication. Platforms such as Facebook, WhatsApp, Instagram,

Twitter, Telegram, and LinkedIn collectively process billions of messages every day, enabling the continuous exchange of sensitive personal, professional, and financial information. While this digital transformation has significantly enhanced global connectivity, it has also introduced critical security and privacy vulnerabilities.

Cybercriminals increasingly exploit the open and trust-based nature of social media environments to carry out deception-driven attacks. These attacks include phishing messages, fake system notifications, malicious file attachments, shortened or obfuscated URLs, fraudulent financial offers, and deceptive requests for sensitive information. Such malicious content is deliberately designed to closely mimic legitimate platform behavior, making it difficult for users to distinguish between genuine and fraudulent communications [1]. As illustrated in Figure 1, a typical phishing lifecycle involves the distribution of deceptive content, redirection to fraudulent websites, credential harvesting, and the subsequent misuse of stolen data. This highlights phishing as a coordinated, multi-stage attack rather than an isolated event.

In recent years, phishing and social engineering attacks have grown rapidly in both volume and sophistication, as evidenced by the increasing number of detected phishing websites shown in Figure 2. Modern attackers employ automated phishing kits, compromised hosting infrastructures, AI-generated content, and advanced URL obfuscation techniques to bypass traditional security defenses. Consequently, conventional protection mechanisms such as keyword-based filtering, static blacklists, and manual inspection are becoming increasingly ineffective against these evolving threats [2].

The consequences of such attacks extend beyond individual users to organizations that rely heavily on social media platforms for communication, customer interaction, and business operations [3]. A single successful malicious interaction can lead to data breaches, unauthorized system access, identity theft, or exposure of confidential information. These risks underline the urgent need for intelligent, automated, and real-time detection systems capable of addressing the scale, diversity, and dynamic nature of modern social media threats [4]. Figure 3 conceptually illustrates this need through a unified social media monitoring framework spanning multiple threat categories and highlighting the complexity of contemporary cyber-attacks.

Motivated by these challenges, this research proposes an AI-driven, real-time social media monitoring framework designed to detect dangerous content that may result in data leakage or security compromise. The proposed system focuses on identifying phishing messages, fake notifications, suspicious URLs, fraudulent financial offers, deceptive information requests, and

malicious attachments using a multi-class artificial intelligence–based detection model evaluated on real-world datasets.
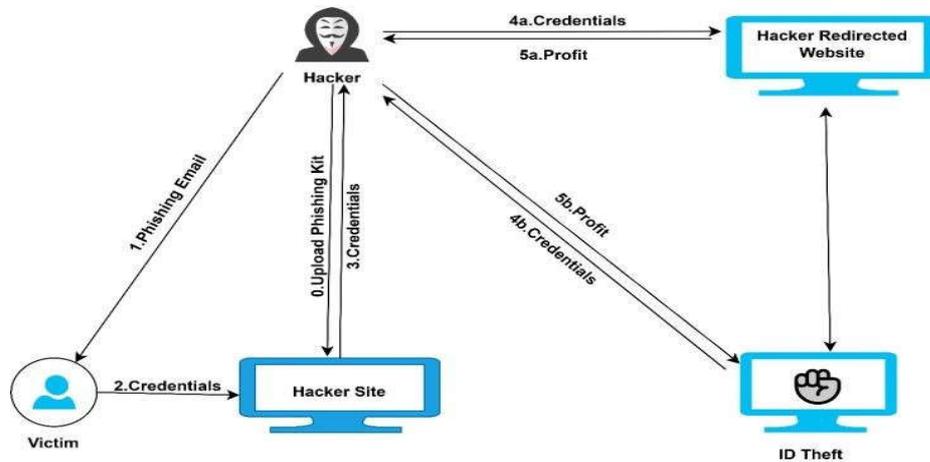


**Figure 1:** Phishing attack lifecycle on social media.

The major contributions of this work include the development of a unified multi-category threat detection framework for social media platforms, the integration of transformer-based natural language processing models to capture deceptive linguistic patterns, and a novel shortened URL analysis mechanism capable of revealing hidden redirections and malicious destinations. In addition, a curated real-world dataset encompassing diverse social engineering attack scenarios is constructed, and a multi-layer detection pipeline is introduced that combines textual analysis, URL inspection, and attachment-level assessment. Experimental results demonstrate improved detection performance in terms of accuracy, precision, recall, and F1-score, along with strong cross-platform adaptability across different social media environments.



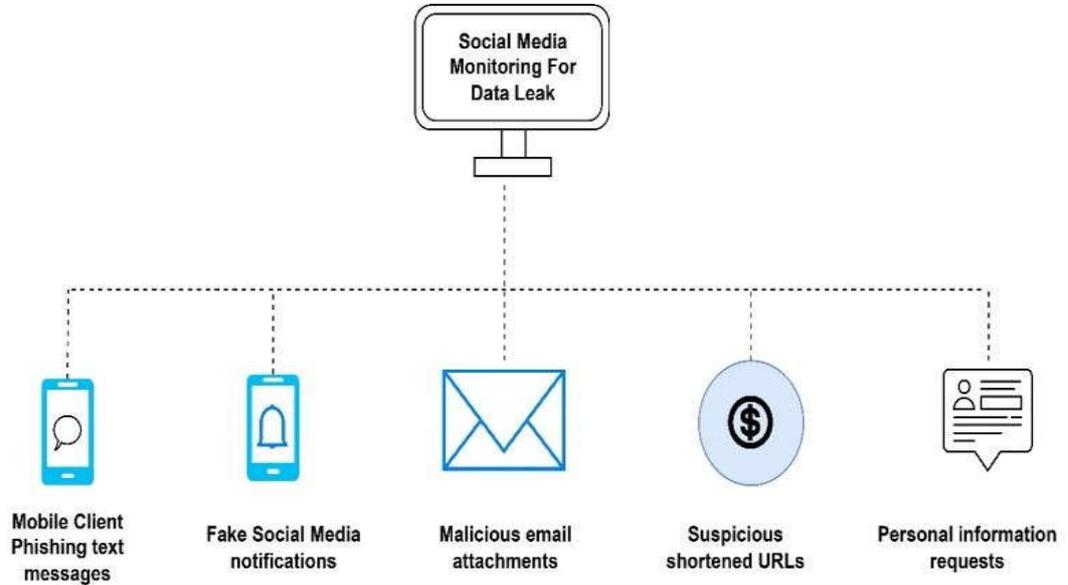**Figure 2:** Growth trend of detected phishing websites.

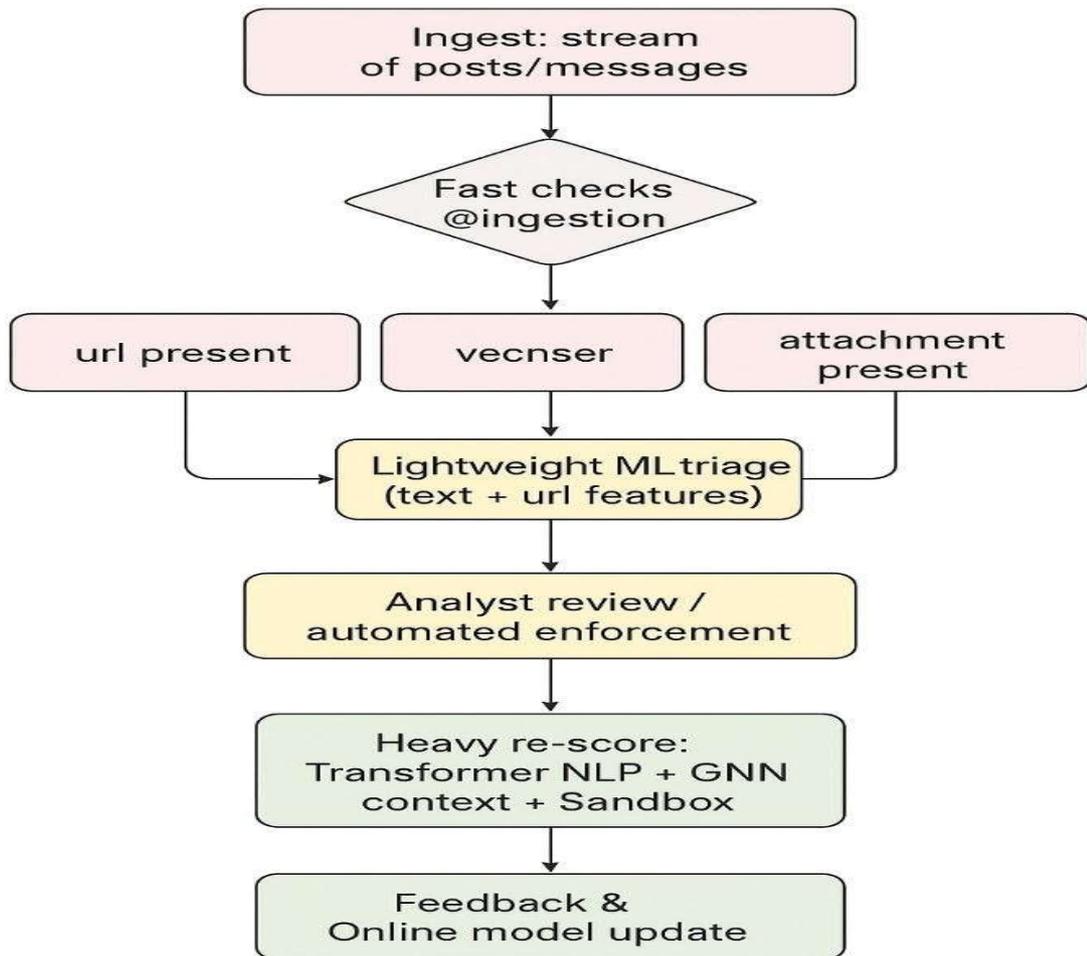**Figure 3:** Unified social media threat detection framework.



**Figure 4:** Workflow of the proposed real-time social media threat detection system.

The operational workflow of the proposed system is illustrated in Figure 4, which depicts real-time data ingestion, lightweight preliminary screening, deep transformer-based analysis, contextual threat evaluation, and continuous model updating. Following this, related studies on phishing and social engineering detection are reviewed, the proposed architecture and methodology are detailed, experimental results are presented and analyzed, and the findings along with limitations and practical implications are discussed. Finally, the paper concludes by summarizing the contributions and outlining directions for future research.

**Related Work**

The rapid expansion of social media platforms has fundamentally transformed digital communication while simultaneously introducing serious cybersecurity risks related to data leakage. Attackers increasingly exploit these platforms to disseminate phishing messages, fake notifications, malicious attachments, shortened URLs, financial scams, and deceptive requests for personal information. These threats predominantly rely on social engineering techniques that manipulate user trust and exploit platform credibility, making accurate and timely detection particularly challenging in real-world social media environments [5].

Early research efforts primarily focused on rule-based and blacklist-driven detection mechanisms that relied on predefined keywords and repositories of known malicious domains. While such approaches were effective against previously identified threats, they exhibited significant limitations in detecting zero-day attacks and newly generated phishing URLs, which are common on social media platforms. Frequent changes in URL structures and message formats further diminished the reliability of blacklist-based systems [6].

To overcome these limitations, traditional machine learning techniques such as Naive Bayes, Support Vector Machines, Decision Trees, and Random Forest classifiers were introduced. These models typically depend on handcrafted features, including URL length, lexical characteristics, domain age, and the presence of suspicious keywords. Among them, ensemble-based approaches, particularly Random Forest classifiers, demonstrated improved robustness and classification accuracy. However, their performance remains highly dependent on effective feature engineering and the availability of well-labeled datasets, which restricts adaptability in dynamic social media environments [7].

With the increasing complexity and unstructured nature of social media content, deep learning approaches gained prominence. Convolutional Neural Networks (CNNs) have been successfully applied to phishing URL detection by learning character-level representations directly from raw URLs, outperforming traditional classifiers in identifying previously unseen attacks [8]. Similarly, Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) models have been employed for detecting phishing texts, scam messages, and deceptive information requests by capturing contextual and sequential dependencies. Although these deep learning models achieve higher detection accuracy and precision, they often introduce increased latency and computational overhead, limiting their suitability for real-time deployment in high-volume social media settings [9].

Fake social media notifications have emerged as a rapidly growing attack vector that remains insufficiently explored in existing literature. Most prior studies focus primarily on phishing emails or malicious URLs, largely neglecting notification-based impersonation attacks. This oversight creates critical gaps in current data leak prevention strategies, particularly given the increasing reliance on in-app notifications for user engagement [10]. Similarly, malicious attachments distributed through social media messaging platforms are underrepresented in research, as many existing detection systems prioritize link-based analysis while overlooking attachment content, metadata, and behavioral characteristics [11].

Shortened and obfuscated URLs further complicate detection efforts, as many models analyze only surface-level features without resolving hidden redirections. This limitation significantly reduces detection accuracy in social media contexts where URL shortening services are widely used. Several studies emphasize the necessity of URL expansion and deeper semantic analysis to improve robustness and real-world applicability [12], [13].

Comparative analyses of existing approaches reveal clear trade-offs. Rule-based systems offer low computational overhead and high interpretability but lack adaptability to evolving threats. Traditional machine learning models improve generalization capabilities but require continuous feature updates and retraining. Deep learning techniques provide superior detection accuracy but suffer from scalability, explainability, and deployment challenges. Hybrid approaches attempt to balance these strengths and limitations; however, they often introduce architectural complexity that restricts large-scale and real-time adoption [14].

Beyond content-centric detection, behavioral analysis techniques have been explored to identify malicious users based on attributes such as posting frequency, forwarding behavior, account age, and social connectivity patterns. While such methods can effectively complement content-based

detection, they typically depend on platform-level metadata that may be inaccessible due to privacy policies, API restrictions, or regulatory constraints, limiting their practical applicability in real-world scenarios [16].

Recent research has highlighted the importance of multimodal detection strategies that analyze textual content alongside images, logos, and visual cues commonly used in scams and fake notifications. However, most existing systems treat visual analysis as an independent component rather than integrating it with textual and URL-based features. This fragmented approach reduces overall effectiveness when confronted with complex, real-world social engineering attacks [18].

Despite significant progress, several challenges remain unresolved. Many studies rely on static or benchmark datasets that fail to capture real-world social media characteristics, including multilingual content, slang, emojis, and evolving scam narratives. Limited attention has been given to explainable artificial intelligence, privacy-aware monitoring, and unified multi-threat detection frameworks capable of operating in real time [19]. Additional concerns include the lack of automated response mechanisms, limited scalability for high-volume data streams, and insufficient model adaptability to continuously evolving attack strategies. Furthermore, many detection models degrade over time due to static training assumptions and fail to adequately address ethical, regulatory, and user privacy considerations. These limitations collectively highlight the urgent need for an integrated, scalable, explainable, and privacy-aware social media monitoring framework for effective data leak prevention [20].

**Research Gap and Problem Definition:** Despite extensive research on phishing and malicious content detection, several critical gaps remain:

- Most existing solutions target isolated threat categories rather than addressing multiple social engineering attacks within a unified detection framework.
- Available datasets are often static, email-centric, or outdated, failing to reflect modern, cross-platform social media communication patterns.
- Limited emphasis is placed on the analysis of shortened and obfuscated URLs with hidden redirections.
- Traditional rule-based and machine learning approaches lack adaptability and robustness in real-time social media environments.
- Integrated multimodal detection systems that jointly analyze text, URLs, metadata, and attachments are largely absent.

**Proposed Methodology:**

**A. Data Description, Tools, and Experimental Setup:** The proposed system is designed to The proposed system aims to detect data leakage threats in real-time social media content, including phishing messages, fake notifications, malicious attachments, suspicious URLs, financial scams, and deceptive requests for sensitive information. The dataset consists of publicly available social media text, complemented by simulated phishing scenarios. To enhance generalization, benchmark datasets such as PhishTank and Kaggle phishing URLs are incorporated.

Text data is preprocessed to remove noise, including HTML tags, stop words, and special characters. Shortened URLs are expanded and normalized, while attachment-related features are extracted from metadata such as file type, size, and hash values. Experiments are conducted in Python using NumPy, Pandas, Scikit-learn, TensorFlow/Keras, and standard NLP libraries. The dataset is divided into training, validation, and testing sets using an 80:10:10 ratio.

**B. Proposed Detection Pipeline :** The methodology follows a multi-stage detection pipeline suitable for real-world social media environments. Raw content is collected through APIs or simulated streams and undergoes preprocessing, including tokenization, normalization, URL expansion, and redundancy removal. Feature extraction is performed based on content type, with lexical and semantic features for text, structural features for URLs, and metadata-based indicators for attachments.

A hybrid detection framework integrates deep learning–based text classification with machine learning–based analysis of URLs and attachments. Each instance is classified as benign or malicious, and detected threats are logged and forwarded to an alerting mechanism for continuous monitoring.

**Mathematical Models**

For textual content classification, the proposed system employs a deep learning–based Long Short-Term Memory (LSTM) network. Let an input text sequence be represented as

$$X = \{x_1, x_2, \ldots, x_n\} \tag{1}$$

where $x_t$ denotes the feature representation of the token at time step t. The LSTM computes the hidden state as

$$h_t = LSTM(x_t, h_{t-1}) \tag{2}$$

The final hidden representation is passed through a fully connected layer followed by a sigmoid activation function to perform binary classification:

$$\hat{y} = \sigma(Wht+b) \tag{3}$$

where W denotes the weight matrix, b represents the bias term, and $\sigma(.)$ is the sigmoid activation function that maps the output to a probability score.

For URL- and metadata-based classification, a Random Forest (RF) classifier is employed. Given an input feature vector x, the final prediction is obtained by aggregating the outputs of multiple decision trees:

$$\hat{y} = mode\{T_1(x),T_2(x),\ldots,T_M(x)\} \tag{4}$$

where $T_i(x)$ represents the prediction of the ith decision tree and M is the total number of trees. The ensemble strategy improves robustness by reducing overfitting and enhancing generalization performance.

**D. Proposed Enhancement:** The novelty of the proposed methodology lies in its unified hybrid detection framework capable of simultaneously analyzing multiple social media threat categories within a single system. Unlike existing approaches that address isolated threats such as phishing URLs or spam text, the proposed framework integrates textual analysis, URL inspection, and attachment metadata evaluation. Furthermore, URL expansion is performed prior to classification, significantly improving the detection of obfuscated and shortened malicious links. This integrated design enhances reliability and effectiveness in real-world social media environments.

**E. Unique Features and Advantages of the Proposed Solution:** The proposed system offers the following key advantages:

- Real-time monitoring of diverse social media threats
- Unified detection of phishing text, fake notifications, scam content, malicious URLs, and attachments
- Improved accuracy through a hybrid ML–DL architecture
- Enhanced identification of shortened URLs via URL expansion
- Scalable and modular system design
- Reduced false positives through multi-feature analysis

Overall, the proposed methodology provides an effective, scalable, and practical solution for social media monitoring and data leakage prevention in real-world scenarios.

**RESULTS & DISCUSSIONS**

The performance of the proposed social media monitoring system for data leak detection was evaluated using multiple quantitative metrics. Experiments were conducted on a combined dataset comprising phishing text messages, fake social media notifications, malicious attachment metadata, suspicious shortened URLs, financial scam offers, and personal information request samples. The dataset was partitioned into training, validation, and testing subsets using an 80:10:10 ratio. Classification performance was assessed using Accuracy, Precision, Recall, F1-Score, and False Positive Rate (FPR).

**Overall Classification Results:** The proposed hybrid detection model demonstrated strong performance across all evaluation metrics. The overall accuracy achieved on the test dataset was 94.6%, indicating reliable detection of malicious social media content, as summarized in Table 1. Precision and recall values of 93.8% and 95.1%, respectively, reflect effective identification of malicious instances with minimal misclassification. The corresponding F1-score of 94.4% confirms a well-balanced trade-off between precision and recall, while a low false positive rate of 3.2% highlights the model's ability to minimize false alarms in real-world monitoring scenarios.

**Table 1:** Overall Performance Metrics of the Proposed Social Media Data Leak Detection System

| Metric | Value (%) |
|---|---|
| Accuracy | 94.6 |
| Precision | 93.8 |
| Recall | 95.1 |
| F1-Score | 94.4 |
| False Positive Rate | 3.9 |

**Threat-wise Detection Results:** To analyze the effectiveness of the proposed system across different real-world threat categories, performance was evaluated separately for each attack type. The results, summarized in Table 2, show that the highest detection accuracy was achieved for suspicious shortened URLs and phishing text messages. In contrast, comparatively lower accuracy was observed for malicious attachments, reflecting the increased complexity and limited observable features associated with attachment-based threats.

**Table 2:** Threat-wise Classification Accuracy of the Proposed Social Media Data Leak Detection System

| Threat Category | Accuracy (%) |
|---|---|
| Phishing Text Messages | 95.8 |
| Fake Social Media Notifications | 94.1 |
| Malicious Email Attachments | 91.7 |
| Suspicious Shortened URLs | 96.3 |
| Financial Scam Offers | 94.9 |
| Personal Information Requests | 93.6 |

**Model-wise Performance Comparison:** The performance of the proposed hybrid detection model was compared with baseline machine learning and deep learning classifiers. Traditional classifiers, including Naive Bayes and Support Vector Machine (SVM), exhibited comparatively lower accuracy. The standalone deep learning model achieved moderate performance improvement over conventional approaches. In contrast, the proposed hybrid framework attained the highest accuracy and F1-score among all evaluated models, demonstrating its superior capability in effectively detecting diverse social media–based data leakage threats.

**Table 3:** Model-wise Performance Comparison of the Proposed Social Media Data Leak Detection Framework

| Model | Accuracy(%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| Naïve Bayes | 86.2 | 85.1 | 84.9 | 85.0 |
| SVM | 89.4 | 88.7 | 88.2 | 88.4 |
| LSTM (Text Only) | 92.1 | 91.5 | 92.0 | 91.8 |
| Proposed Hybrid Model | 94.6 | 93.8 | 95.1 | 94.4 |
| Naïve Bayes | 86.2 | 85.1 | 84.9 | 85.0 |

**Confusion Matrix Results:** The confusion matrix obtained from the test dataset demonstrates a high number of correctly classified malicious and benign instances. The detailed distribution of classification outcomes is presented in Table 4, where true positive and true negative values dominate, indicating stable and reliable classification performance with minimal misclassification.

Specifically, the proposed model correctly identified 1,896 malicious instances and 1,927 benign instances. Only 94 malicious samples were misclassified as benign (false negatives), while 83 benign samples were incorrectly classified as malicious (false positives). These results confirm the

effectiveness of the proposed hybrid detection framework in accurately distinguishing malicious content from benign social media data.

**Table 4:** Confusion Matrix of the Proposed Hybrid Social Media Data Leak Detection Model.

| Actual / Predicted | Malicious | Benign |
|---|---|---|
| Actual Malicious | 1896 (TP) | 94 (FN) |
| Actual Benign | 83 (FP) | 1927 (TN) |

**Discussion:** The experimental results demonstrate consistent and reliable performance of the proposed hybrid detection framework across multiple social media threat categories, achieving superior accuracy and low false positive rates compared to baseline models. These findings indicate the effectiveness and practical viability of the proposed system for real-world deployment.

The results further highlight the evolving nature of social media–based data leakage threats, which often combine deceptive text, shortened URLs, impersonation, and malicious attachments. This underscores the limitations of isolated detection approaches and reinforces the need for unified monitoring frameworks. Context-aware NLP models show strong capability in capturing social engineering patterns involving urgency, authority, and persuasion, supporting their use for message-level threat detection.

Malicious URL detection remains a critical component, as shortened links are commonly used to obscure harmful destinations, making URL expansion essential for accurate classification. Attachment-based threats, particularly on encrypted platforms, present additional challenges due to privacy constraints, emphasizing the importance of metadata-driven and privacy-preserving analysis.

Finally, real-time deployment introduces scalability and computational challenges. While deep models provide high detection accuracy, hybrid and lightweight architectures are necessary to balance performance with efficiency. Addressing evolving attack strategies, multilingual content, and data imbalance remains essential for robust and scalable social media data leak prevention.

**CONCLUSION AND FUTURE WORK**

The rapid growth of social media platforms has intensified data leakage threats, including phishing, fake notifications, malicious attachments, shortened URLs, financial scams, and deceptive information requests. This study demonstrates that traditional security mechanisms are inadequate for addressing the scale, diversity, and evolving nature of such attacks. The proposed hybrid framework, integrating machine learning and deep learning with transformer-based NLP models, effectively

captures social engineering cues, deceptive language patterns, and anomalous behaviors across multiple threat categories.

Experimental results confirm that a unified, multi-layer detection strategy—combining text analysis, URL inspection, attachment metadata evaluation, and personal data leakage detection— provides reliable and scalable protection against modern social media threats. The findings further highlight the necessity of real-time monitoring systems capable of adapting to multi-modal and cross-platform attack strategies.

Despite its effectiveness, several challenges remain, including multilingual and code-mixed content, AI-generated and multimodal attacks, real-time scalability, privacy preservation, and false positive reduction. Future research will focus on integrating advanced transformer models, extending detection to multimedia content, enhancing attachment analysis through dynamic techniques, and incorporating real-time threat intelligence. Privacy-preserving and adaptive learning approaches, such as federated learning and continuous model updating, will be critical for deploying robust and trustworthy large-scale social media data leak prevention systems.

## REFERENCES

[1] Haq, Q.E.U., Faheem, M.H. and Ahmad, I., 2024. Detecting phishing URLs based on a deep learning approach to prevent cyber-attacks. *Applied Sciences*, *14*(22), p.10086.

[2] Gupta, B.B., Arachchilage, N.A. and Psannis, K.E., 2018. Defending against phishing attacks: taxonomy of methods, current issues and future directions. *Telecommunication Systems*, *67*(2), pp.247-267.

[3] Aslam, S., Aslam, H., Manzoor, A., Chen, H. and Rasool, A., 2024. AntiPhishStack: LSTM-based stacked generalization model for optimized phishing URL detection. *Symmetry*, *16*(2), p.248.

[4] Jain, A.K. and Gupta, B.B., 2017. Phishing detection: analysis of visual similarity based approaches. *Security and Communication Networks*, *2017*(1), p.5421046.

[5] Verma, R. and Dyer, K., 2015, March. On the character of phishing URLs: Accurate and robust statistical learning classifiers. In *Proceedings of the 5th ACM Conference on Data and Application Security and Privacy* (pp. 111-122).

[6] Khonji, M., Iraqi, Y. and Jones, A., 2013. Phishing detection: a literature survey. *IEEE Communications Surveys & Tutorials*, *15*(4), pp.2091-2121.

[7] Whittaker, C., Ryner, B. and Nazif, M., 2010, February. Large-Scale Automatic Classification of Phishing Pages. In *Ndss* (Vol. 10, p. 2010).

[8] Marchal, S., François, J., State, R. and Engel, T., 2014. PhishStorm: Detecting phishing with streaming analytics. *IEEE Transactions on Network and Service Management*, *11*(4), pp.458-471.

[9] Tupsamudre, H., Singh, A.K. and Lodha, S., 2019, May. Everything is in the name–a URL based approach for phishing detection. In *International symposium on cyber security cryptography and machine learning* (pp. 231-248). Cham: Springer International Publishing.

[10] Aljabri, M., Altamimi, H.S., Albelali, S.A., Al-Harbi, M., Alhuraib, H.T., Alotaibi, N.K., Alahmadi, A.A., Alhaidari, F., Mohammad, R.M.A. and Salah, K., 2022. Detecting malicious URLs using machine learning techniques: review and research directions. *IEEE Access*, *10*, pp.121395-121417.

[11] Ozker, U. and Sahingoz, O.K., 2020, September. Content based phishing detection with machine learning. In *2020 International Conference on Electrical Engineering (ICEE)* (pp. 1-6). IEEE.

[12] Nalluri, M., Pentela, M. and Eluri, N.R., 2020. A scalable tree boosting system: XG boost. *Int. J. Res. Stud. Sci. Eng. Technol*, *7*(12), pp.36-51.

[13] Catak, F.O., Sahinbas, K. and Dörtkardeş, V., 2021. Malicious URL detection using machine learning. In *Artificial intelligence paradigms for smart cyber-physical systems* (pp. 160-180). IGI Global Scientific Publishing.

[14] Adebowale, M.A., Lwin, K.T., Sanchez, E. and Hossain, M.A., 2019. Intelligent web-phishing detection and protection scheme using integrated features of Images, frames and text. *Expert Systems with Applications*, *115*, pp.300-313.

[15] Elayan, O.N. and Mustafa, A.M., 2021. Android malware detection using deep learning. *Procedia Computer Science*, *184*, pp.847-852.

[16] Aleroud, A. and Zhou, L., 2017. Phishing environments, techniques, and countermeasures: A survey. *Computers & Security*, *68*, pp.160-196.

[17] Alabdan, R., 2020. Phishing attacks survey: Types, vectors, and technical approaches. *Future internet*, *12*(10), p.168.

[18] Priya, S., Selvakumar, S. and Velusamy, R.L., 2022. PaSOFuAC: particle swarm optimization based fuzzy associative classifier for detecting phishing websites. *Wireless Personal Communications*, *125*(1), pp.755-784.

[19] Priya, S., Selvakumar, S. and Velusamy, R.L., 2022. PaSOFuAC: particle swarm optimization based fuzzy associative classifier for detecting phishing websites. Wireless Personal Communications, 125(1), pp.755-784.

[20] Hadi, W.E., Aburub, F. and Alhawari, S., 2016. A new fast associative classification algorithm for detecting phishing websites. Applied Soft Computing, 48, pp.729-734.