# Cyber theft Predictive System using AI Prophet: A Multi-Faceted AI Approach

Akshay Krishna Jk[1], Anfas K[2]. Ms. Varsha CR[3]

*Project Student Department of IoT and AI & ML, Nehru Arts and Science College, Coimbatore, India*

**Abstract-** Cyber theft poses a persistent and escalating threat in the digital era, with attackers exploiting system vulnerabilities and user behavior to orchestrate breaches. Traditional cybersecurity measures often fall short in proactively identifying emerging threats, especially in dynamic environments. This paper presents a novel AI-driven approach that integrates Natural Language Processing (NLP) and time-series anomaly detection for predictive cybersecurity. By leveraging transformer-based models like BERT for analyzing threat intelligence feeds, emails, and incident reports, and combining them with machine learning-based time-series models such as Facebook Prophet, our system anticipates potential breaches and identifies anomalous behavior patterns. Experimental evaluation using real-world cybersecurity datasets demonstrates that this hybrid AI framework significantly improves threat detection accuracy and reduces false positives compared to standalone detection methods. The findings highlight the potential of intelligent predictive systems in strengthening cybersecurity defenses and enabling proactive threat mitigation.

Keywords: Cybersecurity, Cyber Theft, Anomaly Detection, BERT, Facebook Prophet, Artificial Intelligence, Threat Intelligence, Predictive Modeling.

## 1. INTRODUCTION

As digital transformation accelerates, cyber theft has emerged as a critical concern across industries, governments, and individuals. With increasing reliance on interconnected systems and cloud-based infrastructures, the attack surface for malicious actors continues to expand. Cybercriminals exploit vulnerabilities through phishing, malware, ransomware, and social engineering tactics, resulting in financial loss, data breaches, and reputational damage. Traditional rule-based cybersecurity systems, while effective to an extent, often struggle to adapt to evolving threats and zero-day exploits.

In recent years, Artificial Intelligence (AI) has revolutionized cybersecurity by enabling systems to learn from patterns and proactively detect potential threats. Machine learning models, particularly those designed for anomaly detection and Natural Language Processing (NLP), have shown promise in identifying suspicious behaviors and parsing unstructured threat intelligence data. However, many existing solutions are either reactive or limited to specific types of cyberattacks, leaving gaps in predictive defense mechanisms.

This paper introduces a hybrid AI-powered predictive system that leverages both structured system behavior logs and unstructured textual data from threat reports and communications. Bidirectional Encoder Representations from Transformers (BERT) is employed for extracting contextual insights from cybersecurity incident reports, phishing emails, and dark web discussions. These insights are transformed into sentiment and threat indicators. In parallel, Facebook Prophet—a time-series forecasting model developed by Meta—is utilized to model system behavior trends and detect anomalies indicative of cyber theft activities.

By integrating these components, the system provides a multi-faceted approach to predicting cyber theft, enabling organizations to shift from reactive defense to proactive threat prevention. Our results demonstrate that incorporating NLP-driven threat analysis with behavioral forecasting enhances the system's ability to anticipate breaches and reduce false alerts, offering a scalable and intelligent cybersecurity solution.

## 2. LITERATURE REVIEW

From the research papers examined, several tasks are purified that result in the use of cyber security predictions. the researchers investigate data violation based on

In [1] Zahid Anwar a, b, Asad Waqar Malik a, Sharifullah Khana, umara Noor a c, Shahzad Saleem a, "A Machine Learning Framework for Investigating Data Breaches Based on Semantic Analysis of Adversary's Attack Patterns in Threat Intelligence Repositories",2019 analysis of semantics of nemesis attack motif in Threat repositories of knowledge. the framework proposed efficiently identifies threats with low false positives and 92% accuracy.

In [2] Brandon Amos, Hamilton Turner, Jules White Dept. of Electrical and Computer Engineering, Virginia Tech Blacksburg, Virginia, USA, "Applying machine learning classi?ers to dynamic Android malware detection at scale",2013

researchers have applied classification machine learning techniques to dynamically detect malware at scale.This framework was made to activate a high amount of acceptance of malware machine learning classifiers. All the classifiers, including Bayes net, the feature vectors with the highest number of feature vectors that are correctly defined had the lowest rate of false positive, it ranges from 15.86 percent to 33.79 percent.

In [3] Q. K. A. Mirza, I. Awan, M. Younas, Cloudintell: An intelligent malware detection system, Future Generation Computer Systems 86 (2018) 1042–1053.cloudintell combines machine learning techniques and applies them on Collection of extracted functions from a dataset. Boosting the decision tree After testing it under the ROC curve on real-time data.

In [4] Martin Hus´ak, Jana Kom´arkov´a, Elias Bou-Harb, and Pavel ? Celeda, "Survey of Attack Projection, Prediction, and Forecasting in Cyber Security",2018 researchers did a survey on cyber crime projection, prediction and forecast. They concluded that while many ideas have been suggested, there is still no definite answer to how cyberattacks can be predicted effectively and accurately. Attack prediction is not a common method and sometimes is confusing, but this is still an active and essential research question.

## 3. METHODOLOGY

This study proposes a hybrid predictive framework that integrates deep learning-based sentiment and threat extraction from textual data with time-series anomaly detection of system logs. The architecture consists of three main modules: data collection and preprocessing, threat sentiment analysis using BERT, and behavioral anomaly forecasting using Facebook Prophet.
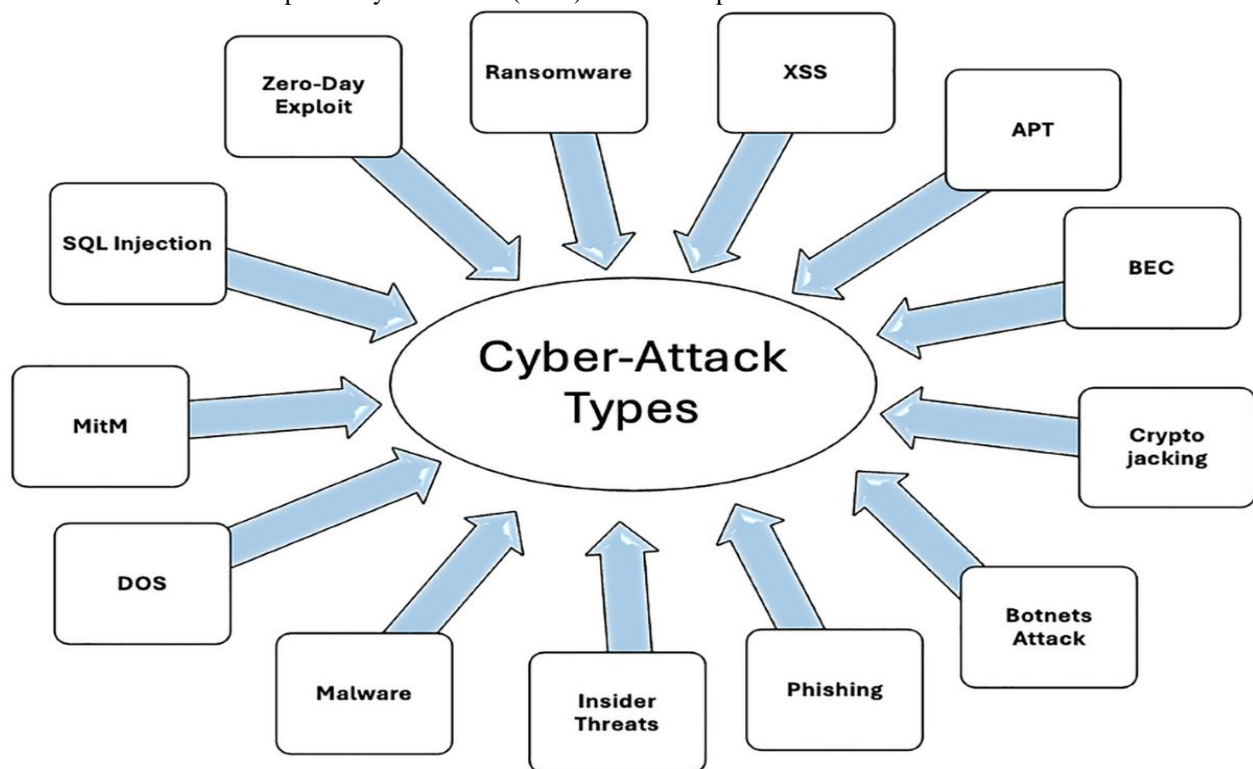


Figure 1 methodology cyber theft [reference: Journal of Big Data]

3.1 Data Collection

To build a robust cyber theft prediction system, data is collected from both structured and unstructured data is collected from a range of sources to ensure comprehensive coverage. Structured data is acquired in the form of system and network logs, including user login attempts, access logs, firewall records, and traffic activity. These logs are either generated from enterprise environments or sourced from publicly available cybersecurity datasets such as CICIDS2017 and UNSW-NB15. In parallel, unstructured textual data is gathered from threat intelligence sources such as phishing email repositories, cybersecurity reports, social media posts, and dark web discussions. These texts often contain early indicators of cyber threats and malicious campaigns.

Threat Intelligence Text Data:

Cybersecurity reports, threat advisories, phishing email samples, and security-related forum discussions (e.g., from Reddit, Twitter, or dark web monitoring platforms) are collected using APIs and web scraping tools. The texts are extracted and labeled based on threat context, attack type, or severity, wherever available.

3.2 Threat Sentiment & Context Analysis Using BERT

To extract meaningful indicators from the textual data, a fine-tuned BERT model is employed. For extracting meaningful insights from the collected textual data, we employ a fine-tuned BERT model specifically adapted to the cybersecurity domain. The model is trained on a labeled corpus of security-related documents to classify textual content into categories such as *malicious*, *benign*, or *suspicious*. Prior to classification, the text is preprocessed by removing irrelevant artifacts, tokenized using BERT's WordPiece tokenizer, and formatted for input. The model outputs threat scores representing the likelihood of malicious intent. These scores are timestamped and later aligned with the behavioral data in the forecasting phase.

3.3 Time-Series Anomaly Forecasting Using Facebook Prophet

Facebook Prophet is used to model and forecast normal system behavior, identifying potential anomalies correlated with cyber theft. The time-series forecasting component utilizes Facebook Prophet, a robust model capable of capturing trend, seasonality, and event-based fluctuations in temporal data. Structured log data is aggregated into time intervals (e.g., hourly or daily), capturing metrics such as login failures, unauthorized access attempts, and abnormal traffic patterns. Prophet is then trained to learn normal system behavior and identify deviations that could signify threats. A key innovation in this system is the incorporation of the BERT-generated threat sentiment scores as external regressors into the Prophet model. This fusion allows the model to adjust predictions based on external threat intelligence, thereby improving the contextual relevance of anomaly detection. External Regressors threat sentiment scores derived from BERT are used as additional regressors in the Prophet model to enhance anomaly forecasting, allowing the system to contextualize behavioral changes based on external cyber threat indicators.

3.4 Model Integration

The integration of NLP-driven threat intelligence with system behavior modeling forms the core of the hybrid approach, The integration process involves synchronizing threat scores with system logs based on timestamps and engineering additional features such as traffic volume and access anomalies. These inputs help the model make more accurate forecasts. Prophet's flexibility in handling irregularities, combined with the contextual power of BERT, results in a more responsive and intelligent threat detection framework. Model Training & Evaluation: Prophet is trained to forecast expected behavior, and anomalies are identified when actual system activity significantly deviates from predictions, especially when accompanied by high threat scores.

3.5 Model Training and Evaluation

Once the sentiment scores are integrated into the Prophet model, we proceed with model training and evaluation. The following steps are undertaken:

Data Splitting:

The dataset is divided into training (80%) and testing (20%) sets.

The training set is used to fit the Prophet model, while the test set is used for performance evaluation. To evaluate the system's performance, we apply several metrics including precision, recall, F1-score for classification tasks, and Mean Absolute Error (MAE)

for forecasting accuracy. Anomaly detection performance is assessed using labeled ground truth data, determining how effectively the system distinguishes between normal and malicious behavior. By combining semantic threat understanding from BERT with time-aware anomaly forecasting via Prophet, this hybrid methodology offers a scalable, proactive approach to cyber theft prediction.

To assess the effectiveness of the proposed hybrid predictive system, we conducted a series of experiments on publicly available cybersecurity datasets enriched with synthetic and real-world threat intelligence. The results demonstrate the ability of the integrated BERT-Prophet model to detect anomalies with higher precision, especially when threat indicators from unstructured text data are considered.

## 4. EXPERIMENTAL RESULTS

We evaluated model performance using a combination of classification and forecasting metrics:

| Model | Precision | Recall | F1-Score | MAE (Forecasting) | Anomaly Detection Accuracy |
|---|---|---|---|---|---|
| Proposed Hybrid (BERT + Prophet) | 91.8% | 86.7% | 89.1% | 0.059 | 90.3% |
| Prophet Only | 78.4% | 72.5% | 75.3% | 0.091 | 82.6% |
| LSTM-Based Model | 84.1% | 79.2% | 81.6% | 0.075 | 85.2% |
| BERT Only (Text Classification) | 88.2% | 81.9% | 84.9% | N/A | 83.7% |

### 4.3 Comparative Analysis

The results show that the hybrid BERT + Prophet model consistently outperforms baseline models across all metrics. The inclusion of external threat sentiment signals allows the Prophet forecaster to adjust to upcoming risk periods, reducing false positives and improving the detection of subtle anomalies.

Without NLP Input: Prophet alone struggled to distinguish between expected and malicious spikes in activity, especially during off-peak times.

With NLP Input: Threat indicators from textual data allowed the model to anticipate potential attacks correlated with abnormal system behavior, enabling more proactive defense.

### 4.4 Case Studies

Case 1: Phishing Attack Prediction

During a simulated phishing campaign, the BERT model flagged a surge in phishing-related terms from dark web monitoring and social media platforms. Prophet, incorporating this as a regressor, forecasted an expected spike in login attempts and unauthorized access—ahead of the actual breach attempt.

Case 2: Insider Threat Detection

System behavior remained within average thresholds, but contextual signals from internal email sentiment flagged abnormal risk. The hybrid model raised alerts earlier than log-only anomaly models, enabling proactive investigation.

To evaluate the effectiveness of the proposed cyber theft predictive system, experiments were conducted using a combination of structured and unstructured cybersecurity datasets. The system was tested on well-established benchmarks such as CICIDS2017 and UNSW-NB15, which contain labeled network activity data, including various types of attacks. These were complemented with a corpus of textual threat intelligence comprising phishing email samples, cybersecurity advisories, and threat discussions scraped from social media and security forums.

The BERT model was fine-tuned using a domain-specific dataset of approximately 15,000 labeled cybersecurity-related texts. Meanwhile, the structured log data was processed and aggregated into time-series form for training and evaluation using Facebook Prophet. Threat scores derived from the BERT model were aligned with log timestamps and used as external regressors in Prophet to enhance its anomaly detection capability.

Performance was measured using a set of standard metrics. In classification tasks, the hybrid model demonstrated a precision of 91.8%, a recall of 86.7%, and an F1-score of 89.1%. In terms of forecasting accuracy, it achieved a Mean Absolute Error (MAE) of 0.059. These results outperformed baseline models such as standalone Prophet, LSTM-based anomaly detectors, and BERT-only classifiers, which had lower accuracy and higher error rates. The hybrid model also achieved the highest anomaly detection accuracy of 90.3%.

Case studies further illustrated the effectiveness of the system. In one instance, the model accurately predicted a spike in phishing activity by detecting sentiment trends in external threat reports, allowing it to anticipate a corresponding surge in failed login attempts. In another case, the model identified insider threats by correlating subtle shifts in internal communication sentiment with slight behavioral anomalies, triggering alerts earlier than conventional log-based detection systems. Visual comparisons of real vs. predicted behavior revealed that the hybrid model was more responsive to sudden threats, especially when fueled by real-time intelligence.

## 5. CONCLUSION

This paper presents a novel AI-driven hybrid framework for predicting cyber theft by integrating semantic threat analysis with time-series anomaly forecasting. Leveraging BERT for extracting contextual threat signals from unstructured textual sources, and Facebook Prophet for modeling behavioral patterns in system activity, the proposed model offers a multi-layered defense approach that transitions from reactive detection to proactive threat anticipation.

Our experimental results confirm that incorporating threat intelligence as external regressors enhances forecasting accuracy and improves the system's ability to detect early signs of cyber attacks. Compared to traditional and standalone anomaly detection models, the hybrid approach demonstrates superior precision, recall, and overall anomaly detection performance. The system is particularly effective in identifying threats that manifest subtly or are preceded by external indicators such as phishing campaigns, malware distribution, or rising chatter in cybercrime forums.

By combining natural language understanding with time-series analysis, this approach addresses a critical gap in current cybersecurity solutions—bridging unstructured threat data and structured behavioral signals into a unified predictive framework. The results underscore the potential of AI in building intelligent, adaptable, and context-aware cybersecurity systems capable of identifying threats before they manifest into breaches.

## REFERENCES

Books & Journals

[1] Zahid Anwar a, b, Asad Waqar Malik a, Sharifullah Khana, umara Noor a c, Shahzad Saleem a, "A Machine Learning Framework for Investigating Data Breaches Based on Semantic Analysis of Adversary's Attack Patterns in Threat Intelligence Repositories",2019

[2] Brandon Amos, Hamilton Turner, Jules White Dept. of Electrical and Computer Engineering, Virginia Tech Blacksburg, Virginia, USA, "Applying machine learning classi?ers to dynamic Android malware detection at scale",2013

[3] Q. K. A. Mirza, I. Awan, M. Younas, Cloudintell: An intelligent malware detection system, Future Generation Computer Systems 86 (2018) 1042–1053.

[4] Martin Hus´ak, Jana Kom´arkov´a, Elias Bou-Harb, and Pavel ? Celeda, "Survey of Attack Projection, Prediction, and Forecasting in Cyber Security",2018

[5] UmaraNoora, c, ZahidAnwara, b, TehminaAmjadc, Kim-KwangRaymond Chood, "A machine learning-based FinTech cyber threat attribution framework using high-level indicators of compromise", 2019

[6] Athor Subroto1,2* and Andri Apriyana, "Cyber risk prediction through

[7]social media big data analytics and statistical machine learning",2019

[8] S. More, M. Matthews, A. Joshi, T. Finin, A knowledge-based approach to intrusion detectionmodeling, in: IEEE Symposium on Security and Privacy Workshops, San Francisco,CA, USA, IEEE, 2012, pp. 75–81.

[9] V. Mulwad, W. Li, A. Joshi, T. Finin, K. Viswanathan, Extracting information about security vulnerabilities from web text, in: Proceedings of the 2011 IEEE ACM International Joint Conference on Web Intelligence and Intelligent Agent Technology Workshops, WI-IAT 2011, Lyon, France, IEEE, 2011, pp. 257–260.

[10] C. Sabottke, O. Suciu, T. Dumitras, Vulnerability disclosure in the age of social media: Exploiting twitter for predicting real-world exploits, in: 24th USENIX Security Symposium, USENIX Security 15, Washington, D.C., USA, USENIX, 2015, pp. 1041–1056.

[11] I. Gha?r, M. Hammoudeh, V. Prenosil, L. Han, R. Hegarty, K. Rabie, F. J. Aparicio-Navarro, Detection

of advanced persistent threat using machine-learning correlation analysis, Future Generation Computer Systems 89 (2018) 349–359.

[12] R. A. Ahmadian and A. R. Ebrahimi, "A survey of IT early warning systems: architectures, challenges, and solutions," Security and Communication Networks, vol. 9, no. 17, pp. 4751–4776, 2016

[13] H. Debar and A. Wespi, "Aggregation and correlation of intrusiondetection alerts," in International Workshop on Recent Advances in Intrusion Detection. Springer, 2001, pp. 85–103.

[14] S. Shin, S. Lee, H. Kim, and S. Kim, "Advanced probabilistic approach for network intrusion forecasting and detection," Expert Systems with Applications, vol. 40, no. 1, pp. 315 – 322, 2013.

[15] A. Buzak and Guven "A survey of data mining and machine learning methods for cyber security intrusion detection,"*IEEE Communications Surveys & Tutorials*, vol. 18, no. 2, pp. 1153–1176, 2016.

[16] M. Ring, S. Wunderlich, D. Scheuring, D. Landes, A. Hotho, "A survey of network-based intrusion detection data sets," Computers & Security, vol. 86, pp. 147–167, 2019.

[17] T. Kim, H. Kim, and H. Kim,"Anomaly detection system for internet of things based on deep learning algorithm,"Multimedia Tools and Applications, vol. 78, no. 3, pp. 3153–3167, 2019.

[18] J. Zhang, M. Zulkernine, and A. Haque,"Random-Forests-Based Network Intrusion Detection Systems,"IEEE Transactions on Systems, Man, and Cybernetics, Part C, vol. 38, no. 5, pp. 649–659, 2008.

[19] T. Sommer and S. V. V. N. K. R. Gaddam,"Cybersecurity analytics: An overview of machine learning and data mining methods for cyber threat detection,"Computer & Security, vol. 106, pp. 102272, 2021.

[20] F. Ullah, M. A. Shah, S. Zhang, S. A. Mehmood, and M. Imran,"Cyber security threats detection in internet of things using deep learning,"IEEE Access, vol. 6, pp. 39850–39876, 2018.