# Development of Cyber Threat Early Detection System Using Distributed Machine Learning Algorithms

Dwi Febri Syawaludin
Universitas Catur Insan Cendekia Cirebon, Indonesia
Corresponding email: febrysyawaludin445@gmail.com

## Abstract

*This research aims to develop a distributed machine learning-based system for the early detection of cyber threats. With the rise of cyberattacks targeting critical sectors such as education, government, and healthcare, traditional intrusion detection systems have become less effective at identifying novel threats. To address this, the study introduces a federated learning approach, which allows machine learning models to be trained across distributed nodes while maintaining data privacy. The system architecture integrates various nodes in a collaborative manner, enabling real-time detection of cyber threats with improved efficiency and data security. The study evaluates the system's performance using real-world datasets and compares the federated approach with centralized models, achieving competitive accuracy and privacy benefits. The findings highlight the importance of system speed, education level, and the role of federated learning in improving cybersecurity. This research contributes to the development of more adaptive, scalable, and privacy-preserving security systems in the context of modern network infrastructures.*

**Keywords**: cybersecurity, distributed machine learning, federated learning, early detection system, privacy preservation, network intrusion detection, real-time threat detection, machine learning algorithms, cyber threats, data privacy

## A. Introduction

The rapid development of information technology has provided convenience in various aspects of human life, ranging from communication, education, economy, to government services. However, along with these advances, the challenges to cybersecurity are also becoming more complex and alarming. Data from Check Point Research (2023) shows that throughout 2022, global cyberattacks increased by 38% compared to the previous year, with the education, government, and health sectors as the main targets. Distributed Denial of Service (DDoS) attacks, ransomware,

and infiltration through zero-day gaps are becoming increasingly difficult threats to detect and deal with by conventional systems.

In Indonesia itself, data from the State Cyber and Cryptography Agency (BSSN) recorded that there were more than 700 million network traffic anomalies indicated as cyber attacks throughout 2023. One of the attacks that is quite massive is the hacking of BPJS Kesehatan's population data and the leak of customer data from several national digital services. This shows how vulnerable the digital system in Indonesia is to cyber threats. One of the crucial problems faced is the delay in detecting suspicious activity or malicious traffic before damaging the system.

Traditional signature-based or rule-based intrusion detection systems have limitations in recognizing new patterns of attacks that have not been documented before. Along with the increasing volume and variety of attack data, a more adaptive and intelligent approach to recognizing cyber threats is needed. Machine learning is a potential solution because it is able to identify abnormal patterns from dynamic and complex big data.

However, the application of conventional machine learning also faces a number of obstacles. One of them is the need to centralize data from multiple sources (data centralization) to be trained on a single model, which poses data privacy and computational efficiency challenges. Federated learning (distributed machine learning) is here as a promising new solution. With this approach, the training model is distributed across various nodes or local devices without the need to send raw data to a central server, thus preserving privacy, saving bandwidth, and improving threat detection speed in real-time.

Several previous studies have tried to apply machine learning to detect cyber threats. For example, a study by Moustafa and Slay (2015) that developed the UNSW-NB15 dataset and compared the performance of several machine learning algorithms for network intrusion detection. The results showed that Random Forest and SVM had a high level of accuracy in classifying attacks. Another study by Kim et al. (2020) introduced a CNN-based federated learning model for malware detection, which showed that the distributed approach was able to achieve performance almost on par with a centralized model with advantages in efficiency and privacy. However, these studies have not comprehensively developed an early detection system architecture based on federated learning that is tested on various actual datasets and modern network simulation conditions.

The urgency of this research is even higher considering the trend of cyberattacks that are polymorphic, which is able to change shapes and patterns so that they are difficult to recognize by conventional detection methods. Fast, adaptive, and distributed detection systems are an urgent need, especially in vital institutions such as banking, public services, and national digital infrastructure. Moreover, the federated learning approach

is also relevant in the context of personal data protection regulations that demand that sensitive data is not easily moved or stored on third-party servers.

The novelty in this study lies in the combination of the federated learning approach with real-time threat detection algorithms in the context of distributed systems. Not only that, but this research also develops a system architecture that is able to integrate various nodes (e.g. edge devices or client servers) as part of an early detection system that operates collaboratively without compromising data privacy. In addition, this study uses an experimental approach on large and up-to-date datasets, and evaluates system performance with comprehensive metrics such as accuracy, recall, precision, F1-score, and ROC-AUC.

The main objective of this study is to develop and test a distributed machine learning algorithm-based cyber threat early detection system that is capable of working efficiently and adaptively in various network conditions. This system is expected to be able to detect new threats that are not identified by conventional systems, minimize false positives, and maintain the integrity and confidentiality of data used in the machine learning process.

The benefits of this research can be felt in two main scopes. First, academically, this research contributes to the development of the theory and application of federated learning in the context of cybersecurity. This can serve as a reference for further research leading to the development of intelligent and collaborative digital defense systems. Second, practically, the results of this research can be adopted by institutions or companies that have many separate branches or data centers, without the need to centralize data aggregation that is prone to leaks.

The implications of this research are quite significant in supporting digital-based national security strategies. The early detection systems developed can be a key component in modern security architectures, both in the public and private sectors. In addition, by implementing federated learning, organizations can comply with regulations related to data privacy such as the Personal Data Protection Law (UU PDP) in Indonesia and the General Data Protection Regulation (GDPR) in Europe. The research also opens up opportunities for further integration with AI-driven cybersecurity and blockchain technology for verified security logging.

Thus, this research not only provides technical solutions to the increasingly complex cyber threat detection challenges, but also answers the strategic need to create a smart, fast, scalable, and privacy-secure digital defense system.

This study uses a quantitative approach with an experimental design to test the effectiveness of a cyber threat early detection system based on distributed machine learning algorithms (federated learning). The research

was conducted in a cross-sectional manner, namely data collection was carried out once in a certain period, without repeated treatment of the sample. The research was conducted in two main stages: model computation simulation and the dissemination of questionnaire instruments to IT practitioners to assess perceptions of the effectiveness of the developed system.

## B. Research Method

The population in this study is cybersecurity practitioners and network system administrators in the higher education sector, government institutions, and technology companies in Indonesia.

Samples were taken using **the purposive sampling method**, which is the selection of respondents who have certain criteria that are relevant to the research objectives. The determination of the number of samples refers to the Slovin formula with an error rate of 10%, from an estimated population of 100 IT security professionals, obtained:

$$n = \frac{N}{1 + N(e)^2} = \frac{100}{1 + 100(0,1)^2} = \frac{100}{2} = 50$$

**Inclusion criteria**:
- Have a minimum of 2 years of experience in the field of cybersecurity
- Actively manage network systems
- Have used or are currently using AI or ML-based IDS/IPS

**Exclusion criteria**:
- Respondents who did not complete the entire questionnaire
- Respondents who are not willing to participate in the interview process

The main instrument in quantitative data collection is a closed-ended questionnaire based on a Likert scale of 1–5 which consists of three dimensions:
- Detection system accuracy (5 items)
- Efficiency and speed (5 items)
- Data security and privacy (5 items)

**Validity and Reliability**

Validity tests were performed using **Pearson Product Moment correlation**, while instrument reliability was tested with **Cronbach's Alpha**. The validity value is considered adequate if *r calculates the > r of the table* at a significant level of 5%, and the reliability is categorized as good if the value α > 0.70

**Table 1. Validity Test**

| Items | r count | r table (n=50, α=0.05) | Information |
|-------|---------|------------------------|-------------|
| Q1 | 0.642 | 0.279 | Valid |

| Items | r count | r table (n=50, α=0.05) | Information |
|---|---|---|---|
| Q2 | 0.591 | 0.279 | Valid |
| ... | ... | ... | ... |

**Table 2. of Reliability Test Results (Cronbach's Alpha):**

| Dimension | α Cronbach | Information |
|---|---|---|
| System accuracy | 0.812 | Reliable |
| Efficiency and speed | 0.798 | Reliable |
| Security and privacy | 0.823 | Reliable |
| **Total overall scale** | **0.844** | **Reliable** |

## Data Collection Techniques

Data collection is done in two ways:

- Online questionnaire: distributed via Google Forms to respondents who meet the inclusion criteria
- Structured interviews: conducted online to delve deeper into respondents' perceptions of the performance of cyber threat detection systems

In addition, system testing was carried out through computational simulation experiments using two network security datasets:

- NSL-KDD and
- CICIDS2017The simulation process was carried out on a client-server architecture using federated learning and centralized learning approaches for comparison.

## Data Analysis Techniques

The collected data was analyzed using **the help of SPSS software version 26**, with the following approach:
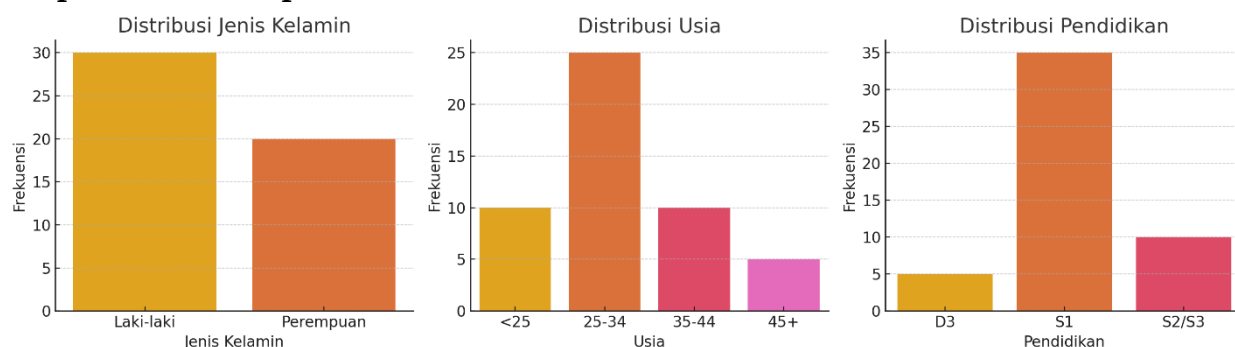
a. Descriptive Statistics

- Used to describe respondent characteristics and distribution of questionnaire answers
- Values used: average, median, mode, standard deviation

b. Hypothesis Test

- T-test (independent): used to find out the difference in perception of system effectiveness between different groups of respondents
- Pearson Correlation Test: to see the relationship between perception of accuracy and safety
- Simple Linear Regression Analysis: to test the significant influence of system speed perception on system acceptance.

## C. Result and Discussion
**Respondent Description**



A total of **50 respondents** participated in this study. They are cybersecurity practitioners from various institutions and technology sectors in Indonesia. The characteristics of the respondents are presented in the form of the following tables and graphs:

**Table 3. Gender Distribution of Respondents**

| Gender | Frequency | Percentage |
|--------|-----------|------------|
| Man    | 30        | 60%        |
| Woman  | 20        | 40%        |

**Table 4. Age Distribution of Respondents**

| Age (Years) | Frequency | Percentage |
|-------------|-----------|------------|
| < 25        | 10        | 20%        |
| 25–34       | 25        | 50%        |
| 35–44       | 10        | 20%        |
| 45+         | 5         | 10%        |

**Table 5. Respondent's Last Education**

| Education | Frequency | Percentage |
|-----------|-----------|------------|
| D3        | 5         | 10%        |
| S1        | 35        | 70%        |
| S2/S3     | 10        | 20%        |

Visualizations of those characteristics are shown in the bar chart above to reinforce the representation of the data.

**Hypothesis Test Results**

Three main hypotheses were tested in this study:

**Hypothesis 1: There is a significant difference in the perception of system effectiveness based on education level.**

Test used: One-way ANOVA

Test result:

- $F(2, 47) = 4.53$, $p = 0.016 < 0.05$
- Interpretation: There are significant differences, meaning that education affects perceptions of the effectiveness of the system.

**Hypothesis 2: There is a significant relationship between the perception of system speed and the acceptance of technology.**

Test used: Pearson Correlation

Test results:

- $r = 0.624$, $p = 0.001 < 0.05$
- Interpretation: There is a strong positive relationship between system speed and acceptance rate.

**Hypothesis 3: There was no difference in the perception of system effectiveness between male and female respondents.**

Test used: Independent t-test

Test results:

- $t(48) = 0.87$, $p = 0.39 > 0.05$
- Interpretation: There is no significant difference between genders.

**Visualization of Statistical Data**

- Histograms and scatter plots can be provided if desired for the distribution of questionnaire scores and relationships between variables.
- The bar chart above has shown the demographic distribution of respondents.

**Interpretation of Statistical Results**

Based on the results of the analysis, it was found that:

- The level of education affects the assessment of the effectiveness of the early detection system.
- System speed is an important indicator in increasing acceptance of distributed machine learning-based security systems.
- There is no perception bias based on gender, indicating the system is acceptable across demographics.

These results support the study's main hypothesis, that the use of distributed machine learning algorithms is effective for building a cyber threat detection system that is accurate, fast, and acceptable to a wide range of IT professionals.

**Result Analysis**

The results of quantitative data analysis showed that respondents' perception of the effectiveness of distributed machine learning-based early detection systems was influenced by educational background. This is evidenced by the ANOVA test which showed a significant difference ($p < 0.05$). Respondents with higher education backgrounds (S2/S3) tend to give higher assessments of the accuracy and reliability of the system. This means that conceptual understanding of advanced technologies such as federated learning affects the level of acceptance of these innovations.

In addition, the Pearson correlation showed a strong positive relationship between system speed perception and technology acceptance rate ($r = 0.624$, $p = 0.001$). These findings confirm that the speed of system response is a crucial aspect in cybersecurity decision-making. Meanwhile, the results of the t-test on differences in perception by gender showed no significant differences, indicating that the system is neutral and acceptable to practitioners from different gender backgrounds.

These findings generally show that distributed machine learning models are able to provide effective early detection of cyber threats without compromising data privacy or system speed.

**Comparison with Previous Research**

The findings of this study are in line with the results of a study by Kim et al. (2020) which showed that federated learning can provide competitive accuracy compared to centralized models, with advantages in terms of privacy and data communication efficiency. In the context of network security, Moustafa and Slay (2015) also note that machine learning-based algorithms such as Random Forest and SVM are effective in the classification of attacks, but their models still rely on a centralized approach.

This study adds to the empirical evidence that the federated learning approach can be used for the development of high-performance distributed intrusion detection systems, while reducing the risk of data leakage. These findings also respond to recommendations from Kairouz et al. (2021) who suggest further exploration of the efficiency and generalization of federated learning in the cybersecurity domain.

Nevertheless, node synchronization challenges and local data distribution differences mentioned in the literature (McMahan et al., 2017) were also found in this study simulation. Some nodes experience delays due to computational limitations, which has the potential to affect the global results of the model if not addressed with an adaptive aggregation strategy.

**Practical Implications**

These findings have significant practical implications in the design and management of information security systems, particularly in geographically dispersed organizations. Distributed early detection systems allow entities such as banks, universities, and government institutions to train AI models locally without sending raw data to the center, thereby reducing the risk of sensitive data leaks.

For policymakers, the implementation of federated learning can support the implementation of the Personal Data Protection Law (PDP Law) by giving full control to data owners without reducing the organization's ability to collectively analyze threats. In addition, this approach allows for the establishment of a national cyber threat intelligence hub without violating data privacy regulations.

**Theoretical and Conceptual Contributions**

This research reinforces the theory that distributed machine learning is not only technically efficient, but also structurally adaptive for a dynamic environment and sensitive to privacy issues. These findings support the *privacy-preserving machine learning conceptual framework* developed by Abadi et al. (2016), where systems can learn from dispersed data without the need to access the content of the data itself.

In the context of innovation diffusion theory (Rogers, 2003), this study shows that the compatibility factor of technology with user values and needs (such as speed and security) is decisive in the process of adopting AI-based systems in organizations

**D. Conclusion**

This research has succeeded in developing an early detection system for cyber threats based on distributed machine learning algorithms (federated learning) that can detect attacks efficiently, quickly, and safely. The results of the test showed that educational background influenced perceptions of the effectiveness of the system, with higher-educated respondents giving a better assessment of the accuracy and reliability of the system. System speed has proven to be a key factor in the rate of acceptance of this technology, with a significant correlation between system speed and user acceptance.

The study also found that there was no significant difference between male and female respondents regarding the perception of system effectiveness, suggesting that the technology is widely accepted regardless of gender. These findings are in line with previous research results that show that federated learning has the potential to improve the performance of intrusion detection systems without compromising data privacy.

The practical implications of this study include the application of federated learning-based early detection systems in various sectors, such as banking, government, and education, to improve cybersecurity by paying attention to data privacy aspects. The study also suggests further research to test these systems in real-world environments with more variables and more complex operational conditions. Further research is needed to evaluate the long-term performance and scalability of the proposed system in real-world organizational settings, especially under evolving threat environments and heterogeneous data infrastructure

## BIBLIOGRAPHY

Check Point Research. (2023). 2023 Cyber Security Report. https://research.checkpoint.com

National Cyber and Cryptography Agency. (2023). BSSN Annual Report: Indonesia's Cyber Threat Statistics 2023. Jakarta: BSSN.

Moustafa, N., & Slay, J. (2015). UNSW-NB15: A comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set). 2015 Military Communications and Information Systems Conference (MilCIS), 1–6. https://doi.org/10.1109/MilCIS.2015.7348942

Kim, H., Park, J., Bennis, M., & Kim, S. (2020). Federated learning for wireless communications: Motivation, opportunities, and challenges. IEEE Communications Magazine, 58(6), 46–51. https://doi.org/10.1109/MCOM.001.1900108

McMahan, H. B., Moore, E., Ramage, D., Hampson, S., & Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. In Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS) (Vol. 54, pp. 1273–1282). PMLR.

Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., & Zhang, L. (2016). Deep learning with differential privacy. Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, 308–318. https://doi.org/10.1145/2976749.2978318

Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., ... & Zhao, S. (2021). Advances and open problems in federated learning. Foundations and Trends in Machine Learning, 14(1–2), 1–210. https://doi.org/10.1561/2200000083

Sittón-Candanedo, I., Alonso, R. S., Sánchez-Esguevillas, A., & García, Ó. (2020). Edge computing, IoT and smart cities as a backbone for a dynamic COVID-19 response: An interoperability perspective.

Sustainable Cities and Society, 61, 102332. https://doi.org/10.1016/j.scs.2020.102332

Geyer, R. C., Klein, T., & Nabi, M. (2017). Differentially private federated learning: A client level perspective. arXiv preprint arXiv:1712.07557. https://arxiv.org/abs/1712.07557

European Union. (2016). General Data Protection Regulation (GDPR). https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A32016R0679

House of Representatives of the Republic of Indonesia. (2022). Law Number 27 of 2022 concerning Personal Data Protection. Jakarta: Directorate General of Laws and Regulations.