

Fingerprinting Deep Learning Models via Network Traffic Patterns in Federated Learning

Md Nahid Hasan Shuvo
mshuvo@gmu.edu
George Mason University
Fairfax, Virginia, USA

Moinul Hossain
mhossa5@gmu.edu
George Mason University
Fairfax, Virginia, USA

ABSTRACT

Federated Learning (FL) is increasingly adopted as a decentralized machine learning paradigm due to its capability to preserve data privacy by training models without centralizing user data. However, FL is susceptible to indirect privacy breaches via network traffic analysis—an area not explored in existing research. The primary objective of this research is to study the feasibility of fingerprinting deep learning models deployed within FL environments by analyzing their network-layer traffic information. In this paper, we conduct an experimental evaluation using various deep learning architectures (i.e., CNN, RNN) within a federated learning testbed. We utilize machine learning algorithms, including Support Vector Machines (SVM), Random Forest, and Gradient-Boosting, to fingerprint unique patterns within the traffic data. Our experiments show high fingerprinting accuracy, achieving 100% accuracy using Random Forest and around 95.7% accuracy using SVM and Gradient Boosting classifiers. This analysis suggests that we can identify specific architectures running within the subsection of the network traffic. Hence, if an adversary knows about the underlying DL architecture, they can exploit that information and conduct targeted attacks. These findings suggest a notable security vulnerability in FL systems and the necessity of strengthening it at the network level.

CCS CONCEPTS

• **Computing methodologies** → **Distributed artificial intelligence**; **Machine learning**; • **Security and privacy**;

KEYWORDS

Federated Learning, Fingerprint, Machine Learning, Deep Learning, CNN, RNN, Network Traffic

ACM Reference Format:

Md Nahid Hasan Shuvo and Moinul Hossain. 2025. Fingerprinting Deep Learning Models via Network Traffic Patterns in Federated Learning. In

This is the author's version of the work. It has been accepted for publication in the Proceedings of the 2025 ACM Workshop on Wireless Security and Machine Learning (WiseML 2025), July 3, 2025, Arlington, VA, USA. The final version will be available in the ACM Digital Library.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference'17, July 2017, Washington, DC, USA

© 2025 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

Proceedings of ACM Conference (Conference'17). ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

1 INTRODUCTION

Federated learning (FL) has improved the training process of artificial intelligence (AI) models by allowing multiple devices to collaborate without sharing raw data [12]. This method limits data access, hence improving privacy; it is especially advantageous in fields like the Internet of Things (IoT) networks, autonomous robotics, and healthcare [6]. Federated learning (FL) primarily shares model updates and enables local model training, hence obviating the necessity of transmitting private information to a centralized server. However, FL is not immune to security concerns despite its privacy-preserving characteristics. Although extensive research has concentrated on direct attacks, such as data poisoning, membership inference attacks, and model inversion, there has been little focus on indirect privacy concerns, especially vulnerabilities associated with network traffic analysis [7],[21].

Federated learning showcases unique features, particularly its dependency on clients and servers for sharing model updates throughout training phases. Although these updates are encrypted, their network traffic features (e.g., packet sizes, transmission direction, and interarrival time) may reveal private information. Previous studies have demonstrated that utilizing network traffic makes it possible to do device fingerprinting in federated learning. Consequently, an attacker could actively monitor and identify the unique devices from their communication patterns [5]. However, a critical and unexplored issue is whether one can fingerprint the deep learning models by analyzing network traffic patterns in FL.

Due to the different computation properties of deep learning (DL) architectures in FL, Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN) utilize system resources differently, especially when sending data over networks. At present, there is no prior research exploring whether FL traffic leakage information can be used to fingerprint the DL models being trained in the FL systems. This research gap is critical because if an adversary can identify these neural network structures from network leakage, then it will enable the malicious actor to develop and launch highly effective architecture-specific targeted cyberattacks, such as creating adversarial attacks explicitly designed for CNNs [23] and RNNs [22]. Thereby, this study aims to address this research gap by investigating the possibility of DL architecture fingerprinting using network traffic data, such as timing information, packet size distribution, and traffic direction. If these core models in FL systems can be fingerprinted this way, it would introduce a new class of security risk, potentially undermining the privacy benefits of FL.

This study investigates the aforementioned research gap of how standard machine learning techniques can effectively fingerprint deep learning architectures deployed in federated learning environments by analyzing network layer traffic patterns. We aim to validate the feasibility of fingerprinting deep learning models using easily accessible network metadata. To achieve this goal, we perform a systematic experiment using NVIDIA GPU-equipped clients and servers running various deep-learning architectures and sniff their network traffic using a packet sniffing tool, Wireshark. We deployed popular machine-learning classification algorithms to fingerprint the DL architecture, including Random Forest, Support Vector Machine, and Gradient Boosting. Our experimental results confirm that deep learning architectures exhibit distinguishable network traffic patterns in the FL system, making them vulnerable to fingerprinting attacks.

In this paper, we make the following key contributions: (1) We propose a novel method to fingerprint deep learning architecture in federated learning using network-layer traffic patterns. (2) We developed a controlled FL testbed to analyze the network pattern using ideal and noisy network conditions. (3) Using traditional machine learning algorithms, we show the feasibility of this fingerprinting attack, which poses security threats to the FL environments. (4) Finally, we discuss the implications of our findings, highlight current limitations, and propose future research directions to improve both fingerprint attacks and potential defense mechanisms.

2 RELATED WORKS

Traffic fingerprinting (TF) is a widely studied traffic analysis approach for identifying objects such as mobile applications, website browsing, or devices according to the distinctive characteristics or behavioral patterns of network traffic. Network managers use TF for security monitoring and filtering [10], [11], whereas adversaries leverage TF to intercept sensitive information [1], [19].

Website fingerprinting has garnered significant attention in the fields of web security and privacy. Despite the safeguards provided by encryption technologies, such as HTTPS and Tor, Panchenko et al. demonstrated the potential to analyze encrypted online traffic to determine the websites accessed by users [14]. In [15], Rahman et al. showed that a passive local eavesdropper could utilize website fingerprinting to reveal the web browsing activities of Tor users. Furthermore, a significant area of study is device fingerprinting, which aims to recognize individual devices based on their distinct hardware and software characteristics. Laperdrix et al. examined browser-based device fingerprinting and demonstrated that seemingly benign information—such as installed plugins, screen resolution, and system fonts—can uniquely identify devices, thereby posing significant privacy concerns [8]. In [16], researchers have extensively discussed fingerprinting industrial IoT devices from network traffic information.

On top of devices and websites, fingerprinting has been extensively applied in mobile application identification, where network flow analysis can be used to identify mobile applications. In [18], Taylor et al. explored how mobile apps generate unique traffic signals recognizable via network-layer data, thereby permitting attackers to determine app types and usage behaviors. In addition, Li et al. show that the application's traffic data contains multiplexed user traffic. Hence, they provide a robust app fingerprinting method

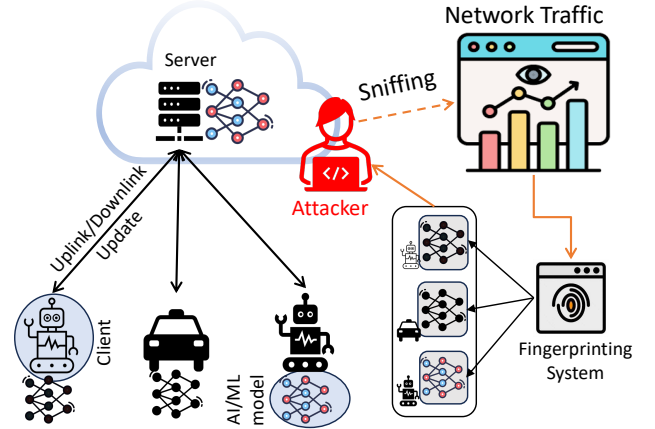


Figure 1: Framework for fingerprinting attack on FL.

[9]. In [20], researchers have shown that mobile activity can also be fingerprinted using control channel information in 5G networks.

While these studies focus on fingerprinting user activities and applications, a parallel research direction explores fingerprinting in distributed learning environments such as federated learning. The advancement of distributed machine learning has led to the integration of fingerprinting with FL. In [13], Melis et al. showed the capability of fingerprinting an FL client device using communication patterns, which reveals the client's identification. Song et al. investigated how attackers can utilize network information in FL systems to identify and monitor specific clients [17].

However, despite having tremendous progress in various domains, fingerprinting deep learning architectures in FL systems is a critical research gap. FL relies heavily on server-client communication while running different neural networks. Hence, network traffic may leak sensitive information about these models, and despite having potential security risks, no research has explicitly explored whether adversaries can fingerprint DL architectures. In this work, we explicitly ventured into this unexplored research area by analyzing network layer traffic generated by federated learning systems to fingerprint DL architectures.

3 METHODOLOGIES

In this section, we will discuss the experimental design, the configuration of our federated learning testbed, the procedure for collecting data, the feature engineering techniques, and the fingerprint approach that we utilize to conduct fingerprint attacks on federated learning systems.

3.1 Experimental Design

The primary goal of our research is to identify various DL architectures by evaluating network traffic data generated during federated learning. To accomplish this goal, we concentrate on the training phase of the FL framework, where multiple devices work together to train a global model in a distributed and secure manner.

A federated learning system with star topology generally contains two fundamental components: client and server. Each client individually trains its own local model, producing local updates. The updates are then sent to the central server, integrating all client

updates to train a global model. The improved global model parameters are later propagated to clients, thereby triggering another cycle of local training. This repeated exchange between the server and clients continues until convergence or a specified training criterion is achieved. Federated learning depends on the network connection, making data transmission between the server and clients an extremely vulnerable point. Despite the encryption of the model parameters, the network traffic inevitably reveals important meta-data, including packet sizes, transmission direction, and interarrival times. A skilled attacker with access to packet-sniffing tools can passively acquire this compromised information during the parameter update phase.

Figure 1 shows a high-level overview of our system model designed for fingerprinting DL networks in FL scenarios. In this approach, the adversary intercepts network packets transmitted during FL training sessions via conventional packet-sniffing methods. The collected network data is further processed and analyzed using our proposed fingerprinting system. After processing by the fingerprinting framework, an adversary can identify both the model architecture (e.g., CNN, RNN, Transformer) and the data modality type (e.g., text, audio, video) used for training. When the attacker obtains this information, then they can launch an effective, targeted attack on those models and modalities.

3.2 Threat Model

In this research, we consider a passive network adversary who can monitor communication between the FL clients and the server. The adversary does not actively interfere with the FL training process but aims to fingerprint deep learning architecture based on network-layer traffic patterns. We assume the adversary has the ability to (1) sniff layer-3 traffic using widely available packet capturing tools, like Wireshark; (2) extract basic packet metadata, including packet size and packet direction; and (3) observe the communication timing patterns based on the interarrival times. However, the adversary does not have access to the flow level or application layer data, such as payload information or FL model updates. Also, the adversary cannot access the higher-layer control information, such as the TCP sequence. Given those constraints, the adversary aims to identify distinct traffic signatures associated with the different DL architectures and infer the DL structure and its modality. This fingerprinting capability could enable further targeted adversarial attacks on the FL system by exploiting model-specific vulnerabilities.

3.3 Testbed Setup

In our experiment, we develop a basic federated learning testbed to assess our fingerprinting approach. To achieve this, we select two popular deep learning architectures, such as CNN and RNN, due to their fundamental differences in computational and data processing structures. CNN models are commonly used for spatial data (e.g., images and videos), while RNNs are used for sequential data with temporal features. These unique computational patterns of CNNs and RNNs introduce different training update patterns, which directly influence the communication pattern in federated learning.

Since this is a preliminary study to determine whether deep learning models are fingerprintable, we conducted our experiments

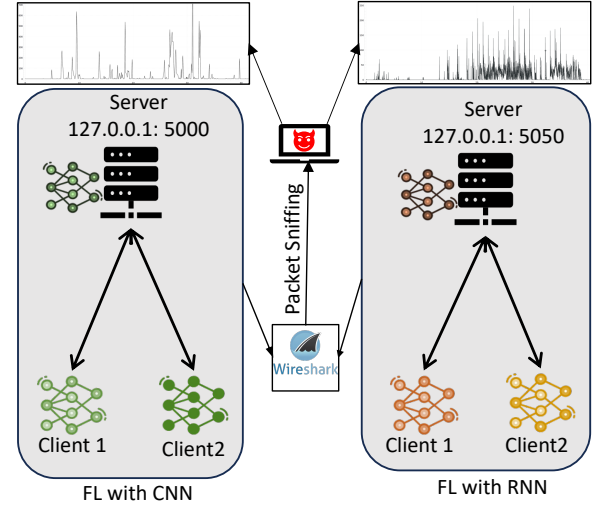


Figure 2: Testbed design for the fingerprint attack.

using a local host setup on a single machine instead of executing training over a wireless network. The primary reason for employing a local host setup instead of a proper wireless configuration is to provide a controlled federated learning environment. We reduce outside noise by isolating external network traffic and maintaining consistency between experiments. Therefore, it helps to accurately identify variations in traffic patterns in deep learning architectures rather than network circumstances or environmental influences. However, in the future, we plan to assess the feasibility of this threat under more realistic conditions.

We conduct these experiments on a PC with an NVIDIA RTX 3060 GPU to enable faster iteration, optimize troubleshooting, and allow efficient packet capturing. Local client instances are simulated via unique ports within the localhost environment (e.g., 127.0.0.1 with unique port numbers). Then, we collected network traffic using packet sniffing tools (e.g., Wireshark). Figure 2 illustrates the initial setup of our federated learning testbed and the data-collection methodology, focusing on the server-client layer-3 communication. The left section of this figure depicts the server aggregating updates from numerous clients executing the CNN model, while the right section exhibits clients using the RNN model. Furthermore, the upper region of the figure displays the type of packet collected by Wireshark.

3.4 Data Collection

To extensively assess the success of our fingerprinting method, we carefully collected network traffic data in controlled and realistic conditions. Data collection is conducted in two independent environments: an ideal setup with no external traffic and an alternative setting that includes external noise from internet browsing activities on clients. For each DL architecture, we use distinct datasets and tasks to ensure variational workloads.

We choose various custom CNN architectures by varying the number of layers and connections and train them on CIFAR-10 and Fashion MNIST image datasets. These databases vary in complexity, resolution, and image type. We also utilize the Sunspot Forecasting

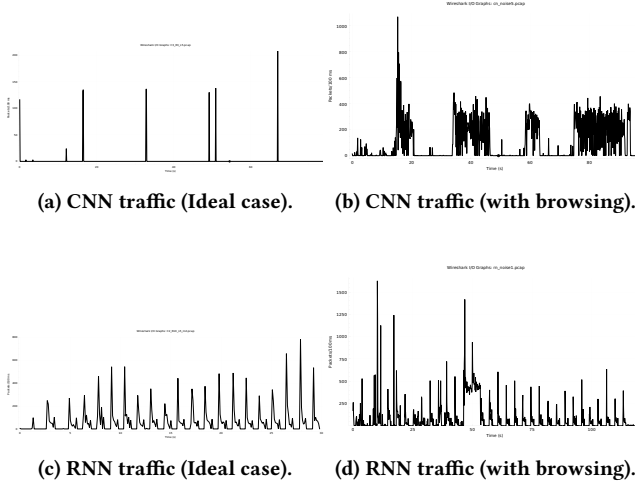


Figure 3: Network traffic patterns for CNN and RNN architectures under ideal and noisy experimental conditions.

time series dataset to train customized RNN models. The reason behind choosing these customized models is to identify the patterns that relate specifically to CNN or RNN.

In order to capture variations in the dataset, we systematically modified the federated learning hyperparameters, including local and global epoch numbers and training parameters. By modifying these parameters, we change the computation of the models, hence the variation in communication patterns, which aids in obtaining fingerprint robustness and generalization.

After that, we collect the data using Wireshark, where each traffic session is saved as an individual packet capture (pcap) file. We use these PCAP files for feature extraction and analysis and then utilize those features to fingerprint the DL models. For training, we use 8 CNN and 8 RNN-based traffic data (including both ideal and noisy data) and a total of 23 PCAP files (11 for RNN, 12 for CNN) for testing our approach.

In Figure 3, we show some collected data from both ideal and noisy cases to demonstrate how web browsing impacts network traffic.

3.5 Feature Engineering

In this section, we discuss the data processing and feature selection procedures of this work.

3.5.1 PCAP to CSV conversion. We use Wireshark to capture raw traffic data in a PCAP file. Then we use Tshark and Python to transform the PCAP into a structured CSV file. The CSV files contained important information, including timestamps, packet length, packet direction (uplink/downlink), and interarrival times.

3.5.2 Derived Features. To capture higher-level traffic behaviors, we computed statistical features such as packet-length statistics (mean, standard deviation, minimum, maximum, number of peaks), direction statistics (average direction, uplink/downlink proportions), interarrival-time statistics (mean, standard deviation, minimum, maximum, number of peaks), and overall traffic duration. However, we later discard the traffic duration because traffic duration depends on how long we are capturing the packets; it doesn't

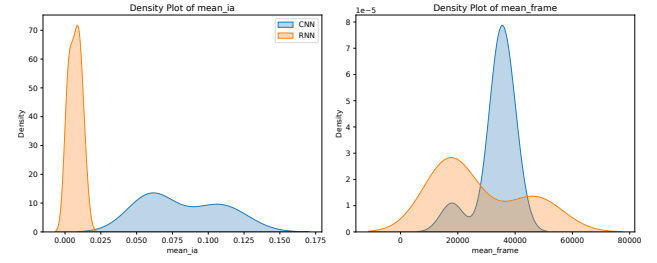


Figure 4: KL divergence plots of selected network traffic features (left: interarrival time; right: packet size) for CNN and RNN architectures.

have a direct impact on the FL traffic. We also measure the Kullback–Leibler (KL) divergence to compare CNN and RNN feature distributions. Figure 4 shows two distributions: mean interarrival time (mean_ia) on the left and mean packet size (mean_frame) on the right. Strong discriminative strength is indicated by a significant KL divergence, as seen in mean_ia, where the orange and blue curves hardly overlap. The left plot illustrates that CNN packets have more interarrival time variability, whereas RNN packets have a more concentrated distribution, indicating a more regular transmission pattern. Features having higher overlap, such as mean_frame, have lower KL values, indicating less successful class separation. The right plot demonstrates that CNN-generated packets are larger and more uniform than RNN-generated packets, which are more dispersed. These differences show that deep learning architectures have different network traffic patterns, facilitating network-layer fingerprinting.

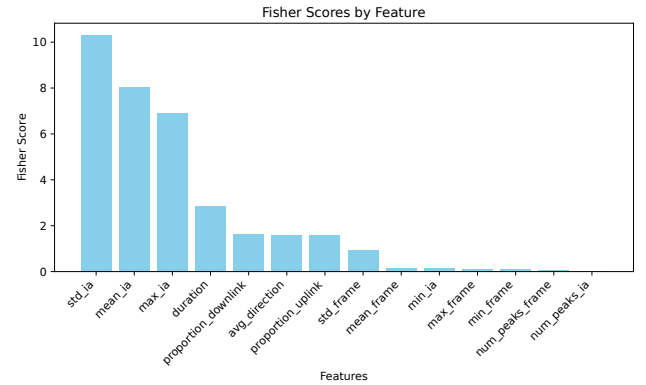


Figure 5: Fisher score analysis from the data features.

3.5.3 Feature Selection. Fisher score-based feature selection was employed to identify the most discriminative features. We use Fisher scores to determine the most important and discriminative features from the derived features. Fisher scores analyze the discriminative capability of features by evaluating the differences in means between classes relative to the variances within classes. Here, the classes are two: class A, which is the data from CNN training, and class B, which is from the RNN training. Figure 5 shows the computed Fisher scores for all candidate features. From here, I took the

highest-ranked features for the next step. This selection mainly helps to reduce dimensionality, minimize overfitting, and increase overall classification performance.

4 EVALUATION

In this section, we discuss the fingerprinting methods of DL models and the performance evaluation of our experiments.

4.1 Fingerprinting Approach

After processing the data, we selected the essential features using the methods discussed in the feature engineering section to perform the model fingerprinting. We employed well-established machine learning classifiers.

- **Random Forest (RF):** This is an ensemble method that contains multiple decision trees trained using bagging and random feature selection. It can handle non-linear data and give intrinsic feature importance, making it suitable for fingerprinting tasks [2].
- **Support Vector Machines (SVM):** SVM detects an optimal hyperplane that discriminates between classes by maximizing the margin between them; it performs better in binary classification tasks, even with a limited training sample, which makes it ideal for our fingerprinting task of identifying CNN and RNN-based tasks [4].
- **Gradient Boosting (XGBoost):** This is an ensemble learning technique that makes a robust classifier by iteratively combining weaker predictive models [3].

For each classifier, hyperparameter tuning was conducted using grid-search cross-validation on the training dataset. The performance of classifiers was evaluated on an independent test dataset, employing standard metrics such as accuracy, precision, recall, and F1-score to measure fingerprinting effectiveness comprehensively.

4.2 Evaluation Metrics

To evaluate the fingerprinting performance, we use the following widely accepted metrics:

- **Accuracy:** It measures the fraction of correctly predicted instances among all tested instances.
- **Precision:** It quantifies the proportion of true positive predictions among all predicted positives, thus reflecting false-positive sensitivity
- **Recall:** It measures how effectively a classifier identifies actual positive instances, capturing sensitivity to false negatives.
- **F1-score:** It combines precision and recall into a single metric, providing a balanced measure of classification effectiveness, particularly useful when classes are imbalanced or when precision and recall are both crucial.

These metrics are calculated and reported separately for each classifier (Random Forest, SVM, and XGBoost) and each class (CNN and RNN), providing comprehensive insight into their respective fingerprinting capabilities.

Table 1: Fingerprint performance for CNN and RNN architectures.

Method	Class	Precision	Recall	F1-score	Accuracy
Random Forest	CNN	1.00	1.00	1.00	100%
	RNN	1.00	1.00	1.00	
SVM	CNN	1.00	0.92	0.96	95.65%
	RNN	0.92	1.00	0.96	
XGBoost	CNN	0.92	1.00	0.96	95.65%
	RNN	1.00	0.91	0.95	

4.3 Result Analysis

Table 1 summarizes the fingerprinting performance using the well-established three machine learning models. Using these models, we distinguish between the CNN and RNN architectures that we utilize in federated learning. We modeled the fingerprinting methods as binary classification tasks, where each classifier was trained and evaluated using traffic data generated by training CNN and RNN models independently on the client devices in a simulated FL environment. Among the classifiers, Random Forest algorithms show performance by accurately classifying all test instances perfectly, achieving perfect scores across all other metrics.

The SVM classifier gained an overall average accuracy of 95.65%, misclassifying only one CNN instance as RNN. Regardless of a slight reduction in CNN recall to 92%, the classifier achieved perfect CNN precision at 100% and a high overall F1-score of 95.65%. The XGBoost classifier attained an accuracy of 95.65%, misclassifying a single RNN instance as CNN. This led to a CNN precision drop to 92% while maintaining perfect recall at 100% and a strong F1 score of 95.65%.

These results confirm the capability of accurately recognizing DL architectures using only network-layer traffic. The performance of Random Forest represents its robustness and suitability for this fingerprinting task, while slight errors observed in SVM and XGBoost suggest a potential sensitivity to noise or overlapping features.

5 DISCUSSION AND FUTURE DIRECTIONS

5.1 Discussion

Our study demonstrates that deep learning architectures can be effectively fingerprinted by analyzing network traffic patterns in federated learning environments. By leveraging network-layer meta-data and statistical traffic features, our approach successfully distinguishes between CNN and RNN architectures with high accuracy across different classifiers. This finding highlights a significant privacy vulnerability in federated learning, where an adversary can infer critical details about a deployed deep learning model through passive network observation.

While our results confirm the feasibility of such an attack, they also reveal several limitations. First, our experiments were conducted in a controlled environment where network conditions were ideal. Although we introduced some noise by browsing the web during data collection, real-world federated learning systems operate in far more complex and unpredictable conditions. Network packet loss, retransmissions, congestion, and multiplexed traffic from multiple sources may reduce the effectiveness of our attack.

Additionally, our fingerprinting approach assumes that only one type of deep learning model is being trained per traffic instance, whereas real-world federated learning environments often involve multiple architectures running simultaneously, making traffic patterns significantly more complex. Another limitation is the dataset size and model diversity. We tested our fingerprinting method on a small sample of deep learning architectures, whereas federated learning scenarios include a much wider variety of deep learning structures. This limited dataset likely led to overfitting to specific traffic patterns, reducing the generalizability of our fingerprinting method.

Despite these limitations, our results strongly indicate that deep learning architectures exhibit unique network traffic patterns, and with a more robust framework, it should be possible to identify different architectures even in realistic, complex network conditions.

5.2 Future Research Directions

To enhance the practicality and robustness of our fingerprinting approach, several key research directions should be pursued:

- **Fingerprinting in Real-World Network Conditions:** Future research should examine how packet loss, congestion, encryption, and background traffic affect fingerprinting accuracy. To manage multiplexed traffic, where many deep learning models train on the same network, models should be refined.
- **Expanding Scope of Attack and Defenses:** Modern federated learning deployments incorporate transformers (BERT, ViTs), RLs, and GANs. Checking if these designs have unique network traffic patterns is crucial. This new risk in federated learning must be mitigated by creating effective countermeasures such as network-layer traffic concealment, adversarial perturbations, and safe aggregation.
- **Building a Scalable and Generalized Framework:** Future research should train fingerprint models on vast, diversified datasets with different architectures, network conditions, and traffic patterns to improve generalizability. Fine-tuned classification models that adapt to real-world federated learning settings and ensure reliable fingerprinting across deployment scenarios will result.

By addressing these areas, future work can enhance attack effectiveness while strengthening security defenses, ensuring greater privacy and robustness in federated learning systems.

6 CONCLUSIONS

We proposed a novel method to fingerprint deep learning architectures by analyzing network traffic patterns. Our approach consists of a multi-step procedure, including traffic preprocessing, feature engineering, and classification using machine learning techniques. While this study presents an initial exploration, our findings highlight a previously overlooked privacy threat in federated learning environments.

This vulnerability is particularly concerning for safety-critical applications where deep learning models operate as core components, such as autonomous driving and healthcare systems. Our experimental results demonstrate that DL models are susceptible

to fingerprinting attacks through encrypted network leakage information, emphasizing the need for both stronger defenses and more robust attack strategies to assess vulnerabilities further. Future research should focus on increasing attack resilience in more diverse network conditions and developing countermeasures to mitigate this security risk effectively.

REFERENCES

- [1] Abbas Acar, Hossein Fereidooni, Tigist Abera, Amit Kumar Sikder, Markus Miettinen, Hidayet Aksu, Mauro Conti, Ahmad-Reza Sadeghi, and Selcuk Uluagac. 2020. Peek-a-boo: I see your smart home activities, even encrypted!. In *Proceedings of the 13th ACM Conference on Security and Privacy in Wireless and Mobile Networks*. 207–218.
- [2] Leo Breiman. 2001. Random forests. *Machine learning* 45 (2001), 5–32.
- [3] Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*. 785–794.
- [4] Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning* 20 (1995), 273–297.
- [5] Zhimin He, Jie Yin, Yu Wang, Guan Gui, Bamidele Adebisi, Tomoaki Ohtsuki, Haris Gacanin, and Hikmet Sari. 2021. Edge device identification based on federated learning and network traffic feature engineering. *IEEE Transactions on Cognitive Communications and Networking* 8, 4 (2021), 1898–1909.
- [6] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Ben-Nis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. 2021. Advances and open problems in federated learning. *Foundations and trends® in machine learning* 14, 1–2 (2021), 1–210.
- [7] Vishal Kaushal and Sangeeta Sharma. 2025. Securing the collective intelligence: a comprehensive review of federated learning security attacks and defensive strategies. *Knowledge and Information Systems* (2025), 1–39.
- [8] Pierre Laperdrix, Walter Rudametkin, and Benoit Baudry. 2016. Beauty and the beast: Diverting modern web browsers to build unique browser fingerprints. In *2016 IEEE Symposium on Security and Privacy (SP)*. IEEE, 878–894.
- [9] Jianfeng Li, Zheng Lin, Jian Qu, Shuohan Wu, Hao Zhou, Yangyang Liu, Xiaobo Ma, Ting Wang, Xiapu Luo, and Xiaohong Guan. 2024. Robust App Fingerprinting Over the Air. *IEEE/ACM Transactions on Networking* (2024).
- [10] Xiaobo Ma, Jian Qu, Jianfeng Li, John CS Lui, Zhenhua Li, and Xiaohong Guan. 2020. Pinpointing hidden IoT devices via spatial-temporal traffic fingerprinting. In *IEEE INFOCOM 2020-IEEE conference on computer communications*. IEEE, 894–903.
- [11] Xiaobo Ma, Jian Qu, Jianfeng Li, John CS Lui, Zhenhua Li, Wenmao Liu, and Xiaohong Guan. 2021. Inferring hidden IoT devices and user interactions via spatial-temporal traffic fingerprinting. *IEEE/ACM Transactions on Networking* 30, 1 (2021), 394–408.
- [12] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*. PMLR, 1273–1282.
- [13] Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. 2019. Exploiting unintended feature leakage in collaborative learning. In *2019 IEEE symposium on security and privacy (SP)*. IEEE, 691–706.
- [14] Andriy Panchenko, Fabian Lanze, Jan Pennekamp, Thomas Engel, Andreas Zinnen, Martin Henze, and Klaus Wehrle. 2016. Website Fingerprinting at Internet Scale. In *NDSS*, Vol. 1. 23477.
- [15] Mohammad Saidur Rahman, Payap Sirinam, Nate Mathews, Kantha Girish Gangadhara, and Matthew Wright. 2019. Tik-tok: The utility of packet timing in website fingerprinting attacks. *arXiv preprint arXiv:1902.06421* (2019).
- [16] Chuan Sheng, Wei Zhou, Qing-Long Han, Wanlun Ma, Xiaogang Zhu, Sheng Wen, and Yang Xiang. 2025. Network Traffic Fingerprinting for IIoT Device Identification: A Survey. *IEEE Transactions on Industrial Informatics* (2025).
- [17] Mengkai Song, Zhibo Wang, Zhifei Zhang, Yang Song, Qian Wang, Ju Ren, and Hairong Qi. 2020. Analyzing user-level privacy attack against federated learning. *IEEE Journal on Selected Areas in Communications* 38, 10 (2020), 2430–2444.
- [18] Vincent F Taylor, Riccardo Spolaor, Mauro Conti, and Ivan Martinovic. 2017. Robust smartphone app identification via encrypted network traffic analysis. *IEEE Transactions on Information Forensics and Security* 13, 1 (2017), 63–78.
- [19] Tao Wang. 2020. High precision open-world website fingerprinting. In *2020 IEEE Symposium on Security and Privacy (SP)*. IEEE, 152–167.
- [20] Gunwoo Yoon and Byeongdo Hong. 2024. Scalable and Robust Mobile Activity Fingerprinting via Over-the-Air Control Channel in 5G Networks. *arXiv preprint arXiv:2409.12572* (2024).
- [21] Yifei Zhang, Dun Zeng, Jinglong Luo, Xinyu Fu, Guanzhong Chen, Zenglin Xu, and Irwin King. 2024. A survey of trustworthy federated learning: Issues, solutions, and challenges. *ACM Transactions on Intelligent Systems and Technology* 15, 6 (2024), 1–47.

- [22] Yiwen Zhang and Weilin Zeng. 2024. Local adversarial attack of time series forecasting based on mutual information. In *Third International Conference on Machine Vision, Automatic Identification, and Detection (MVAID 2024)*, Vol. 13230. SPIE, 467–473.
- [23] Man Zhou, Wenyu Zhou, Jie Huang, Junhui Yang, Minxin Du, and Qi Li. 2024. Stealthy and effective physical adversarial attacks in autonomous driving. *IEEE Transactions on Information Forensics and Security* (2024).