

# Genetic Algorithm-Based Dynamic Backdoor Attack on Federated Learning-Based Network Traffic Classification

Mahmoud Nazzal\*, Nura Aljaafari<sup>§</sup>, Ahmed Sawalmeh<sup>||</sup>, Abdallah Khreishah\*,  
Muhammad Anan<sup>||</sup>, Abdulelah Algosaibi<sup>§</sup>, Mohammed Alnaeem<sup>§</sup>, Adel Aldalbahi<sup>§</sup>, Abdulaziz Alhumam<sup>§</sup>,  
Conrado P. Vizcarra<sup>§</sup>, and Shadan Alhamed<sup>§</sup>

\* New Jersey Institute of Technology, Newark, NJ 07102, USA

<sup>§</sup> King Faisal University, Al Hofuf 31982, KSA

<sup>||</sup> Alfaisal University, Riyadh 11533, KSA

E-mails: mahmoud.nazzal@ieee.org, naaljaafari@kfu.edu.sa, asawalmeh@alfaisal.edu, abdallah@njit.edu  
manan@alfaisal.edu, {aalgosaibi, naeem, aaldalbahi, aahumam, cvizcarra, ssalhamed}@kfu.edu.sa

**Abstract**—Federated learning enables multiple clients to collaboratively contribute to the learning of a global model orchestrated by a central server. This learning scheme promotes clients’ data privacy and requires reduced communication overheads. In an application like network traffic classification, this helps hide the network vulnerabilities and weakness points. However, federated learning is susceptible to backdoor attacks, in which adversaries inject manipulated model updates into the global model. These updates inject a salient functionality in the global model that can be launched with specific input patterns. Nonetheless, the vulnerability of network traffic classification models based on federated learning to these attacks remains unexplored. In this paper, we propose GABAttack, a novel genetic algorithm-based backdoor attack against federated learning for network traffic classification. GABAttack utilizes a genetic algorithm to optimize the values and locations of backdoor trigger patterns, ensuring a better fit with the input and the model. This input-tailored dynamic attack is promising for improved attack evasiveness while being effective. Extensive experiments conducted over real-world network datasets validate the success of the proposed GABAttack in various situations while maintaining almost invisible activity. This research serves as an alarming call for network security experts and practitioners to develop robust defense measures against such attacks.

**Index Terms**—Backdoor attack, trigger design, federated learning, network traffic classification, genetic algorithm.

## I. INTRODUCTION

Network traffic classification (NTC) categorizes network traffic observations based on their features and characteristics to infer certain properties. NTC is an integral part of network management commonly employed by network administrators and service providers. A direct outcome of NTC is giving an insight into the types of activities, applications, and protocols used network-wide. This information is essential for optimizing the network resources and identifying any potential threats or malicious actions. Classical NTC approaches include port-based classification, deep packet inspection [1], and statistical classification [2]. Similar to the case with many other application areas, machine learning (ML) models have

been successfully utilized in NTC due to their outstanding performances as data-driven approaches [3].

Conventional ML models used in NTC share a common limitation; depending on manually crafted features that typically require experts’ knowledge, practice, and time. As an example, [4] utilizes support vector machines (SVM) operating on 250 network flow features proposed in [5]. A common drawback of these approaches is also requiring a reasonable size of training data to train the model.

To resolve the limitations of standard ML models for NTC, recent research considers a growing interest in federated learning (FL) [6], [7] as a framework for model training. FL aggregates user-end ML model contributions to obtain a global model. Since FL allows local training on the client side, it achieves two main advantages; promoting the client’s privacy, and saving the network bandwidth as only model coefficients are communicated [8]. These are attractive features for an application like NTC since NTC data can easily reveal network vulnerabilities and weaknesses. Accordingly, several works have recently considered FL for NTC [9]–[13].

Despite the attractive advantages of FL, it is widely believed to be inherently vulnerable to adversarial and poisoning attacks [8] similar to the case of virtually all ML settings due to their data dependency. In a security-critical application like NTC, malicious actors have strong incentives to attack NTC models to tweak their cyber attacks [14]. A key incentive is to limit their cyber attacks to be within legitimate network traffic [15] thereby bypassing NTC-based network intrusion detection and firewalls. This highlights the importance of understanding what creates this vulnerability and developing efficient countermeasures accordingly.

Even though FL keeps data at clients, FL NTC models can enable multiple clients to access model parameters and thus empower malicious intervention. While recent literature has a few works on establishing the vulnerability of NTC models to adversarial attacks, their susceptibility to backdoor attacks is

yet studied. Particularly, the vulnerability of the recent FL-based NTC approach is not addressed. A backdoor injects a salient functionality into a target model. This functionality is only activated if a certain trigger pattern is available in a test input. To keep the attack evasive, the target model should behave normally with benign inputs [16]. While various existing backdoor attacks perform relatively well on effectiveness, evasiveness is a more challenging goal [17], [18].

**Contributions** Motivated by the above discussion; we present the following contributions.

- Establishing the vulnerability of FL-based NTC models to backdoor attacks characterized by specially crafted triggers injected in training and activated in inference.
- GABAttack: a new algorithm for backdoor attack against FL optimizing the locations and values of trigger patterns in network traffic data based on genetic algorithm optimization. GABAttack offers attack transferability and produces input-tailored dynamic triggers enhancing attack evasiveness.
- A comprehensive set of experiments to investigate the performance of GABAttack in real-world NTC attack in terms of effectiveness evasiveness.

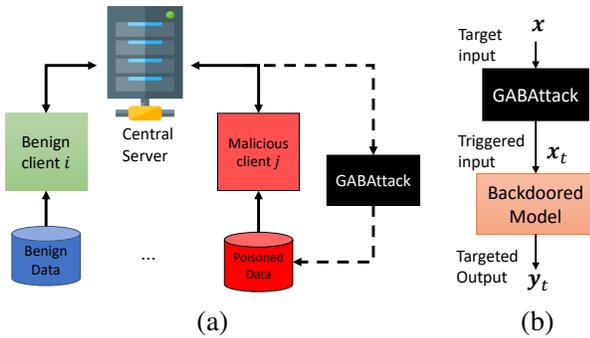


Fig. 1. GABAttack: BD pattern injection in (a) where the adversary modifies the training data of compromised clients to inject a backdoor functionality in the model updates submitted to the server. Attack launching in (b).

## II. BACKGROUND

**FL-based NTC FL** [19] is a setting for distributed ML where  $n$  clients collaboratively learn a global model  $\phi$ . In a training round  $t \in \{1, \dots, T\}$ , each client  $i \in \{1, \dots, k\}$  trains a model  $\phi_i$  on its local data  $D_i$  with based on the previous global model  $\phi_{t-1}$ . The client model (or its update over the previous global model) is then communicated to the central server. This server aggregates all client model updates to obtain an updated global model, and the process is repeated in the next FL round. Since FL is based on communicating model coefficient updates rather than data points, it naturally promotes data privacy while substantially reducing the communication overhead compared to standard distributed learning. Thus, FL has gained popularity in a wide spectrum of applications ranging from healthcare [20] to autonomous driving [21] and, more recently, NTC [9]–[13]. In an NTC context, FL is used for data privacy reasons [9]. FL is enhanced with an attention mechanism for better client model

aggregation in [22]. Other works [10], [23] use FL-based NTC for device identification. These works assume honest clients and overlook the possibility of having some clients potentially compromised by malicious actors as depicted in Fig. 1.

**Genetic algorithm (GA)** [24] is a meta-heuristic search algorithm that has been widely used in many application areas. The idea behind GA is based on alternating between the generation of new candidate solutions and selecting the best among them. These operations are inspired by the processes of evolution and natural selection in evolutionary theory. The deployment of GA requires first encoding data (either inputs or solutions) into *chromosomes*. A chromosome represents the information of a given solution. GA then starts from an arbitrary set of initial candidate solutions referred to as the parents. Then, the parent set is iteratively refined by yielding new solutions referred to as the offspring. Fundamentally, GA uses the processes of mutation and crossover to obtain offspring from parents. Next, in each iteration, the set of parents is updated by the inclusion of specific members of the offspring, replacing specific elements in the parent set. This is done based on a certain fitness function. For the sake of diversity and exploration, GA balances between greedily selecting the best candidates and keeping a few “bad” ones.

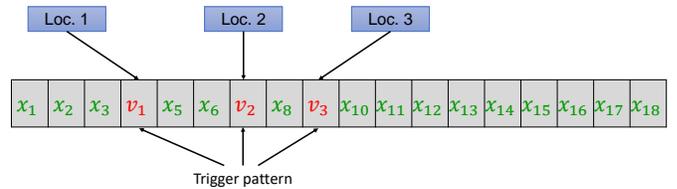


Fig. 2. Chromosome construction and encoding, where  $x_1$  through  $x_{18}$  are the elements of the target feature input and  $v_1$  through  $v_3$  are trigger values at locations *Loc1* through *Loc3*.

## III. THREAT MODEL

We consider an FL-based NTC system [9]–[13] composed of  $N$  clients orchestrated by a server as shown in Fig. 1(a). FL is operated in rounds where in each round  $t$ , the server randomly selects  $k = C.N$ ,  $C \leq 1$  clients to participate in the training process. We characterize the threat model in terms of the adversary’s objectives, knowledge, and capabilities. The objective of the adversary is to craft *effective* and *evasive* (salient) backdoor attacks on the target model. The adversary is assumed to possess a few of the FL clients and to use them to inject the backdoor during training and launch it during testing times. So, the adversary controls the training data and training operations for its clients only. Specifically, the adversary can manipulate the training data points and their labels at these clients as shown in Fig. 1(a). Besides, it knows the initial model broadcast to them by the server. For launching the attack as shown in Fig. 1(b), the adversary needs to know what network data to target. We assume that the adversary can observe legitimate packets by “sniffing the network” and therefore record the feature values of legitimate traffic as

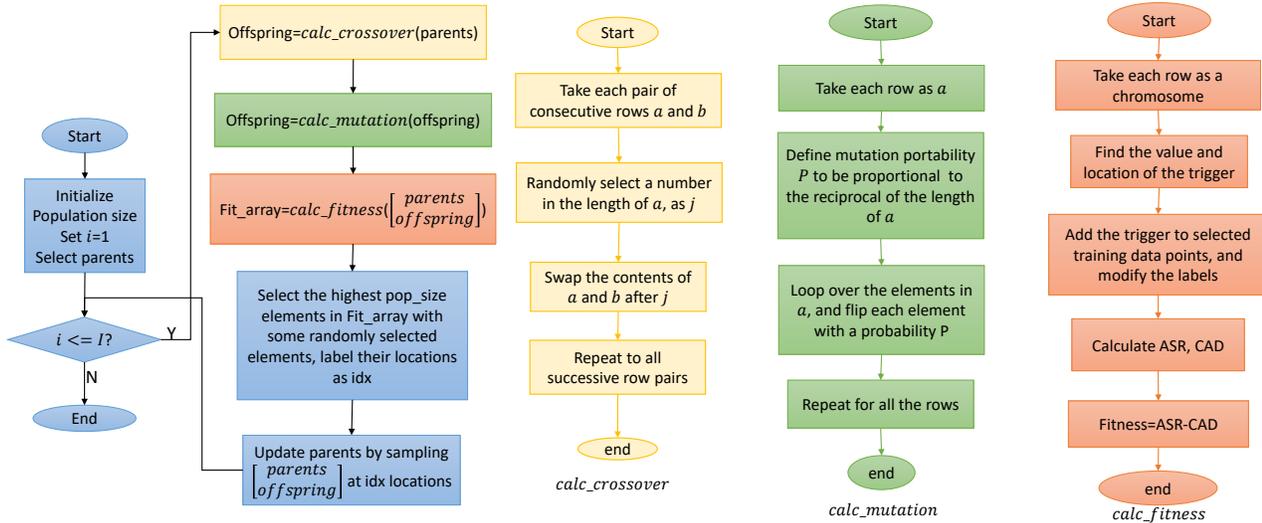


Fig. 3. GABAttack in flowchart description.

commonly assumed in adversarial attack works against NTC models [14], [25].

#### IV. THE PROPOSED GABATTACK

Fig. 1 illustrates the workflow of the GABAttack. First, GABAttack is used to inject a backdoor functionality in the target model during training (part a). Then, it is used to launch an attack during the run time of the model (part b). The goals GABAttack aims at when generating backdoor patterns are effectiveness and evasiveness, as represented below.

$$\phi(x) = \begin{cases} \phi(x_t) = y_t & \text{effectiveness} \\ \phi(x) = \phi_0(x) & \text{evasiveness} \end{cases}$$

where  $\phi$  is a backdoored model corresponding to a benign model  $\phi_0$ ,  $x$  is a benign input,  $x_t$  is a triggered input, and  $y_t$  is the targeted label of  $x_t$ . According, the backdoor attack problem can be formulated as follows.

$$\arg \min_{\phi^*, D_m} \mathcal{L}_M(D_m; \phi^*) + \rho \|\Delta_m^{t+1} - \bar{\Delta}_{\text{ben}}^t\|_2 \quad (1)$$

where  $D_m$  denotes the data of a malicious client,  $\rho$  is a hyperparameter,  $\mathcal{L}_M$  is Cross-entropy loss for main task, and  $\bar{\Delta}_{\text{ben}}^t$  denotes benign clients' averaged model updates. In this paper,  $\bar{\Delta}_{\text{ben}}^t$  is estimated using a shared global model. We assume that the aggregated global model is similar to the benign client's local model as the shared global model converges to a point with a high testing accuracy. It should be noted that the addition of a regularization term is not sufficient to ensure that the malicious weight update is close to that of the benign agents since there could be multiple local minima with similar loss values.

GABAttack employs a GA-based approach for optimizing both the values of the added trigger elements and their placement in the data. This requires first developing a chromosome encoding of the GA. Fig. 2 shows the proposed chromosome encoding process along with the backdoor pattern placement.

A direct measure of an attack's effectiveness is the attack success rate (ASR) defines as follows.

$$\text{ASR} = \frac{I(\phi(x_t) = y_t)}{\text{num}} \quad (2)$$

where  $x_t$  is a triggered input,  $y_t$  is the targeted model outcome,  $\text{num}$  is the total number of attacked inputs, and  $I$  is an indicator function.

As for evasiveness, one can quantify it as the drop in model accuracy working on benign data after having a trigger functionality. Thus, it can be calculated as

$$\text{CAD} = \text{acc}(\phi_0(X), Y) - \text{acc}(\phi(X), Y) \quad (3)$$

where  $\text{acc}$  is an accuracy function,  $X$  is a set of benign training data with  $Y$  true labels,  $\phi_0$  is a benign trained model, and  $\phi$  is its backdoored version. Accordingly, we define the following GA fitness function to incorporate ASR and CAD.

$$f = \text{ASR} - \gamma \text{CAD} \quad (4)$$

where  $\gamma$  balances the trade-off between ASR and CAD.

The main steps of GABAttack are outlined in the flowchart of Fig. 3. Along with maintaining attack effectiveness and evasiveness, there are several advantages of GA in optimizing the triggers. First, GA is naturally a global optimization algorithm that is likely to generate globally optimized outcomes. Second, it allows searching over the set of possible candidates without the need for establishing this set ahead of time.

#### V. EXPERIMENTS

We examine the performance of GABAttack first in a centralized setting assuming white-box model access and then in the intended usage in an FL setting with the same datasets and models.

### A. Experimental setup

We consider the widely used Moore [5] dataset in all the experiments. This dataset uses 216 features used in training and inference. As for the classes, the top six classes are assumed; *WWW*, *MAIL*, *FTP-DATA*, *FTP-CONTROL*, *DATABASE*, *SERVICES*, and *ATTACK*. As for the simulation platform, experiments are run over Google’s Colab with Keras 2.8.0 with Tensorflow Backend, and Python 3.7.13.

For the centralized setting, we assume the poisoning occurs on the data the model directly uses for training, where it is combined with clean data with the correct label. The model is trained for 3 epochs, with a batch size of 100 and a learning rate of 0.001. For the FL experiments, the data is distributed across the clients following an IID nature. In each FL round, 10 clients are assumed to be participating in the model training. The benign clients are assumed to follow the traditional procedure of training the model and that they are training the model over 3 epochs with a batch size of 100 data points. Malicious clients perform two operations; first, they run GABAttack on their whole training data to get the best trigger values and locations. Then only 50% at maximum is cut from their training data and poisoned with the trigger. The second operation is model training, in which the malicious clients combine the clean data and the poisoned data and then use the combined data to train their models. Following similar training parameters to the benign clients, it is assumed that each malicious client trains its model for 3 epochs, with a batch size of 100, and a learning rate of 0.001.

The model accuracy and ASR are recorded to measure the attack’s success. In the FL setting, it is reported for malicious clients, and the same metrics are reported for the global model. Furthermore, the global model accuracy on the benign data is recorded at the end of each FL round. However, since this is a dynamic attack, where each malicious client has its trigger and poisons the data in a different location, this presents a challenge of how to measure the attack success on the global model. Thus, the attack success in the FL is measured in the global model for each trigger in an attack round.

**Classification models** We use the following three models. The first model, shown in Fig. 4(a) is a neural network with one long short-term memory (LSTM) layer of 100 hidden units, following a fully connected layer. The second model is a neural network that contains a layer of 1D convolution layer, then an LSTM layer with 100 units, followed by a fully connected layer. The structure of *Convolutional<sub>LSTM</sub>*

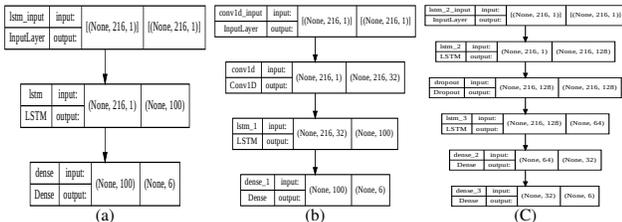


Fig. 4. The architectures of the *Simple<sub>LSTM</sub>*, *Convolutional<sub>LSTM</sub>*, and *Complex<sub>LSTM</sub>* models in (a), (b), and (c), respectively.

is shown in Fig. 4(b). The third model as shown in Fig. 4(c), is a neural network that has an LSTM layer with 100 units, then another LSTM layer with 64 units. The model includes 2 fully-connected layers. The three models are created incrementally in terms of the depth and the number of parameters in each model, i.e., their complexity.

### B. Experiments

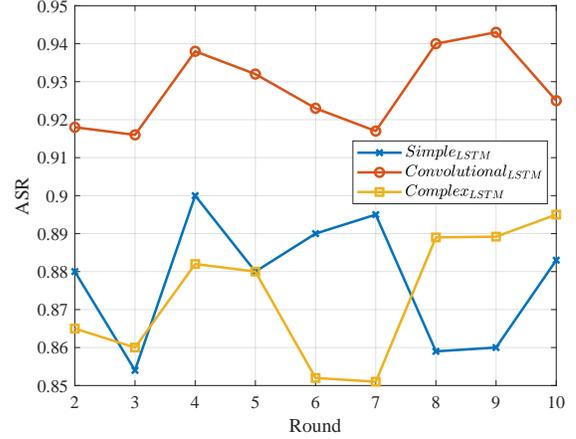


Fig. 5. Attack performance when the percentage of adversary clients is 20%.

The first point that is investigated is the attack success over different models, that is, testing if the attack can succeed in deep and complex models and the attack transferability to complex models. It is worth noting that several experiments are conducted using traditional ML, such as SVM, and decision trees, but we do not present their results as they perform poorly on the main classification task, and are thus of low practical interest. For this experiment, in the FL setting, it is assumed that 10% of the clients participating in the pool are malicious and the attack is initiated on the fourth round. It is also assumed that adversary clients run independently.

Table. I demonstrates the centralized model accuracy on the aforementioned models as well as the ASR. For the training setting here, it is assumed that the model is trained from scratch. It is seen that the attack succeeds in the cases without compromising the model accuracy on the benign data. It is noted that the models considered show better results on the data than most models presented in the literature.

TABLE I  
CENTRALIZED MODEL RESULTS.

	<i>Simple<sub>LSTM</sub></i>	<i>Convolutional<sub>LSTM</sub></i>	<i>Complex<sub>LSTM</sub></i>
Accuracy	93.16%	95.75%	95.59%
ASR	99.9%	100.0%	100.0%

Next, we examine the attack performance metrics in the FL setting. Fig 5 shows the ASR versus FL round with the three models considered. It can be seen that the ASR success is communcerate with model’s compelexity. This can be interprestred in view of the fact that model complexity

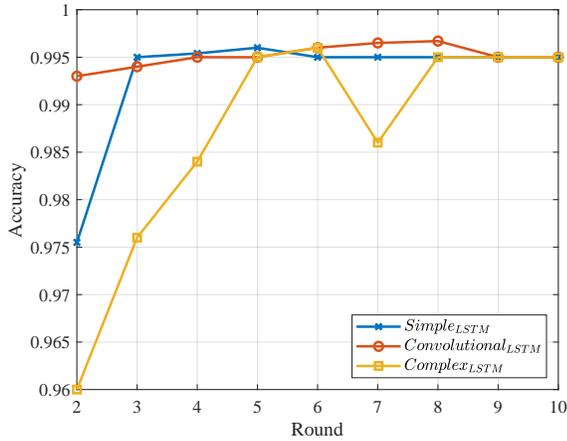


Fig. 6. Attack performance when the percentage of adversary clients is 50%.

opromises for accomodating for both the main task and the backdoor functionality. Moreover, the attack is generally highly successful for the three models considered. As for attack evasiveness, Fig. 6 shows the global model accuracy with the three models. It can be seen that the most complex model exhibits less model accuracy in the first initial rounds. However, the model accuracy is maintained for the three models. This result assures the evasiveness of the attack.

As for the centralized setting, the ASR is at its highest on the attack round, which is around 99.98% on average, but it slightly decreases on the following rounds, reporting 98.73% and 98.20%, respectively. The results presented in this section show that the GABAttack’s attack is successful in terms of both attack’s effectiveness and evasiveness.

## VI. RELATED WORK

**Backdoor Attack on FL Models** In backdoor attacks, the adversary manipulates the local models of compromised clients to obtain poisoned models to be then aggregated into the global model. There are many works of backdoor attacks on FL. Several recent works demonstrate the susceptibility of the FL model to backdoor attacks. These works assume that some FL clients are compromised by the adversary and under its control [26], [27]. Some works concentrate on scaling the impact of malicious clients’ contributions to the global model in what is known as model replacement attacks [28]. Other works focus on keeping the attack as stealthy as possible [29].

**Adversarial attack on NTC models** is still in its early stages. The current research body in this area focuses mainly on developing adversarial examples to manipulate NTC outcomes. Along this line, [30] employs the well-known Carlini and Wanger method [31] to generate such adversarial examples. Next, [32] investigates a range of established adversarial attacks in targeting NTC models thus demonstrating their vulnerability. Subsequently, [33] leverages mutual information to identify the optimal features to perturb to manipulate NTC classification. It can be seen that this research body focuses on interference-time evasive attacks and overlooks backdoor

attacks. Besides, to the best of our knowledge, there are no attacks on FL-based NTC.

## VII. CONCLUSION

In this paper, we propose GABAttack, a new backdoor attack against FL-based NTC models. The proposed GABAttack uses a genetic algorithm as a means for tuning the values and locations of backdoor trigger patterns to best fit the input and the model if known. Thus, this is an input-tailored dynamic attack promising improved attack evasiveness. Extensive experiments conducted over real-world NTC data with varying model complexities validate the success of the proposed GABAttack in terms of attack effectiveness and evasiveness. This work establishes the vulnerability of FL-based NTC models to backdoor attacks and calls for devising viable defense measures against such attacks. Future work will include the design of a coordinated attack across adversary-compromised clients and better ways of balancing the effectiveness-evasiveness trade-offs in the attack.

## ACKNOWLEDGEMENT

The authors extend their appreciation to the Deputyship for Research & Innovation, Ministry of Education in Saudi Arabia for funding this research work through the project number 1120.

## REFERENCES

- [1] W. Song, M. Beshley, K. Przystupa, H. Beshley, O. Kochan, A. Pryslupskyi, D. Pieniak, and J. Su, “A software deep packet inspection system for network traffic analysis and anomaly detection,” *Sensors*, vol. 20, no. 6, p. 1637, 2020.
- [2] C.-T. Su and J.-H. Hsu, “An extended chi2 algorithm for discretization of real value attributes,” *IEEE transactions on knowledge and data engineering*, vol. 17, no. 3, pp. 437–441, 2005.
- [3] R. Alshammari and A. N. Zincir-Heywood, “Identification of voip encrypted traffic using a machine learning approach,” *Journal of King Saud University-Computer and Information Sciences*, vol. 27, no. 1, pp. 77–92, 2015.
- [4] W. Zhongsheng, W. Jianguo, Y. Sen, and G. Jiaqiong, “Retracted: Traffic identification and traffic analysis based on support vector machine,” *Concurrency and Computation: Practice and Experience*, vol. 32, no. 2, p. e5292, 2020.
- [5] A. Moore, D. Zuev, and M. Crogan, *Discriminators for use in flow-based classification*, 2013.
- [6] V. Mothukuri, R. M. Parizi, S. Pouriyeh, Y. Huang, A. Dehghantanha, and G. Srivastava, “A survey on security and privacy of federated learning,” *Future Generation Computer Systems*, vol. 115, pp. 619–640, 2021.
- [7] S. Banabilah, M. Aloqaily, E. Alsayed, N. Malik, and Y. Jararweh, “Federated learning review: Fundamentals, enabling technologies, and future applications,” *Information processing & management*, vol. 59, no. 6, p. 103061, 2022.
- [8] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, “Advances and open problems in federated learning,” *Foundations and Trends® in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.
- [9] H. Mun and Y. Lee, “Internet traffic classification with federated learning,” *Electronics*, vol. 10, no. 1, p. 27, 2020.
- [10] Z. He, J. Yin, Y. Wang, G. Gui, B. Adebisi, T. Ohtsuki, H. Gacanin, and H. Sari, “Edge device identification based on federated learning and network traffic feature engineering,” *IEEE Transactions on Cognitive Communications and Networking*, vol. 8, no. 4, pp. 1898–1909, 2021.
- [11] Y. Peng, M. He, and Y. Wang, “A federated semi-supervised learning approach for network traffic classification,” *arXiv preprint arXiv:2107.03933*, 2021.

- [12] Y. Guo and D. Wang, "Feat: A federated approach for privacy-preserving network traffic classification in heterogeneous environments," *IEEE Internet of Things Journal*, vol. 10, no. 2, pp. 1274–1285, 2022.
- [13] C. Sun, B. Chen, Y. Bu, and D. Zhang, "Traffic classification method based on federated semi-supervised learning," in *Proceedings of the 2022 6th International Conference on Electronic Information Technology and Computer Engineering*, 2022, pp. 875–882.
- [14] M. E. Ahmed and H. Kim, "Poster: Adversarial examples for classifiers in high-dimensional network data," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 2017, pp. 2467–2469.
- [15] M. Usama, J. Qadir, A. Al-Fuqaha, and M. Hamdi, "The adversarial machine learning conundrum: can the insecurity of ml become the achilles' heel of cognitive networks?" *IEEE Network*, vol. 34, no. 1, pp. 196–203, 2019.
- [16] R. Pang, H. Shen, X. Zhang, S. Ji, Y. Vorobeychik, X. Luo, A. Liu, and T. Wang, "A tale of evil twins: Adversarial inputs versus poisoned models," in *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, 2020, pp. 85–99.
- [17] B. G. Doan, E. Abbasnejad, and D. C. Ranasinghe, "Februus: Input purification defense against trojan attacks on deep neural network systems," in *Annual Computer Security Applications Conference*, 2020, pp. 897–912.
- [18] Y. Liu, W.-C. Lee, G. Tao, S. Ma, Y. Aafer, and X. Zhang, "Abs: Scanning neural networks for back-doors by artificial brain stimulation," in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 2019, pp. 1265–1282.
- [19] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [20] Y. Chen, X. Qin, J. Wang, C. Yu, and W. Gao, "Fedhealth: A federated transfer learning framework for wearable healthcare," *IEEE Intelligent Systems*, vol. 35, no. 4, pp. 83–93, 2020.
- [21] X. Liang, Y. Liu, T. Chen, M. Liu, and Q. Yang, "Federated transfer reinforcement learning for autonomous driving," *arXiv preprint arXiv:1910.06001*, 2019.
- [22] M.-y. Zhu, Z. Chen, K.-f. Chen, N. Lv, and Y. Zhong, "Attention-based federated incremental learning for traffic classification in the internet of things," *Computer Communications*, vol. 185, pp. 168–175, 2022.
- [23] P. M. Sánchez Sánchez, A. Huertas Celdrán, J. R. Buendía Rubio, G. Bovet, and G. Martínez Pérez, "Robust federated learning for execution time-based device model identification under label-flipping attack," *Cluster Computing*, pp. 1–12, 2023.
- [24] W. Banzhaf, P. Nordin, R. E. Keller, and F. D. Francone, *Genetic programming: an introduction: on the automatic evolution of computer programs and its applications*. Morgan Kaufmann Publishers Inc., 1998.
- [25] G. Wang, T. Wang, H. Zheng, and B. Y. Zhao, "Man vs. machine: Practical adversarial detection of malicious crowdsourcing workers," in *23rd {USENIX} security symposium ({USENIX} security 14)*, 2014, pp. 239–254.
- [26] S. Shen, S. Tople, and P. Saxena, "Auror: Defending against poisoning attacks in collaborative deep learning systems," in *Proceedings of the 32nd Annual Conference on Computer Security Applications*, 2016, pp. 508–519.
- [27] C. Xie, K. Huang, P.-Y. Chen, and B. Li, "Dba: Distributed backdoor attacks against federated learning," in *International conference on learning representations*, 2020.
- [28] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2938–2948.
- [29] H. Wang, K. Sreenivasan, S. Rajput, H. Vishwakarma, S. Agarwal, J.-y. Sohn, K. Lee, and D. Papailiopoulos, "Attack of the tails: Yes, you really can backdoor federated learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 16 070–16 084, 2020.
- [30] G. Verma, E. Ciftcioglu, R. Sheatsley, K. Chan, and L. Scott, "Network traffic obfuscation: An adversarial machine learning approach," in *MILCOM 2018-2018 IEEE Military Communications Conference (MILCOM)*. IEEE, 2018, pp. 1–6.
- [31] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 IEEE Symposium on Security and Privacy (SP)*. Ieee, 2017, pp. 39–57.
- [32] M. Usama, J. Qadir, and A. Al-Fuqaha, "Adversarial attacks on cognitive self-organizing networks: The challenge and the way forward," in *2018 IEEE 43rd Conference on Local Computer Networks Workshops (LCN Workshops)*. IEEE, 2018, pp. 90–97.
- [33] M. Usama, A. Qayyum, J. Qadir, and A. Al-Fuqaha, "Black-box adversarial machine learning attack on network traffic classification," in *2019 15th International Wireless Communications & Mobile Computing Conference (IWCMC)*. IEEE, 2019, pp. 84–89.