

Deep Transfer Learning: A Novel Collaborative Learning Model for Cyberattack Detection Systems in IoT Networks

Tran Viet Khoa, Dinh Thai Hoang, Nguyen Linh Trung, Cong T. Nguyen,
Tran Thi Thuy Quynh, Diep N. Nguyen, Nguyen Viet Ha and Eryk Dutkiewicz

Abstract—Federated Learning (FL) has recently become an effective approach for cyberattack detection systems, especially in Internet-of-Things (IoT) networks. By distributing the learning process across IoT gateways, FL can improve learning efficiency, reduce communication overheads and enhance privacy for cyberattack detection systems. However, one of the biggest challenges for deploying FL in IoT networks is the unavailability of labeled data and dissimilarity of data features for training. In this paper, we propose a novel collaborative learning framework that leverages Transfer Learning (TL) to overcome these challenges. Particularly, we develop a novel collaborative learning approach that enables a target network with unlabeled data to effectively and quickly learn “knowledge” from a source network that possesses abundant labeled data. It is important that the state-of-the-art studies require the participated datasets of networks to have the same features, thus limiting the efficiency, flexibility as well as scalability of intrusion detection systems. However, our proposed framework can address these problems by exchanging the learning “knowledge” among various deep learning models, even when their datasets have different features. Extensive experiments on recent real-world cybersecurity datasets show that the proposed framework can improve more than 40% as compared to the state-of-the-art deep learning based approaches.

Index Terms—Cybersecurity, cyberattack detection, Internet of things (IoT), deep learning, transfer learning, federated learning.

I. INTRODUCTION

IN recent years, the rapid development of various technologies, such as 5G/6G, Industry 4.0, and Internet-of-Things (IoT), has enabled numerous applications to become an integral part in many aspects of our daily lives. However, such ever-fast growth has also led to an unprecedented massive amount of data and the proliferation of interconnected devices, e.g., sensors, smart cars, and cameras, which raises serious security and privacy concerns. Particularly, the increasing

number of emerging applications has also brought forth many new types of cyberattacks. For example, the number of new (zero-day) cyberattacks has increased by 60% from 2018 to 2019 [1]. Besides the dire consequences to the economic, e.g., ransomware alone cost more than \$5 billion globally in 2017 [2], cyberattacks pose serious threats to other areas with highly sensitive information such as healthcare and public security. As a result, cyberattack detection methods play a key role in detecting and promptly preventing consequences of cyberattacks in future IoT networks.

Recently, with outstanding classification ability, Machine Learning (ML) techniques, especially deep learning (DL), have been widely applied for cyberattack detection problems. Particularly, DL models can effectively learn the signatures of various cyberattack types. Moreover, DL models even can detect new types of attacks that have never been learned/trained before [3]. Nevertheless, DL-based cyberattack detection systems are also facing some practical challenges. Particularly, conventional DL approaches usually require a huge amount of data to achieve a high performance. However, in many applications, data are very difficult to collect because they are often stored locally on user devices such as IoT devices, smartphones, and wearable devices. This poses a threat to user privacy because sensitive data (e.g., location and private information) have to be sent over the network and stored at the centralized server for processing. Besides the privacy concerns, transmitting such a collectively large amount of data also imposes an extra communication burden over the network. Consequently, these limitations have been hindering the effectiveness of DL techniques in cyberattack detection systems.

To address these problems, Federated Learning (FL) has emerged to be a highly effective solution. Unlike conventional DL techniques that collect data and train the global model at a central server, FL enables the learning process to be distributed across all devices. Particularly, instead of sending data to a central server, the local data can be used to train a global model locally on each user device. Then, the obtained model weights of each device are periodically sent to a central server for aggregation. Afterward, the aggregated weights are sent back to all devices to update their local models’ weights. Since only the weights are transmitted in FL, both the privacy and communication overhead issues can be mitigated [4].

Despite its effectiveness, FL is still facing some challenges. Particularly, FL only performs well if the training data and

T. V. Khoa is with the School of Electrical and Data Engineering, University of Technology Sydney, Sydney, NSW 2007, Australia and the Advanced Institute of Engineering and Technology (AVITECH), University of Engineering and Technology, Vietnam National University, Hanoi, Vietnam (e-mail: khoa.v.tran@student.uts.edu.au, khoa.v.uet@vnu.edu.vn).

C. T. Nguyen, D. T. Hoang, D. N. Nguyen, and E. Dutkiewicz are with the School of Electrical and Data Engineering, University of Technology Sydney, Sydney, NSW 2007, Australia (e-mail: cong.nguyen@student.uts.edu.au, {hoang.dinh, diep.nguyen, eryk.dutkiewicz}@uts.edu.au).

N. L. Trung (corresponding author), T. T. T. Quynh, and N. V. Ha are with the Advanced Institute of Engineering and Technology (AVITECH), University of Engineering and Technology, Vietnam National University, Hanoi, Vietnam (e-mail: {linhtrung, quynhntt, hanv}@vnu.edu.vn).

the predicting data are independent and identically distributed (i.i.d). Consequently, they are not robust to the changes in the system, e.g., changes in network traffic due to the mobility of users, new types of devices participating in the network, and so on. Moreover, the performance of FL largely relies on the availability of labeled data. However, acquiring sufficient labeled data might be costly and time-consuming. Even if the data are available, the participated user data usually have different structures such as features. This leads to difficulty or even mistakes when FL aggregates the global model. Consequently, they may not be suitable for the intensive training process of FL [4] [5].

To address these limitations, transfer learning (TL) has been emerging as a promising solution, especially for problems related to heterogeneous training data [6]–[8]. Unlike DL and FL techniques that are trained only for a specific problem, TL can utilize “knowledge” from rich resource data to enhance the training process and performance of the ML models. Particularly, by transferring “knowledge” from similar scenarios with a lot of high-quality data, TL can address the lack of labeled data for the target networks. Moreover, the TL can exchange “knowledge” even if the data features of the target and source networks are not very similar [6], [9]. However, if the data features are too different, TL might even make the learning process worse than that without using TL, i.e., negative transfer [6]–[8]. In the context of cyberattack detection for IoT networks, negative transfer might be a serious problem since different networks may have various types of devices generating different data.

In this paper, we propose a novel collaborative learning framework that utilizes the strengths of both TL and FL to address the limitations of conventional DL-based cyberattack detection systems. Particularly, we consider a scenario with two different IoT networks¹. The first network (source network) has an abundant labeled data resource, while the second network (target network) has very little data resource (and most of them are unlabeled). Here, unlike most of the current works that assume that the data at these networks have the same features [10], we consider a much more practical and general case in which data at these two networks may have different features. To address the problem of dissimilar feature spaces of the target and source networks, we propose to transform them into a new joint feature-space. In this case, at each learning round of the federated learning process, trained models of target and source networks can be exchanged through the joint feature-space. Thus, by periodically exchanging and updating the trained model, the target network can eventually achieve the converged trained deep neural network that can predict attacks with high accuracy (thanks to useful “knowledge” transferred from the source network). Besides the exchanging and updating the learning model iteratively, we use a small number of mutual samples between two networks to mitigate the negative transfer learning. More importantly, unlike FL where networks try to train a joint global model, our proposed framework enables the participating networks

to obtain their particular trained models that are specific to their networks, i.e., better predict attacks for particular networks with different data structures. Extensive experiments on recent real-world datasets, including N-BaIoT [11] [12], KDD [13], NSL-KDD [14] and UNSW [15] show that our proposed framework can achieve an accuracy of up to 99% and an improvement of up to 40% over the unsupervised learning approach. The main contributions of this paper can be summarized as follows:

- We propose a novel collaborative learning framework that can effectively detect cyberattacks in decentralized IoT systems. By combining the strengths of FL and TL, our proposed framework can improve learning efficiency and the accuracy of cyberattack detection in comparison with the conventional DL-based cyberattack detection systems.
- We propose an effective transfer learning approach that can allow the deep learning model from the rich-data network to transfer useful knowledge to the low-data network even they have different features for cyberattack detection in IoT networks.
- We perform extensive experiments on recent real-world datasets including N-BaIoT, KDD, NSL-KDD, and UNSW to evaluate the performance of the proposed collaborative learning framework. The results show that our proposed approach can achieve an accuracy of up to 99% and an improvement of up to 40% over the unsupervised learning approach.

The rest of this paper is organized as follows. We first discuss related works in Section II. We then propose the federated transfer learning model for cyberattack detection in Section III. After that, simulation settings and results are discussed in Section IV. Finally, we conclude the paper in Section V.

II. RELATED WORK

A. Deep Learning for Cyberattack Detection

There have been a rich literature proposing DL approaches for cyberattack detection. In [16], a deep neural network (DNN) model is developed to detect zero-day attacks based on two types of data, i.e., network activities and local system activities. The results show that for most of the datasets, the proposed DNN can achieve a higher detection accuracy and lower false-positive rate compared to those of the other conventional machine learning classifiers such as K-Nearest Neighbors (KNN), and Support Vector Machine (SVM). Another DL approach is proposed in [17] to detect cyberattacks in the mobile cloud computing environments. The main difference between [16] and [17] is that the approach in [17] consists of a feature analysis phase before the learning phase. In the analysis phase, the datasets are analyzed to identify meaningful features, thereby reducing the data dimension and computational complexity. Experiments on the KDD [13], NSL-KDD [14], and the UNSW [15] datasets show that the proposed approach can achieve a detection accuracy of up to 97.1%.

¹The cases with multiple networks can be straightforwardly extended, e.g., by scheduling for networks to exchange information in order.

B. Federated Learning for Cyberattack Detection

With the advent of FL, the research focus has recently shifted towards applying this framework for cyberattack detection, especially in environments with numerous devices such as IoT and mobile edge networks. In [18], an FL framework is proposed for cyberattack detection in an edge network. In this network, the data for intrusion detection are stored locally at each edge node. The edge nodes train their data locally and send their models' weights to an FL server for aggregating. After aggregation, the FL server sends the weights back to all edge nodes. In this way, each edge node can benefit from the other nodes' data and training while protecting its privacy and reducing the network's communication burden. Experiments with the NSL-KDD datasets show that the proposed approach can achieve an accuracy of up to 99.2%. Another FL approach is proposed in [19] for attack detection in industrial cyber-physical systems. In the considered setting, there are multiple cyber-physical systems acting as FL nodes. However, unlike the previous frameworks, the authors propose a novel architecture combining a convolution neural network (CNN) and a gated recurrent unit for training at each FL node. Experiments with self-collected data show that the proposed approach can outperform other state-of-the-art approaches, e.g., [20]–[22], with an accuracy up to 99.2%. However, because of the limitations of FL as presented in the previous section, the learning model can only combine data with the same features and labels.

C. Transfer Learning for Cyberattack Detection

Although FL techniques can effectively address the privacy and communication load concerns of conventional ML for cyberattack detection, they are still facing some challenges. Particularly, FL approaches usually require high-quality and labeled data for training. However, collecting and labeling such data is expensive and time-consuming, especially for large-scale systems. On the other hand, unlabeled data are often abundant in environments such as IoT and mobile edge networks. Thus, a deep TL approach is proposed for IoT intrusion detection in [23] based on network activities, which can utilize both labeled and unlabeled data. In this approach, the authors employ two AEs. The first AE is trained with labeled data, while the second AE is trained with unlabeled data. Then, the knowledge is transferred from the first AE to the second AE by minimizing the Maximum Mean Discrepancy (MMD) distances between their weights. Experiments over nine IoT datasets were conducted to show that the proposed approach can achieve higher Area Under the Curve (AUC) scores compared to those of several other approaches.

Besides analyzing network traffic, another approach to detect cyberattacks is to analyze the devices' fingerprints. Particularly, attackers may try to impersonate a device in the system by copying its signal. For this kind of attack, ML techniques can be used to detect if the signals are coming from the real device or the malicious device. TL approaches such as [24]–[27] are proposed to identify cyberattacks based on device fingerprints. Among them, [26] and [27] leverage the environmental effects to classify signals from devices. To

improve the classification accuracy and address the lack of data, these approaches transfer the knowledge from nearby devices (since they share similar environmental effects). On the other hand, [24] and [25] leverage the knowledge from previous experiences, i.e., data collected in the past. These past data are then combined with the current data for training, thereby addressing the lack of fingerprint data.

Unlike all the abovementioned approaches, the collaborative learning framework proposed in this paper can leverage the strengths of both FL and TL to address limitations of ML-based intrusion detection systems, e.g., lack of labeled data, privacy and heterogeneous data feature space. Moreover, in our approach, each IoT network has a separated model that is fine-tuned specifically for that network, therefore the model is more effective for that network's cyberattack detection compared to FL frameworks with a single model for all networks. Furthermore, our proposed system model can utilize knowledge from both source and target data in the network instead of only transferring knowledge from a single source as proposed in most of the mentioned TL frameworks [23]–[25], [28], thereby mitigating the negative transfer problem.

III. PROPOSED FEDERATED TRANSFER LEARNING FRAMEWORK FOR CYBERATTACK DETECTION IN IOT NETWORKS

A. System Model

The conventional FL model requires to use a centralized server to maintain and aggregate all the trained models in the whole learning process. However, this may lead to a high cost to maintain and may not be effective to deploy in IoT networks. Thus, in this work, we propose a federated transfer learning model that allows the learning process to be performed more flexibly and effectively in IoT environments. In particular, we consider a network which has unlabeled data (e.g., Network B as illustrated in Fig. 1), and it wants to learn more knowledge from other networks with abundant labeled data. In this case, this network will connect with a target network (e.g., Network A as illustrated in Fig. 1) and nominate itself as a centralized node which can train its own data as well as perform transfer learning to exchange knowledge with the target network.

We denote a labeled cybersecurity dataset $D_A = \{X^A, Y^A, F^A\}$ of Network A with $(X^A, Y^A) = \{x_1^A, y_1^A, x_2^A, y_2^A, \dots, x_{M_A}^A, y_{M_A}^A\}$ where M_A is the number of samples of dataset A. In contrast, Network B has an unlabeled cybersecurity dataset $D_B = \{X^B, F^B\}$ with $(X^B) = \{x_1^B, x_2^B, \dots, x_{M_B}^B\}$ where M_B is the number of samples of dataset B. F^A, F^B are the feature spaces of Network A and Network B, respectively. The proposed model will perform transfer learning between two neural network by minimizing the total loss J to predict the label $P(z^B)$ for the unlabeled dataset of Network B. In this way, the network can help to improve the accuracy in identifying network traffics by learning useful knowledge from other labeled networks. Each network can be managed by an IoT gateway and possesses its own private dataset. The IoT gateway uses its deep learning model to detect normal and abnormal traffics. It is important

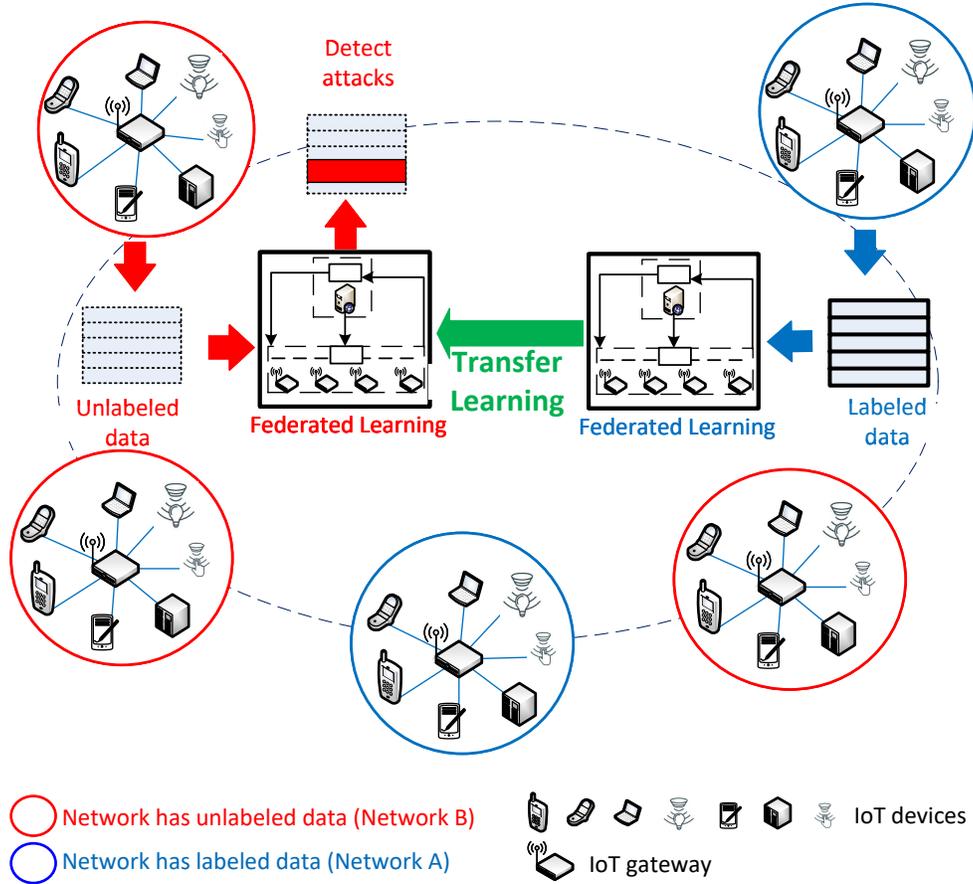


Fig. 1: Illustration of a system model for cyber attack detection in IoT networks.

to note that, unlike conventional FL approaches [29], in this work, we consider a practical scenario in which the datasets of networks may have different features.

B. Proposed Federated Transfer Learning Approach for Cyberattack Detection

In this section, we propose a highly-effective federated transfer learning model that can exchange knowledge between an unlabeled network and multiple networks which may have different features. To better analyze the impact of our proposed approach, we consider a specific scenario in which one labeled network is used as a source network to support an unlabeled network (i.e., target network). The scenario with one unlabeled network and multiple labeled networks can be straightforwardly extended, and we leave it for future study. Fig. 2 describes the training and predicting processes of FTL algorithm that we use in this case. The table of notations is presented in Table I. As described in previous section, Network A, Network B have their dataset D^A, D^B , respectively. They also have their model parameters called W^A and W^B . The outputs of two neural networks are calculated as follows:

$$Z^A = W^A * X^A, \quad (1a)$$

$$Z^B = W^B * X^B. \quad (1b)$$

TABLE I: Notations.

Notation	Description
X	The total samples of a dataset
Y	The labels of a dataset
F	The feature space of a dataset
x	A dataset sample
D	Network
M_A, M_B	The number of samples of dataset A, B, respectively
M_C	The number of predicted labels
M_{AB}	The overlapping samples between dataset A and dataset B
W_A, W_B	The parameter matrices of models A and B, respectively
Z_A, Z_B	The outputs of models A and B, respectively
z	The output of an input sample after learning model
j	The loss of an input sample
J	The loss function
γ, λ	Weight parameters
w	Training parameters

Algorithm 1 Federated Transfer Learning Algorithm: Training Process

```

1: Input: The learning rate  $\eta$ , the weight parameter  $\gamma, \lambda$ , the
   maximum iteration  $T$ , the tolerance  $t$  and Network A and
   Network B initialize model parameters  $W^A, W^B$ ;
2: Output: The trained model parameter  $W^A, W^B$ ;
3:  $iteration = 0$ 
4: while  $iteration \leq T$  do
5:   Network A performs:
6:    $z_i^A = h_i^A * x_i^A$  for  $i \in D_A$ ;
7:   Send  $\{z_i^A, y_i^A\}$  to Network B;
8:   Network B performs:
9:    $z_i^B = h_i^B * x_i^B$  for  $i \in D_B$ ;
10:  Send  $\{z_i^B\}$  to Network A;
11:  Network A performs:
12:  Compute  $\frac{\partial J}{\partial w_i^A}$  and  $J^A$ , then send them to Network B;
13:  Network B performs:
14:  Compute  $\frac{\partial J}{\partial w_i^B}$ ,  $J^B$  and  $J^{AB}$ , then send them to Net-
     work A;
15:  Network A performs:
16:  Update  $w_i^A = w_i^A - \eta \frac{\partial J}{\partial w_i^A}$ ;
17:  Network B performs:
18:  Update  $w_l^B = w_l^B - \eta \frac{\partial J}{\partial w_l^B}$ ;
19:  if  $J_{prev} - J \leq t$  then
20:    Send stop signal to Network B;
21:    Break.
22:  else
23:     $J_{prev} = J$ ;
24:     $iteration = iteration + 1$ ;
25:    continue;
26:  end if
27: end while

```

We need to find the prediction function $P(z_j^B) = P(z_1^A, y_1^A, \dots, z_{M_A}^A, y_{M_A}^A, z_j^B)$ to predict the output of Network B. To find a high-quality predict function, we first need to minimize the loss function using the labeled dataset as follows:

$$\arg \min_{W^A, W^B} J^B = \sum_i^{M_c} j^B(y_i^A, P(z_i^B)), \quad (2)$$

where M_c is the number of predicted labels, and j^B represents the loss of the loss function which depends on the type of output or mechanism, i.e., the logistic loss function [30] with the predicted value \mathbf{z} and the labeled \mathbf{y} :

$$j^B(\mathbf{z}, \mathbf{y}) = \log(1 + \exp(-\mathbf{z} \times \mathbf{y})). \quad (3)$$

In addition, datasets A and B may have some overlapping samples, and thus we can use these samples to optimize the loss function. We denote M_{AB} as the overlapping samples between dataset A and dataset B. We need to minimize the alignment loss function between A and B as follows:

$$\operatorname{argmin}_{W^A, W^B} J^{AB} = - \sum_i^{M_{AB}} j^{AB}(z_i^A, z_i^B), \quad (4)$$

where j^{AB} represents the alignment loss function. The common alignment loss function can be represented in modulus $j^{AB} = \|z_i^A - z_i^B\|^2$ or angle $j^{AB} = -z_i^A * z_i^B$. Lastly, we add the regularization $J_R^A = \sum_l^{L_A} \|w_l^A\|^2$ and $J_R^B = \sum_l^{L_B} \|w_l^B\|^2$ in which L_A and L_B are the numbers of layers in neutron Network A and Network B, respectively, to find the final loss function that needs to be minimized:

$$\operatorname{argmin}_{W^A, W^B} J = J^B + \gamma J^{AB} + \frac{\lambda}{2} (J_R^A + J_R^B), \quad (5)$$

where γ and λ are the weight parameters. The gradient for updating W^A, W^B are calculated by the following formula:

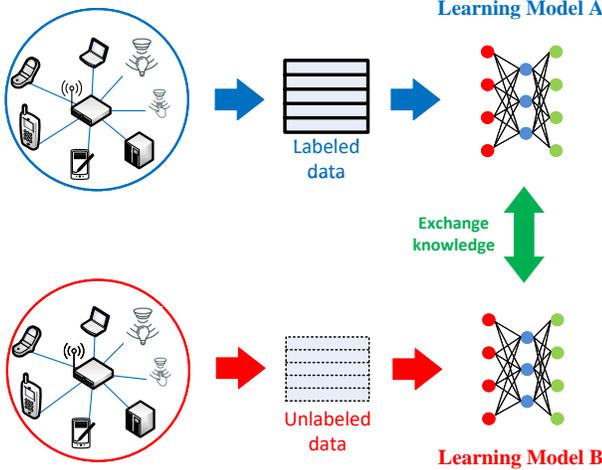
$$\frac{\partial J}{\partial w_l^i} = \frac{\partial J^B}{\partial w_l^i} + \gamma \frac{\partial J^{AB}}{\partial w_l^i} + \lambda w_l^i. \quad (6)$$

The training process is presented in Algorithm 1. Specifically, we first initialize W^A and W^B . Next, we calculate z_i^A and z_i^B from the input samples of dataset A (D^A) and dataset B (D^B) as shown in Equation (1). Then, Network A sends $\{z_i^A, y_i^A\}$ to Network B to calculate J^B , the alignment loss function J^{AB} and the gradients of J^B as shown in Equations (2), (4), (5) and (6), respectively. Similarly, Network B sends $\{z_i^B\}$ to Network A to calculate J^A as in Equation (5). In Equation (4), we use M_{AB} as the mutual samples of two datasets. For example, the same IoT devices are attacked by the same types of cyberattacks in different networks. Each network extracts the attack data with different features, e.g., Network A uses timeslot, packet header, ip address while Network B uses MAC address, error packets, frame header. The number of mutual samples is an important factor that strongly supports the learning process between two networks (we will explain it more details in Section IV). After that, we calculate the final loss function J and the gradient as in Equation (5) and Equation (6). Finally, Network A and Network B update their model parameters based on the gradient and loss functions. This process continuously repeats until the system converges or reaches the maximum number of iterations to minimize the final loss function in Equation (5).

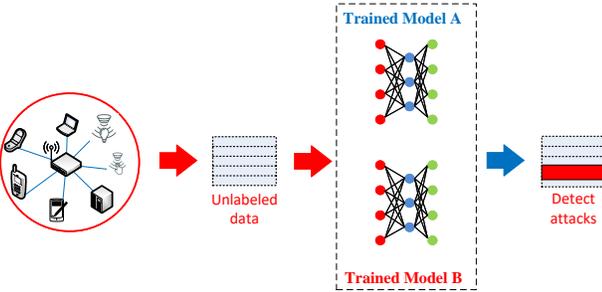
When the training completes, the prediction process described in the Algorithm 2 is called to predict the final result of the unlabeled dataset D_B . In this process, both Network A and Network B have their trained models. Similar to the training process, the dataset D_B firstly goes through the the trained model of Network B to calculate Z^B . Then, Network B sends Z^B to Network A to archive the transfer learning knowledge from trained model of Network A. Network A predicts the results and sends them back to Network B to classify the attack and normal behaviors of the network.

C. Evaluation Methods

As mentioned in [31], [32], the confusion matrix is typically used to evaluate system performance, especially for intrusion detection systems. We denote TP, TN, FP, and FN to be ‘‘True Positive’’, ‘‘True Negative’’, ‘‘False Positive’’, and ‘‘False Negative’’, respectively. The Receiver Operator Characteristic (ROC) is created by plotting the TPR over FPR at different thresholds. Then, we use Area Under the Curve (AUC) to



(a) FTL training process.



(b) FTL predicting process.

Fig. 2: The FTL algorithm.

Algorithm 2 Federated Transfer Learning Algorithm: Predicting Process

- 1: **Input:** The model parameters W^A, W^B and dataset X_B ;
 - 2: **Output:** The prediction Y^B ;
 - 3: Network B performs:
 - 4: $z_i^B = h_i^B * x_i^B$ for $i \in D_B$;
 - 5: Send $\{z_i^B\}$ to Network A;
 - 6: Network A performs:
 - 7: Compute $P(z_i^B) = W^A[z_i^B]$ and send it to Network B.
-

evaluate the performance of the algorithm in the following formula:

$$\xi = \int_{x=0}^1 \text{TP}(\text{FP}^{-1}(x)) dx. \quad (7)$$

In our experiments, we randomly select samples from original dataset to test the algorithm. In this scenario, the p -value is often used to evaluate the results of random tests, and is given by

$$p = F(\xi|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\xi} e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt, \quad (8)$$

in which μ is the mean and σ is the standard deviation.

The results are calculated by the significant number with the following formula:

$$\text{Sig} = F^{-1}(p|\mu, \sigma) = \{\xi : F(\xi|\mu, \sigma) = p\}, \quad (9)$$

where Sig is the significant number that represents the results of 30 random runs and the confidence of this number is calculated by $\text{conf} = 1 - p$. In a normal situation, p is considered confidence when it has values around 0.01 and 0.05, corresponding to the confidence of significant numbers is around 99% and 95%.

IV. PERFORMANCE ANALYSIS

A. Datasets

In this experiment, we use four popular cybersecurity datasets namely the N-BaIoT [11] [12], KDD [13], NSL-KDD [14] and UNSW [15] datasets to evaluate the performance of the proposed method. The Network-based Detection of IoT Botnet Attacks (N-BaIoT) dataset [11] [12] includes the information collected in the setup network about the normal and attack situation. The attack was performed by servers to nine IoT devices and the total network behavior was captured by the sniffer server to extract dataset. This dataset is characterized by 115 features for both normal and attack behaviors. In this dataset, the attack type is the Distributed Denial of Service (DDoS) which was implemented by two well-known botnets, namely Mirai and BASHLITE. The BASHLITE botnet includes 5 types of attacks, i.e., network scanning (scan), spam data sending (junk), UDP flooding (udp), TCP flooding (tcp), and the join of sending spam data and opening port to specific IP address (combo). Besides BASHLITE, the Mirai botnet also includes 5 types of attacks, i.e., scan, ACK flooding (ack), SYN flooding (syn), udp, and optimized UDP flooding (udpplain).

In addition to IoT datasets, we also want to evaluate our proposed solution on some classical intrusion detection datasets, i.e., KDD [13], NSL-KDD [14] and UNSW [15] datasets. The KDD dataset [13] includes many different kinds of network attacks simulated in military network environment. The KDD dataset has 41 features and it classifies attacks into 4 groups including Denial of Service (DoS), Probe, User to Root (U2R), Remote to Local (R2L). The NSL-KDD dataset [14] inherits the properties from KDD [13] dataset such as the features and types of attacks but eliminates the redundant samples in the training dataset and the duplicated samples in the testing dataset. Although both KDD and NSL-KDD datasets are well-known and used in many research works, they were developed long time ago. Thus, some modern attacks were not involved. Therefore, a recent dataset, i.e., UNSW dataset [15], is considered in this work. Unlike KDD and NSL-KDD, the feature space of this dataset includes 42 types and 9 kinds of attacks, namely DoS, Backdoors, Worms, Fuzzers, Analysis, Reconnaissance, Exploits, Shellcode, and Generic.

B. Experiment Setup

In this section, we carry out experiments using all the aforementioned datasets to evaluate the performance of the proposed solution. In this experiment, we denote IoT1-9 as

	FTL	UDL
IoT1	85.771	45.753
IoT2	83.795	63.171
IoT3	94.286	80.453
IoT4	79.241	77.885
IoT5	90.605	81.876
IoT6	91.179	82.703
IoT7	90.670	85.183
IoT8	82.960	65.256
IoT9	83.222	73.072
KDD	99.315	80.477
NSLKDD	98.485	83.025
UNSW	97.072	68.449

(a) The results with $p = 1$.

	FTL	UDL
IoT1	87.398	49.770
IoT2	85.672	65.793
IoT3	94.896	81.070
IoT4	81.672	77.885
IoT5	91.517	82.013
IoT6	92.059	82.703
IoT7	92.030	86.013
IoT8	85.197	68.161
IoT9	85.072	73.078
KDD	99.395	81.304
NSLKDD	98.534	83.450
UNSW	97.141	69.124

(b) The results with $p = 3$.

	FTL	UDL
IoT1	88.259	51.897
IoT2	86.666	67.181
IoT3	95.220	81.397
IoT4	82.959	77.885
IoT5	92.000	82.085
IoT6	92.525	82.703
IoT7	92.750	86.453
IoT8	86.381	69.700
IoT9	86.052	73.082
KDD	99.438	81.742
NSLKDD	98.561	83.675
UNSW	97.177	69.482

(c) The results with $p = 5$.

TABLE II: The results with multiple datasets in CASE 1.

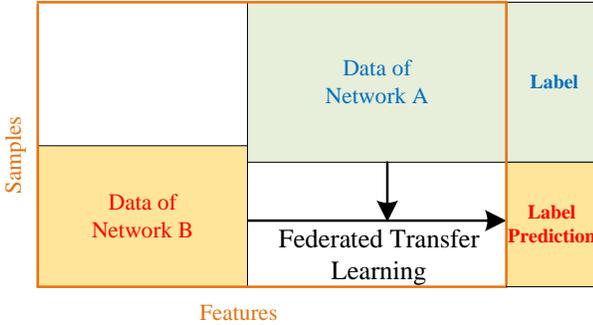


Fig. 3: The data of participated networks used in this experiment.

the dataset names of nine IoT devices. Table III describes the total features and the representative names of datasets that we use in this experiment. Fig. 3 also describes the separated data in each dataset in this experiment. In this experiment, the participated data are randomly selected from the dataset. Then, the selected data are separated into label data (data of Network A) and unlabeled data (data of Network B) with different features as described in Table III. These data have about 10% mutual samples of total dataset samples. We experiment with

two cases, i.e., the first one is with 2000 unlabeled data and 9577 labeled data (CASE 1), the second one is with 10000 unlabeled data and 47893 labeled data (CASE 2).

In this setup, we consider a baseline solution with the state-of-the-art unsupervised deep learning model (UDL) which clusters the unlabeled data into normal and attack behaviors based on autoencoder and k-means techniques [33]. The unsupervised deep learning model includes an autoencoder and k-nearest neighbor to cluster the unlabeled data. In addition, we consider the second baseline solution that uses both supervised and unsupervised datasets to feed the FTL learning models. The FTL will exchange the knowledge from the supervised learning model and the unsupervised learning model to improve the accuracy of learning as well as increase the precise of identifying attack and normal behaviors of the unlabeled data. Then, we measure the AUC of this process 30 times to calculate the signification number of the AUC series results with both baseline solutions. Finally, we plot the reconstruction errors to analyze the convergence of the FTL algorithm for all datasets.

C. Experimental Results

In this section, we show the results of our experiments with different kinds of cybersecurity datasets.

Dataset	Device name	Features of Network A	Features of Network B	Total features
IoT1	Danmini_Doorbell	85	30	115
IoT2	Ecobee_Thermostat	85	30	115
IoT3	Ennio_Doorbell	85	30	115
IoT4	Philips_B120N10_Baby_Monitor	85	30	115
IoT5	Provision_PT_737E_Security_Camera	85	30	115
IoT6	Provision_PT_838_Security_Camera	85	30	115
IoT7	Samsung_SNH_1011_N_Webcam	85	30	115
IoT8	SimpleHome_XCS7_1002_WHT_Security_Camera	85	30	115
IoT9	SimpleHome_XCS7_1003_WHT_Security_Camera	85	30	115
KDD	-	31	10	41
NSLKDD	-	31	10	41
UNSW	-	31	11	42

TABLE III: Dataset preparation

	FTL	UDL
IoT1	90.371	49.783
IoT2	68.193	62.591
IoT3	94.525	83.411
IoT4	87.050	77.725
IoT5	86.535	81.954
IoT6	87.214	82.555
IoT7	97.662	79.517
IoT8	84.609	52.702
IoT9	90.095	63.803
KDD	99.535	84.333
NSLKDD	98.858	81.164
UNSW	97.049	66.329

(a) The results with $p = 1$.

	FTL	UDL
IoT1	91.497	54.079
IoT2	72.573	65.681
IoT3	95.073	83.565
IoT4	88.538	77.781
IoT5	88.150	82.160
IoT6	88.638	82.664
IoT7	97.928	81.400
IoT8	86.691	57.318
IoT9	90.959	65.559
KDD	99.562	84.423
NSLKDD	98.885	81.976
UNSW	97.121	66.901

(b) The results with $p = 3$.

	FTL	UDL
IoT1	92.093	56.354
IoT2	74.892	67.317
IoT3	95.363	83.647
IoT4	89.326	77.811
IoT5	89.006	82.269
IoT6	89.392	82.721
IoT7	98.069	82.397
IoT8	87.793	59.763
IoT9	91.417	66.489
KDD	99.576	84.471
NSLKDD	98.900	82.406
UNSW	97.159	67.203

(c) The results with $p = 5$.

TABLE IV: The results with multiple datasets in CASE 2.

1) *Accuracy Comparison*: In this section, we compare the performance of FTL and the unsupervised deep learning (UDL) method in terms of the significant number of each p as explained in Section III. Table II and Table IV describe the significant number of each dataset with $p = 1, 3, 5$ corresponding to the confidence of 99%, 97%, 95%.

In general, Table II and Table IV show that the significant numbers of all datasets increase as p increases. This is because in (9), we calculate the significant number based on a series of 30 continuous AUC results. When p increases, the AUC results increase in all tables. This demonstrates that most of the AUC results in 30 series are higher than the significant number in the case where $p = 1$.

Table II(c) shows the significant numbers of participated datasets with $p = 5$ in CASE 1. In this table, the IoT1 and UNSW datasets show a significant gap of about 30% and 40% between FTL and UDL. These results show the difficulty of clustering in recognizing the groups of samples and the advantage of collaborative learning in these datasets. The other ten datasets have gaps of around 10-20% between the two methods, which demonstrate the stability of our proposed solution for any cybersecurity dataset.

In addition, Table IV(c) shows the significant numbers of multiple datasets with $p = 5$ in CASE 2. In this table, the significant numbers also have a gap of around 10-40% between the two solutions. It shows the common trend that the significant numbers increase for most datasets when the number of samples increases. However, in IoT2, IoT5, and IoT6 datasets, the significant numbers slightly decrease because of the randomly selected samples from the original dataset. It also can be demonstrated by the high fluctuation of the reconstruction errors of IoT2, IoT5, IoT6 datasets in Fig. 5(b) compared with other datasets. However, in all studied datasets, our proposed solution still performs much better than the state-of-the-art UDL solution. These results demonstrate that our solution can work efficiently in all IoT and conventional cybersecurity datasets in detecting cyberattacks in the network.

2) *Reconstruction Error Analysis*: In this section, we discuss the convergence of the FTL algorithm in each dataset. Fig. 4 describes the reconstruction errors of the nine IoT

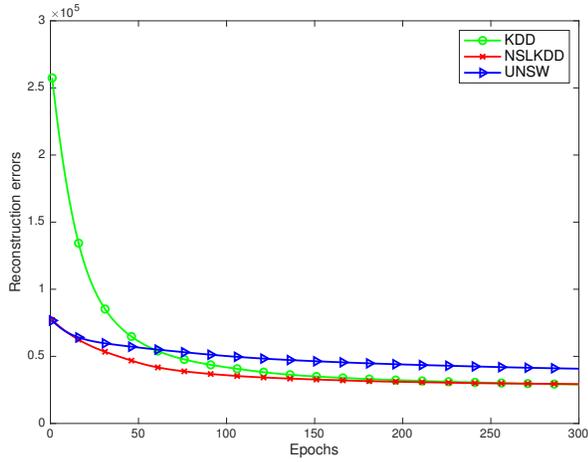
datasets and the conventional datasets like KDD, NSLKDD, and UNSW in CASE 1. Fig. 5 describes the reconstruction errors of study datasets in CASE 2.

In Fig. 4(a) and Fig. 5(a), we can see that at the first few epochs, the errors are very high for KDD (up to 2.6×10^5 in CASE 1 and 12×10^5 in CASE 2), but this error dramatically reduces to 0.3×10^5 in CASE 1 and 1.5×10^5 in CASE 2 after only 200 epochs. For NSLKDD and UNSW, they have very similar trends with 0.75×10^5 in CASE 1 and 3.8×10^5 in CASE 2 at the beginning and gradually reduce to 0.4×10^5 in CASE 1 and 1.9×10^5 in CASE 2 after 200 epochs, respectively. After 200 epochs, the algorithm converges as all the reconstruction error curves are flattened.

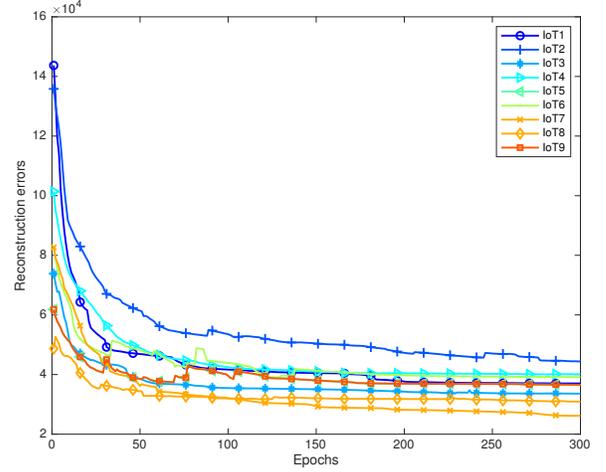
Fig. 4(b) and Fig. 5(b) show the reconstruction errors of nine IoT datasets in both CASE 1 and CASE 2. we can observe the same trend over all datasets, i.e., all errors gradually reduce when the number of epochs increases. However, it can be observed that the trend exhibits some fluctuations in comparison with the trends in Fig. 4(a) and Fig. 5(a) because of the heterogeneous distribution in IoT datasets. The high fluctuation of the reconstruction errors of IoT2, IoT5, IoT6 datasets in Fig. 5(b) also explains why their significant numbers reduce when the number of samples increases in CASE 2. However, the reconstruction errors of all studied datasets in our proposed solution dramatically decrease and become stable after 200 running epochs in both cases.

3) *Mutual Information Analysis*: As mentioned in the previous section, network A and network B may share a number of mutual samples. The FTL algorithm exploits the information of these mutual samples to perform the prediction for unlabeled data of network B. This section provides the analysis results to identify how this mutual information can affect to the results of label prediction. In this section, we perform the simulation in CASE 2 with a larger number of samples than in CASE 1. Fig. 6 gives information about the variation of AUC when the percentage of mutual data increases.

Fig. 6(a) shows the increase of AUC on KDD, NSLKDD, and UNSW datasets when the percentage of mutual samples increases from 0.005% to 10%. The AUC of KDD and UNSW datasets sharply increase and remain stable at around 96% on

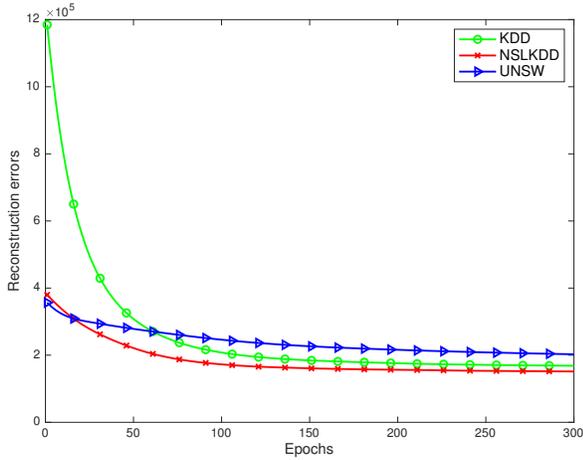


(a) The reconstruction errors of KDD, NSLKDD and UNSW datasets.

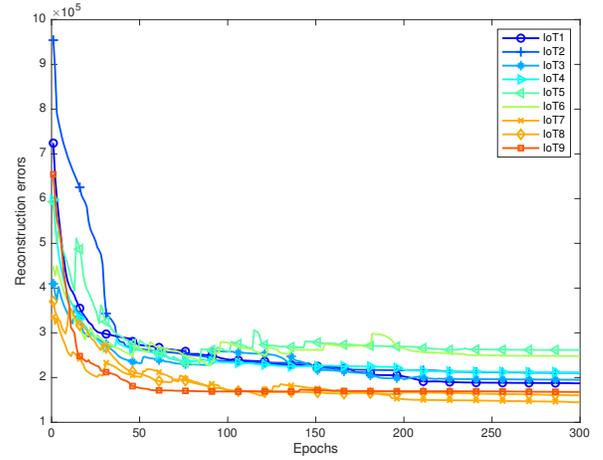


(b) The reconstruction errors of IoT datasets.

Fig. 4: The reconstruction errors in CASE 1.



(a) The reconstruction errors of KDD, NSLKDD and UNSW datasets.



(b) The reconstruction errors of IoT datasets.

Fig. 5: The reconstruction errors in CASE 2.

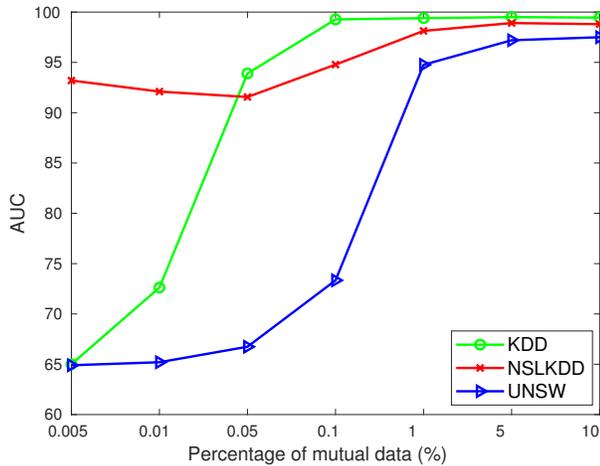
the NSLKDD dataset with about 5% to 10% mutual samples. A similar trend happens with the IoT datasets in Fig. 6(b) when the AUCs of all nine IoT datasets increase and remain stable at approximately 10% of mutual samples. From these results, it can be observed that achieving high efficiency in AUC for IoT datasets may require at least 10% of mutual data.

In summary, the results with 12 cybersecurity datasets show the outperformance of our proposed model in comparison with the state-of-the-art unsupervised deep learning in term of accuracy as shown in Table II for CASE 1 and Table IV for CASE 2, especially with IoT1 and UNSW datasets. Moreover, the reconstruction errors show a fluctuation of the IoT datasets when the number of samples increases due to noise from the collected datasets of some IoT devices. Finally, we vary the amount of mutual data between two networks to evaluate the accuracy of our proposed model. The results show that the proposed model can achieve high performance with 10%

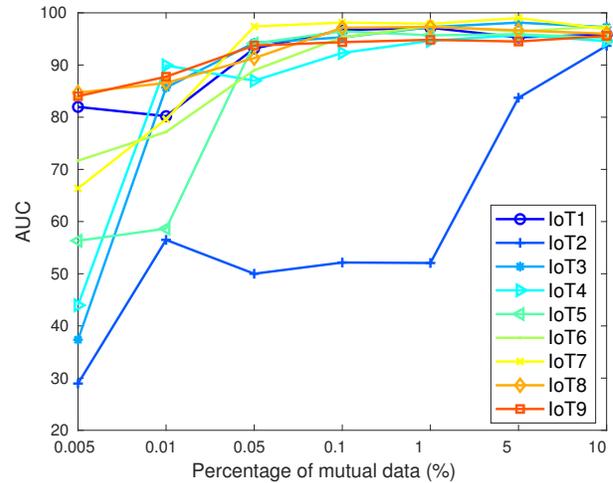
mutual data with all datasets.

V. CONCLUSION

In this work, we have proposed a novel collaborative learning framework to address the limitations of current ML-based cyberattack detection systems in IoT networks. In particular, by extracting and transferring knowledge from a network with abundant labeled data (source network), the intrusion detection performance of the target network can be significantly improved (even if the target has very few labeled data). More importantly, unlike most of the current works in this area, our proposed framework can enable the source network to transfer the knowledge to the target network even when they are different data structures, e.g., different features. The experimental results then show that the accuracy of prediction of our proposed framework is significantly improved in comparison with the state-of-the-art unsupervised deep learning model. In



(a) The percentage mutual information of KDD, NSLKDD and UNSW datasets.



(b) The percentage mutual information of IoT datasets.

Fig. 6: The illustration of AUC with different percentage of mutual information.

addition, the convergence of the proposed collaborative learning model is also analyzed with various cybersecurity datasets. In future work, we can consider using other effective transfer learning techniques to make transfer learning processes more stable and achieve better performance, especially when the amount of mutual information is very limited.

VI. ACKNOWLEDGEMENTS

This work is the output of the ASEAN IVO http://www.nict.go.jp/en/asean_ivo/index.html project “Cyber-Attack Detection and Information Security for Industry 4.0” and financially supported by NICT <http://www.nict.go.jp/en/index.html>.

This work was supported in part by the Joint Technology and Innovation Research Centre – a partnership between the University of Technology Sydney and the University of Engineering and Technology, Vietnam National University, Hanoi, Vietnam.

This research was supported in part by the Australian Research Council under the DECRA project DE210100651.

The work of Cong T. Nguyen was funded in part by Vingroup JSC and supported in part by the Master, PhD Scholarship Programme of Vingroup Innovation Foundation (VINIF), Institute of Big Data, code VINIF.2021.TS.006.

REFERENCES

- [1] Y. Keshet, “Half of the malware detected in 2019 was classified as zero-day threats, making it the most common malware to date,” Mar. 2020. [Online]. Available: <https://www.cynet.com/blog/half-of-the-malware-detected-in-2019-was-classified-as-zero-day-threats-making-it-the-most-common-malware-to-date/>
- [2] S. Morgan, “Global ransomware damage costs predicted to hit \$11.5 billion by 2019,” Mar. 2021. [Online]. Available: <https://cybersecurityventures.com/ransomware-damage-report-2017-part-2/>
- [3] T. V. Khoa, D. H. Son, D. T. Hoang, N. L. Trung, T. T. T. Quynh, D. N. Nguyen, N. V. Ha, and E. Dutkiewicz, “Collaborative learning for cyberattack detection in blockchain networks,” *arXiv preprint arXiv:2203.11076*, Mar. 2022.
- [4] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, and C. Miao, “Federated learning in mobile edge networks: A comprehensive survey,” *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 2031–2063, Apr. 2020.
- [5] M. Aledhari, R. Razzak, R. M. Parizi, and F. Saeed, “Federated learning: A survey on enabling technologies, protocols, and applications,” *IEEE Access*, vol. 8, pp. 140 699–140 725, July 2020.
- [6] C. T. Nguyen, N. Van Huynh, N. H. Chu, Y. M. Saputra, D. T. Hoang, D. N. Nguyen, Q.-V. Pham, D. Niyato, E. Dutkiewicz, and W.-J. Hwang, “Transfer learning for future wireless networks: A comprehensive survey,” *arXiv preprint arXiv:2102.07572*, Feb. 2021.
- [7] S. Niu, Y. Liu, J. Wang, and H. Song, “A decade survey of transfer learning (2010–2020),” *IEEE Transactions on Artificial Intelligence*, vol. 1, no. 2, pp. 151–166, Oct. 2020.
- [8] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He, “A comprehensive survey on transfer learning,” *Proceedings of the IEEE*, vol. 109, no. 1, pp. 43–76, July 2020.
- [9] M. Xu, D. T. Hoang, J. Kang, D. Niyato, Q. Yan, and D. I. Kim, “Secure and reliable transfer learning framework for 6g-enabled internet of vehicles,” *IEEE Wireless Communications*, May 2022.
- [10] M. A. Ferrag, O. Friha, L. Maglaras, H. Janicke, and L. Shu, “Federated deep learning for cyber security in the internet of things: Concepts, applications, and experimental analysis,” *IEEE Access*, vol. 9, pp. 138 509–138 542, Oct. 2021.
- [11] Y. Meidan, M. Bohadana, Y. Mathov, Y. Mirsky, A. Shabtai, D. Breitenbacher, and Y. Elovici, “N-BaIoT-Network-based detection of IoT botnet attacks using deep autoencoders,” *IEEE Pervasive Computing*, vol. 17, no. 3, pp. 12–22, July 2018.
- [12] Y. Mirsky, T. Doitshman, Y. Elovici, and A. Shabtai, “Kitsune: An ensemble of autoencoders for online network intrusion detection,” *arXiv preprint arXiv:1802.09089*, Feb. 2018.
- [13] “KDD dataset,” <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>.
- [14] “NSL-KDD dataset,” <https://www.unb.ca/cic/datasets/nsl.html>.
- [15] N. Moustafa and J. Slay, “UNSW-NB15: a comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set),” in *2015 Military Communications and Information Systems Conference*, Nov. 2015, pp. 1–6.
- [16] R. Vinayakumar, M. Alazab, K. Soman, P. Poornachandran, A. Al-Nemrat, and S. Venkatraman, “Deep learning approach for intelligent intrusion detection system,” *IEEE Access*, vol. 7, pp. 41 525–41 550, Apr. 2019.
- [17] K. K. Nguyen, D. T. Hoang, D. Niyato, P. Wang, D. Nguyen, and E. Dutkiewicz, “Cyberattack detection in mobile cloud computing: A deep learning approach,” in *2018 IEEE Wireless Communications and Networking Conference*, Apr. 2018, pp. 1–6.
- [18] A. Abeshu and N. Chilamkurti, “Deep learning: The frontier for distributed attack detection in fog-to-things computing,” *IEEE Communications Magazine*, vol. 56, no. 2, pp. 169–175, Feb. 2018.
- [19] B. Li, Y. Wu, J. Song, R. Lu, T. Li, and L. Zhao, “Deepfed: Federated deep learning for intrusion detection in industrial cyber-physical systems,” *IEEE Transactions on Industrial Informatics*, Sep. 2020.

- [20] T. D. Nguyen, S. Marchal, M. Miettinen, H. Fereidooni, N. Asokan, and A.-R. Sadeghi, "D̄iot: A federated self-learning anomaly detection system for iot," in *2019 IEEE 39th International Conference on Distributed Computing Systems*, July 2019, pp. 756–767.
- [21] W. Schneble and G. Thamarasu, "Attack detection using federated learning in medical cyber-physical systems," in *28th International Conference on Computer Communications and Networks*, July 2019, pp. 1–8.
- [22] 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, Apr. 2020.
- [23] L. Vu, Q. U. Nguyen, D. N. Nguyen, D. T. Hoang, and E. Dutkiewicz, "Deep transfer learning for IoT attack detection," *IEEE Access*, vol. 8, pp. 107 335–107 344, June 2020.
- [24] Y. Sharaf-Dabbagh and W. Saad, "Transfer learning for device fingerprinting with application to cognitive radio networks," in *2015 IEEE 26th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications*, Aug. 2015, pp. 2138–2142.
- [25] C. Zhao, Z. Cai, M. Huang, M. Shi, X. Du, and M. Guizani, "The identification of secular variation in IoT based on transfer learning," in *2018 International Conference on Computing, Networking and Communications*, Mar. 2018, pp. 878–882.
- [26] Y. Sharaf-Dabbagh and W. Saad, "On the authentication of devices in the internet of things," in *2016 IEEE 17th International Symposium on A World of Wireless, Mobile and Multimedia Networks*, June 2016, pp. 1–3.
- [27] Y. S. Dabbagh and W. Saad, "Authentication of wireless devices in the internet of things: Learning and environmental effects," *IEEE Internet of Things Journal*, vol. 6, no. 4, pp. 6692–6705, Apr. 2019.
- [28] T. Wen and R. Keyes, "Time series anomaly detection using convolutional neural networks and transfer learning," *arXiv preprint arXiv:1905.13628*, May 2019.
- [29] T. V. Khoa, Y. M. Saputra, D. T. Hoang, N. L. Trung, D. Nguyen, N. V. Ha, and E. Dutkiewicz, "Collaborative learning model for cyberattack detection systems in iot industry 4.0," in *2020 IEEE Wireless Communications and Networking Conference*, May 2020, pp. 1–6.
- [30] Y. Liu, Y. Kang, C. Xing, T. Chen, and Q. Yang, "A secure federated transfer learning framework," *IEEE Intelligent Systems*, vol. 35, no. 4, pp. 70–82, Apr. 2020.
- [31] T. Fawcett, "An introduction to ROC analysis," *Pattern Recognition Letters*, vol. 27, no. 8, pp. 861–874, June 2006.
- [32] D. M. Powers, "Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation," *Journal of Machine Learning Technologies*, pp. 37–63, Oct. 2011.
- [33] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.