

Safeguarding Privacy and Integrity of Federated Learning in Heterogeneous Cross-Silo IoRT Environments: A Moving Target Defense Approach

Zan Zhou, Changqiao Xu, *Senior Member, IEEE*, Shujie Yang, Xiaoyan Zhang, Hongjing Li, Sizhe Huang, and Gabriel-Miro Muntean, *Fellow, IEEE*

Abstract—**Bridging the gap between the Internet of Things and collaborative robots, the recent advancements in the Internet of Robotic Things (IoRT) aim at significantly improving production and operation efficiency and quality. As the scope and complexity of IoRT continue to expand, involving also very large numbers of robots, there is a need for employment of innovative solutions such as federated learning. However, this growing demand is accompanied by multiple challenges, including threats to data privacy and model integrity. Besides, the heterogeneity of the robots and their interaction, multiplies these challenges. In this paper, we discuss the key concerns of collaborative training in IoRT, and propose a shuffling-based moving target defense approach for federated learning in heterogeneous cross-silo IoRT environments (FedMTD). Based on a hierarchical training structure with node clustering, FedMTD bounds heterogeneity by domains, thereby minimizing the learning error and privacy loss. It also enhances resistance to poisoning attacks through decentralized credit evaluation. Experimental results show that FedMTD brings significant improvements in learning performance, privacy enhancement, and poisoning resistance.**

Index Terms—Internet of robotic things, Federated learning, Privacy, Moving target defense

I. INTRODUCTION

AS a central facet of Industry 5.0, the Internet of Robotic Things (IoRT) [1] [2] aims to revolutionize manufacturing by amalgamating the Internet of Things (IoT), artificial intelligence, digital twins, human-robot collaboration, and an array of emerging technologies. Recent advancements in 6G integrated sensing, computing, and communication (ISCC) coupled with Artificial Intelligence Generated Content (AIGC) have equipped the latest generation of robots with the capability to collaboratively perceive their environment and make intelligent decisions. In the realm of **IoRT**, a multitude of AI-powered services, such as Robotics as a Service (RaaS), collaborative robots (Cobots), and Robotic Process Automation (RPA), are thriving. According to the “Global Opportunity Analysis and Industry Forecast” report by Allied Market Research, the global IoRT market is anticipated to reach \$2,461.9 billion by 2031, exhibiting a remarkable compound annual growth rate (CAGR) of 28.6% from

Zan Zhou, Changqiao Xu, Shujie Yang, Xiaoyan Zhang, Hongjing Li, and Sizhe Huang are with Beijing University of Posts and Telecommunications, China; Gabriel-Miro Muntean is with the Performance Engineering Laboratory, School of Electronic Engineering, Dublin City University, Dublin, Ireland.

2022 to 2031 (<https://www.alliedmarketresearch.com/internet-of-robotic-things-market-A31839> [Accessed Sep. 8, 2023]).

To accommodate the highly diversified and increasingly complex demands involving growing numbers of robots, employment of innovative solutions such as federated learning (FL) is in urgent need [3]. As shown in Figure 1, taking supervisory control and data acquisition (SCADA) system as an example, geographically dispersed edge robot (ER) nodes, each equipped with optical character recognition (OCR) camera modules from multiple factories, can collaboratively train an intelligent model to facilitate product sorting. Nevertheless, the following long-neglected yet important concerns still hinder the wide deployments of FL in IoRT environments:

- **High heterogeneity:** ERs from varied factories or even production lines can present distinct data distributions [3].
- **Privacy breach:** Although FL avoids direct exposure of local data, new attack paradigms can still restore original data from uploaded gradients or public models [4].
- **Stealthy poisoning:** Due to the unpredictability of heterogeneous ERs, malicious nodes can “poison” the model with crafted uploads and manipulate the FL process [5].

Moreover, the distinctive characteristics of heterogeneous cross-silo IoRT, setting it apart from other FL scenarios, further complicate and necessitate specialized design considerations for FL security measures tailored to this context:

- Despite variations at the global scale, micro-level similarities also exist, exemplified by robots engaged in similar tasks possibly having datasets with closely approximated distributions. These resemblances can enhance learning performance even in heterogeneous environments.
- Likewise, most privacy obfuscation technologies focus on client-level differential privacy (DP) [6], which may introduce substantial redundant privacy loss due to high heterogeneity. In light of this, defenders can revisit privacy protection design to eliminate redundant loss [7].
- The advent of ISCC technologies provides ERs with great processing ability. In contrast, the computing and communication resources of the cloud server are relatively limited [8]. Hence, edge-assisted credit evaluation can improve both efficiency and accuracy, while avoiding potential bottlenecks for IoRT.

Taking the aforementioned characteristics of IoRT into account, in this paper, we attempt to jointly address the pluralistic security-related concerns and develop a moving target

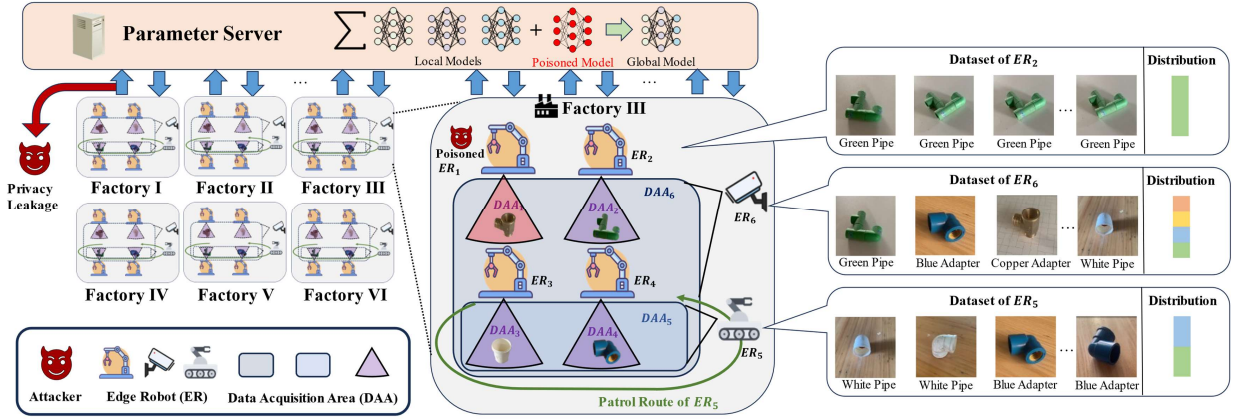


Fig. 1. Federated learning in heterogeneous cross-silo IoRT environments with privacy leakage and covert poisoning attacks

defense (MTD) approach for FL with Heterogeneous Cross-Silo IoRT, named FedMTD. By actively shuffling participant groups and changing the attack surface during the FL training process, FedMTD can protect data privacy and model integrity at the same time. The contributions are four-fold:

- To the best of our knowledge, we provide the first MTD solution for federated learning in heterogeneous cross-silo IoRT environments. By constructing a lightweight hierarchical clustering structure, privacy and integrity could be jointly ensured without too many compatibility requirements or excessively high overhead.
- During the iterative training process, FedMTD integrates a lightweight *cluster-specific sample-level* privacy-enhanced mechanism. Compared to widely-adopted client-level privacy obfuscation methods, the accuracy loss is significantly limited, especially for highly heterogeneous IoRT environments.
- To safeguard the integrity of FL tasks without compromising privacy, we further design a credit-based defense strategy, which can accurately thwart multiple poisoning attacks from stealthy malicious participants.
- Finally, the results of a series of experiments highlight the performance of FedMTD in terms of convergence, model accuracy, privacy cost, and poisoning mitigation in heterogeneous cross-silo scenarios.

The rest of this paper is structured as follows: the key problems of IoRT are analyzed in Section II. The framework of the proposed FedMTD is described in Section III, where the detailed description of cooperative training, privacy enhancement, and poisoning resistance are illustrated, respectively. The experimental results are presented in Section IV. Finally, Section V concludes this paper and briefly discusses the future directions.

II. FEDERATED LEARNING IN IORT: KEY PROBLEMS

A. Heterogeneity

As IoRT continues to expand, it becomes evident that edge robots are prone to exhibit highly diverse distributions. This diversity can result in accuracy drop or convergence failure.

In recognition of this challenge, a straightforward solution is *personalized FL*, which reconciles divergence between edge participants and central aggregator by adding local-adaptive

regularization term into the iterative training process [9]. PWFL [3] proposes a comprehensive task scheduling algorithm that employs proximal policy optimization to find an optimal task scheduling policy in automated warehouses with heterogeneous autonomous robotic systems. SCAFFOLD [10] further adopts variance reduction to correct client drift derived from non-iid data distributions.

Cluster-based FL is another approach to alleviate the performance degradation caused by heterogeneity. Aggregators can adopt pairwise cosine similarity between gradients, data distribution distance, and other criteria to group participants and minimize the inner variance [8].

However, despite substantial research efforts invested in heterogeneous federated learning, there remains room for improvement. As shown in figure 3 (b), most existing methods compromise the convergence speed for better learning accuracy. As edge robots are distributed discretely over a large area, the number of communication rounds contributes the most to the efficiency bottleneck. Hence, in this paper, the cross-silo FL training process should be task-customizable, efficient, and accurate, even in highly heterogeneous IoRT environments.

B. Privacy

By exchanging global models and local gradients, FL avoids over-the-air transmission of sensitive data and thus significantly protects the privacy of participants. Hence, emerging privacy-focused attacks, such as deep leakage and gradient inversion, have revealed that sophisticated adversaries could still restore data samples and exploit sensitive information. Owing to the relatively small processing and communication overhead, differential privacy has emerged as a promising lightweight solution for IoRT environments. For example, DPFL [6] designs a novel differential private model with adaptive gradient descent algorithm. NbAFL [11] further proposes client random scheduling strategy and duplex privacy obfuscation to provide tight and provable privacy guarantees. To ensure an anonymized local model update and counter poisoning attacks, PPAFL[4] integrates with an Autoencoder and a Gaussian mechanism for smart mobile robotic applications.

Nevertheless, existing approaches often overlook the inherent similarities between edge robots, potentially resulting in redundant privacy obfuscation measures. This oversight can

lead to significant performance degradation. Therefore, we attempt to adjust the privacy definition by introducing cluster-specific sample-level DP, to better depict privacy needs.

C. Integrity

In the realm of countering poisoning attacks and preserving the integrity of FL, two predominant strategies have emerged: mitigation and detection.

Mitigation-based strategies focus on reducing the impact of malicious actors by modifying aggregation rules. For instance, the Krum [5] algorithm selects a subset of gradients based on Euclidean distance, but these methods are unable to entirely remove attackers from the distributed learning process.

In contrast, *detection-based* approaches evaluate uploaded gradients to eliminate identified malicious nodes. Effective methods include using validation datasets (e.g., Zeno [12]), employing a feedback loop (e.g., BaFFLe [13]), or utilizing pre-trained anomaly detection models like autoencoders [14].

It's worth noting that these methods are often designed for homogeneous FL scenarios with identical data distributions. In heterogeneous environments, distinguishing deliberate attacks, such as backdoor clients, from benign participants under uniform rules can prove challenging, making these methods less effective.

III. THE PROPOSED FEDMTD FRAMEWORK

A. FedMTD overview

In this section, we will introduce the main framework of FedMTD. Assuming there is a Parameter Server S and multiple factories collaborating to train a universal global model in the IoRT environment. All N ERs form the global robot set $R = \{ER_1, \dots, ER_N\}$, and each needs to identify local data with set $D = \{D_1, \dots, D_N\}$. The normalized distribution of local data categories for ER_i can be represented as P_i . ERs from different factories and production lines have different data distributions, while ERs with the same functionality have similar data distributions. In the cross-silo environment mentioned above, we designed FedMTD as a mobile target defense technology, integrating three modules to jointly solve the heterogeneity, privacy, and security issues of FL. The initial step involves the implementation of a cluster-based scheme, facilitating ER node shuffling—a critical element that underpins subsequent heterogeneous training, privacy protection, and poisoning resistance. ER nodes undergo clustering and regrouping, assuming the role of virtual nodes (VN) for participation in batch-based FL training. Similar to LSSM, MOTAG, and other MTDs [15], each ER node involved in collaborative learning tasks is assigned a dynamically changing credit value. This value serves the purpose of identifying compromised nodes during periodic regrouping and FL interactions. Concurrently, a modified differential privacy noisy obfuscation technique is introduced to enhance privacy guarantees. To be specific, the FedMTD framework comprises the following three core components:

B. Hierarchical cross-silo training process

To provide accurate and fast convergence performance for FL tasks among a mass of heterogeneous and cross-silo edge robots, we first realize the hierarchical cross-silo training process (HCTP) as the foundation of FedMTD. Besides, HCTP also facilitates our scheme to achieve tighter privacy loss and stronger poisoning resistance, which will be illustrated in the following two subsections, respectively.

As shown in figure 2, the workflow of proposed hierarchical cross-silo FL training for heterogeneous IoRT can be explained as follows:

Step 1 (Data acquisition): As depicted in figure 1, ER nodes yield and collect data samples through equipped sensors.

Step 2 (Local training): At each round's beginning, participant ERs calculate gradients with local datasets.

Step 3 (Similarity-based edge robot clustering): Then, heterogeneous ER nodes are divided into nearly homogeneous clusters, to facilitate privacy enhancement and active defense.

Step 4 (Virtual participant node construction): Next, based on specific requirements of current FL task, the composition of VN is determined. Subsequently, ER nodes are randomly selected to construct multiple VNs.

Step 5 (Global aggregation among virtual nodes): Finally, the parameter server aggregates the uploads from all VNs and generates the global model.

Steps 1 and 2 align with the vanilla FL process and require no further elaboration. The subsequent discussion will focus primarily on the detailed implementation and design rationale for steps 3, 4, and 5.

1) *Similarity-based edge robot clustering:* Firstly, considering the high heterogeneity among participants, edge nodes are divided into different clusters based on pairwise similarity to avoid hindering the convergence performance of FL tasks. The *similarity* between ER nodes can be quantified based on *gradient distance, gradient direction, or data distribution* [8]. For simplicity, unless indicated otherwise, we only adopt *data distribution* in the experimental part of this paper. The similarity between distributions is denoted by the L_1 norm.

2) *Virtual participant node construction:* Once ER clusters are determined, the system selects nodes to construct VNs that are closest to the task objectives. As illustrated in figure 2, ER nodes with different data distributions can be combined to minimize the difference between their general data distributions and the target data distribution, subject to multiple constraints (the maximum number of nodes in each cluster). It's evidently an NP-hard problem. The problem can be transformed into a multi-dimensional knapsack problem with upper-bound constraints, further simplified into a 0-1 integer programming problem through binary splitting. Therefore, FedMTD employs heuristic algorithms for efficient solving.

As shown in figure 2, by building VNs close to the global data distribution and isomorphic to each other, and conducting the global aggregation of FL in VN as a unit, HCTP handles the heterogeneity problem in federated learning. Specifically, large-scale FL tasks are difficult to converge because of the huge number of ERs and their heterogeneity in the IOR environment. HCTP reduces the aggregation range of heterogeneous ERs to small batch VNs, alleviates the gradient conflict,

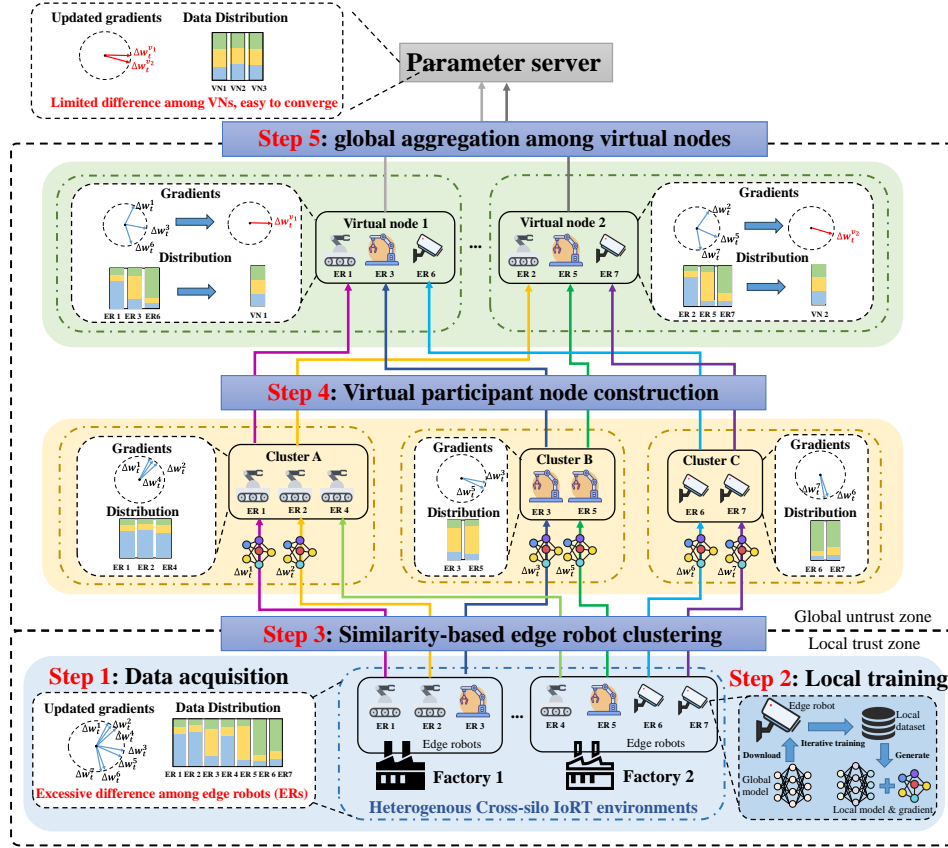


Fig. 2. Design of Hierarchical cross-silo training process for heterogeneous IoRT

and realizes the effective convergence of the FL models. VNs not only have an approximate global data distribution but also have similar gradient directions. Therefore, the global aggregation process of hierarchical cross-silo training will become the aggregation process between isomorphic VNs to achieve the safe and robust aggregation of global models, so as to solve the problem of FL heterogeneity in the IoRT environment.

As revealed by the results in figure 3 (a), it can be observed that finer-grained clustering can better limit the diversity within ER clusters, resulting in a smaller differential privacy sensitivity. However, more clusters may lead to a reduced number of ERs within each cluster, potentially undermining privacy amplification and anonymity in return. The number of cluster centers has a similar effect on VN construction: too few clusters may expand the gap between VNs and optimization objectives, while too many clusters limit the number of virtual nodes that can participate in computations due to insufficient ER nodes within each cluster. Therefore, it is evident that an optimal number of cluster centers exists from both privacy and learning performance perspectives. In FedMTD, we empirically approximate this number with iterative functions.

3) *Global aggregation among virtual nodes:* In the global aggregation phase, the parameter server dynamically scales the number of VNs based on specific task requirements, resource consumption, network communication quality, and other conditions. This optimization process aims to create a cost-effective and highly efficient FL deployment that provides satisfactory learning quality while minimizing unnecessary

overhead. Furthermore, post-global aggregation, the parameter server distributes the global model to participating ER nodes in the next communication round. It also dynamically adjusts or updates the resource utilization of VNs, including local iteration epochs, data scale, and sizes of VNs, based on the regret between the model's current state and task objectives. In addition, incentive mechanisms designed based on game theory principles, such as contract theory, can also be incorporated. These mechanisms can be distributed to participating ER nodes through a feedback loop, motivating ER nodes with highly related data to actively engage in the collaborative computation of FL. This approach can facilitate reaching the optimal matches between multiple different FL tasks and potentially suitable ER nodes.

C. Cluster-specific sample-level privacy-enhanced mechanism

Despite the emergence of various differentially private federated learning algorithms, the excessive privacy overhead resulting from heterogeneous environments remains inadequately addressed. In this part, cluster-specific sample-level privacy-enhanced mechanism (CSPM) focuses on the clustering characteristics of ER nodes in IoRT scenarios, building upon HCTP to further minimize performance losses under the same privacy protection level. To be specific, our privacy enhancement mechanism comprises two modules:

1) *Privacy enhancement:* As the heterogeneity among participants in large-scale cross-silo distributed learning can grow exponentially, leading to substantial accuracy losses when applying conventional DP algorithms, researchers have been

exploring alternative DP obfuscation approaches. For instance, in contrast to most *client-level* DP solutions, Liu *et al.* introduced a *silos-specific sample-level* DP concept in [7], where distinct privacy budgets (i.e., ϵ_i) are assigned to different silos i to facilitate efficient and personalized FLs. Nevertheless, these approaches usually overlook the micro-level similarities among ER nodes. Consequently, we propose the *cluster-specific sample-level* DP method tailored for heterogeneous cross-silo FL. Instead of altering the privacy guarantees, our solution focuses on limiting the sensitivity of gradients from protected ER nodes, achieving a similar level of privacy enhancement. To attain this objective, we employ a “divide and conquer” method to handle heterogeneous ER nodes into different homogeneous clusters, providing an upper bound of obfuscation degree within each cluster. In essence, CSPM significantly reduces the scale of injected noise, promoting the deployment of FedMTD in heterogeneous environments. However, it is important to note that this enhancement comes at the expense of a reduced anonymity space. Therefore, cluster radii should be carefully adjusted to strike the optimal balance between sensitivity and cluster size, particularly in less heterogeneous scenarios.

2) *Privacy amplification*: Although the previous part enhances the privacy performance by reducing sensitivity Δf with cluster-specific similarity. We further find that, under the same obfuscation operation and protected target features (e.g., Δf), FedMTD can still amplify the privacy protection level. Due to the superior compatibility of our hierarchical cross-silo structure, the following modules could be integrated jointly:

Random subsampling uniformly selects a certain proportion of ER nodes within each cluster to participate in every round, instead of involving all nodes every time. Consequently, the uncertainty regarding whether a specific node’s data exists in the set is further increased, reducing the demand for obfuscation under the same privacy budget. Additionally, this approach can reduce the probability of ER nodes engaging in excessive consecutive rounds of computation, partially alleviating the straggler effect caused by differences among the processing capabilities of participants.

Partial concealing usually modifies shared gradients to amplify the privacy. Since most FL tasks employ iterative stochastic gradient descent algorithms for updates, retaining the direction of gradients (i.e., the sign of gradients) is sufficient to guarantee the correct training while significantly reducing interaction overhead and the risk of sensitive information leakage. Nevertheless, everything comes with a price. The convergence rate of FL may be reduced. Therefore, it is better suited for scenarios involving large-scale complex models. In FedMTD, partial concealing is nullified by default; it is only activated when the dimensionality of the model exceeds a certain amount.

D. Compound active defense strategy

After achieving satisfactory FL performance and privacy protection levels, FedMTD further develops a compound active defense strategy. Based on ER nodes’ distributed dynamic credit evaluation, both malicious behaviors and attackers are effectively thwarted.

1) *Credit evaluation*: To accurately and efficiently identify the malicious participants (i.e., compromised ER nodes), the credit evaluation includes two modules:

Evaluation criteria define the maliciousness of participant behaviors. Specifically, the criteria can be divided into two parts: *cosine similarity-based* and *auxiliary verification-based* modules. As the name suggests, the former regards the gradients uploaded from ER nodes as multi-dimensional weighted vectors, forming an evaluation metric through the computation of normalized cosine similarity between vectors. Additionally, the latter makes full use of publicly available data. It samples a small auxiliary verification dataset from non-sensitive global public datasets based on the data distribution information of targeted clusters, facilitating rapid model validation.

It is worth noting that defenders can also set credit evaluation functions arbitrarily according to personal security demands or preferences. Besides, more complex AI-empowered methods could also be adopted, e.g., outlier detection models based on unsupervised learning, graph neural networks, or variational autoencoder.

Evaluation method first calculates four types of credit evaluation values with decaying factors: instantaneous intra-cluster, instantaneous inter-cluster, long-term intra-cluster, and long-term inter-cluster. Subsequently, the inverse-entropy weighting method is employed for the fusion calculation of the above multi-variate credits. Ratings of pairwise credit provided by ERs or clusters with higher similarity are considered more valuable references and are therefore assigned greater aggregation weights. In this mechanism design, attackers are unable to exploit global heterogeneity to conceal deviations between local malicious and benign behaviors.

Furthermore, considering the emergence of large-scale, structurally complex FL tasks, such as those involving AIGC, it is not necessary to utilize information from all dimensions of gradients for identification. For example, one can opt to use only the neural network’s neuron weights before the final softmax layer or employ sparse sampling to select a few significant parameters.

2) *Active defense*: By evaluating the cumulative scores of participant ER nodes, we can effectively cease covert poisoning attacks through active elimination. Our design for eliminating poisoning attacks serves two key purposes. Firstly, it aims to nullify the influence of malicious activities on the FL training process, a focus shared with existing mitigation-based countermeasures like robust FL approaches. Secondly, it seeks to directly remove malicious participants, aligning with the objectives of existing detection-based countermeasures. Consequently, FedMTD enables swift responses to poisoning behaviors, ensuring the maintenance of high FL performance. Additionally, this proactive approach prevents futile participation by compromised nodes, mitigating the risk of resource-exhausting attacks stemming from persistent adversaries. Our active defense strategy encompasses both soft and hard elimination measures:

Soft elimination mainly controls two key parameters: aggregation weights in the global aggregation rule (GAR) and ER selection probabilities. Based on the comprehensive credit assessment values calculated as mentioned above, highly

trustworthy ER nodes will experience an increased probability of selection during the construction of VNs in the subsequent round. Conversely, ER nodes with lower credit ratings will experience a reduction in the probability of being chosen for computational tasks. Similarly, highly trustworthy ER nodes who engage in collaborative computations will have their uploaded gradients assigned larger aggregation weights. All weights and selection probabilities are normalized to ensure compliance with the requirements of FL global aggregation and differential privacy subsampling definitions. Furthermore, in addition to assessing the participating edge nodes, FedMTD considers the phenomenon where malicious data or gradients from existing poisoning attackers often exhibit exceptionally high similarities. Highly homogeneous ER clusters with anomaly maliciousness will be removed. It's worth noting that due to the highly heterogeneous environment, adopting different aggregation weights can slow down convergence. Therefore, in FedMTD's global aggregation phase, the aggregation weights for each VN are averaged.

Hard elimination sets a hard threshold thr to locate and purge the identified adversaries out of the IoRT. As the ERs are randomly selected and reassigned to new VN groups, the accumulated credit values are distinctively different between compromised and benign ER nodes after long enough communication rounds. Therefore, based on Cantelli's inequality, thr is set to be the median value μ of credit evaluation results, adding weighted standard deviation term $\lambda\sigma$ and a bias, as $thr = \mu + \lambda\sigma + bias$, which provides false positive upper bound asymptotic to $O(1/(1+a^2))$. Apparently, λ can balance the poisoning resistance ability and the false recognition rate of FedMTD. With larger λ , the number of ill-treated benign ER nodes is minimized at the expense of robustness decreasing.

The primary computational load of FedMTD lies on the server side, resulting in a relatively minor impact on resource-constrained robot terminals. Moreover, the main computational load of FedMTD is associated with cosine-similarity calculations. Assuming the current epoch involves the selection of n clients and the model parameter quantity is denoted as B , the computational time complexity is $O(n^2B)$. It has a limited impact on server with high computational capabilities.

IV. EXPERIMENTS

To verify the performance of FedMTD, we have established an experimental platform consisting of a workstation equipped with an Intel i9-10940X CPU, 128GB of RAM, and two RTX3090 24GB GPUs. All software components are developed using Python 3.8 and rely on the PyTorch library. We adopted the MNIST dataset as our benchmark dataset. The MNIST dataset encompasses 10 handwritten digit classes, with all images standardized to 28x28 pixels. We allocated data in an 8:2 ratio to construct our training and validation sets. For the learning task on the MNIST dataset, we employed a four-layer CNN model comprising two 5x5 convolutional layers, one fully connected layer, and a softmax output layer. In our experimental setup, we constructed a heterogeneous federated learning environment, consisting of 200 ER nodes and a parameter server. The global communication involves

200 rounds. For the local training of each edge node, we chose a training epoch of $E = 1$, a batch size of $b = 128$, a learning rate of $\eta = 0.1$, and a momentum value of $m = 0.9$.

A. Convergence performance

We first examine the convergence performance of the proposed FedMTD. Figure 3 (a) illustrates the maximum L_1 norm between the mean value (clustering center) and arbitrary node's distribution within each ER cluster. With more clusters, the differences can be reduced, which paves the way for privacy enhancement and accurate adversary identification. Then, we choose the vanilla FedAvg algorithm and personalized FL method SCAFFOLD [10] as benchmarks in Figure 3 (b). The number of clusters has a significant impact on the convergence performance of the global model. Too few clusters make it difficult to build a VN that fits the global data distribution in a fine-grained way. The resulting data drift shifts the VN's optimization objectives and reduces the performance of the global model. However, too many clusters will reduce the participation rate of ERs because the number of ERs in the cluster is not enough to build VN. The overfitting caused by too few training samples will reduce the generalization ability of the global model. The results in Figure 3 (b) suggest that FedMTD can provide a significantly faster convergence rate and better test accuracy compared with existing countermeasures. It is reasonable to expect greater improvements with more heterogeneous environments and more complex models.

B. Privacy enhancement

Subsequently, to elucidate the privacy enhancement achieved by FedMTD, we further compare the average accuracy of FedMTD with two widely adopted privacy-preserving federated learning mechanisms, NbAFL [11] and DPFL [6]. All edge robot nodes construct local heterogeneous datasets following Dirichlet distribution $Fir(\alpha)$ with $\alpha = 0.05$. As depicted in figure 3 (c), the red line (FedMTD) can achieve better learning results with less than 5% loss, while the green line (NbAFL) and orange line (DPFL) suffer about 17% and 55% loss, respectively. Under the same privacy protection level ($\epsilon = 1$), DPFL tailors the gradient directly via median clipping bound, which introduces huge yet unnecessary Gaussian noise. NbAFL can reduce the noise without sacrificing ϵ by duplex noise injection. But the sensitivity of protected gradients is still overlarge. In contrast, FedMTD limits the variance within each robot group via clustering, and eliminates the redundant obfuscation operations accordingly.

C. Poisoning resistance

Finally, we scrutinize the resistance ability of FedMTD against poisoning attacks involving heterogeneous participants, gauged through the average recognition rate after 100 communication rounds. Unless otherwise specified, the proportion of malicious ER nodes is set at 20%. As Byzantine-robust FLs like Krum or GeoMed can only mitigate the influence of malicious uploads, we adopt three detection-based defense schemes as benchmarks: Zeno [12], BaFFLe [13], and an

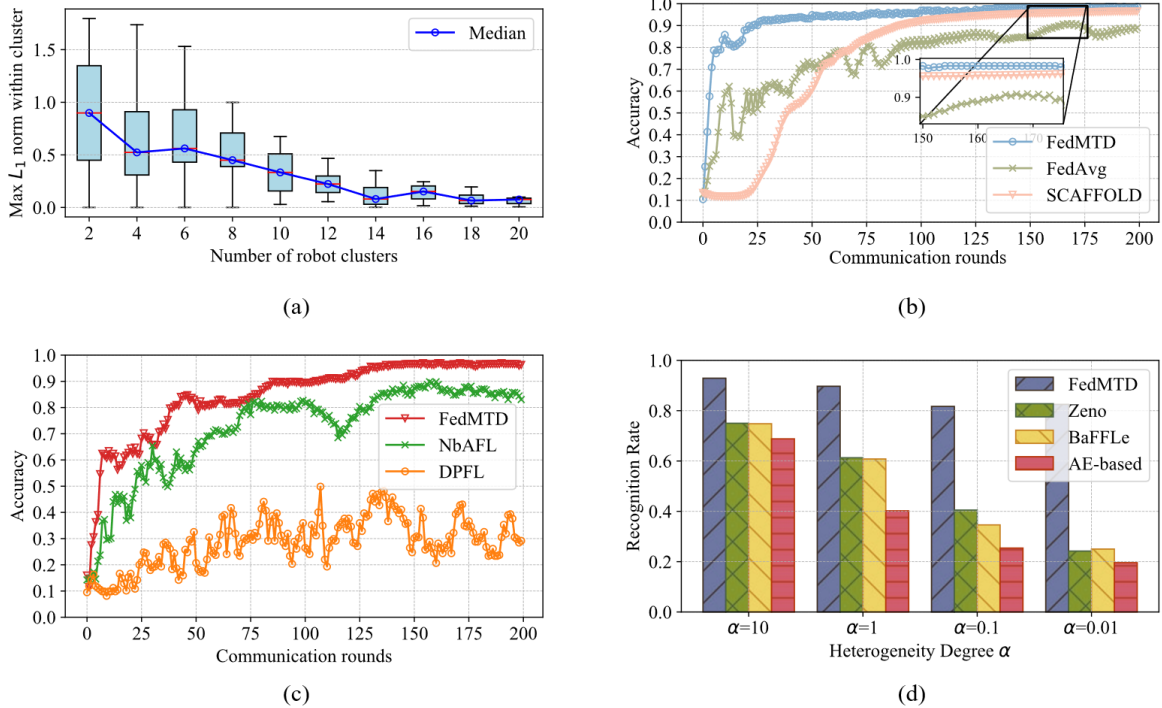


Fig. 3. The evaluation performance of FedMTD: a) heterogeneity minimization; b) average convergence rate and accuracy; c) training performance with privacy noise injection ($\epsilon = 1$); d) average recognition rate of malicious participants.

Autoencoder-based method [14]. Similarly, the data of heterogeneous participants are sampled following $Dir(\alpha)$ from MNIST. Higher value of α denotes an increased degree of heterogeneity. As illustrated in Figure 3 (d), the majority of the defense methods perform admirably when dealing with data distributions that exhibit approximations ($\alpha > 1$). In stark contrast, only FedMTD sustains an accuracy rate exceeding 0.8 when faced with highly heterogeneous participants ($\alpha < 0.1$). This underscores the precision of our defense strategy in identifying malicious adversaries, even in highly heterogeneous environments, while other methods experience a sharp decline.

V. CONCLUSIONS

This paper delves into the intricate security challenges inherent in federated learning within cross-silo Internet of Robotic Things environments. It explores the critical issues of heterogeneity, privacy breaches, and model poisoning. To address these concerns, we introduce an innovative approach known as FedMTD, which employs a moving target defense strategy to simultaneously safeguard data privacy and model integrity during federated learning training, particularly in highly heterogeneous multi-robot systems. Moreover, a series of experiments on publicly available datasets are conducted to assess the effectiveness and feasibility of FedMTD. In terms of future directions, to meet the escalating demands for large-scale model training, we aim to enhance FedMTD by extending its capabilities to encompass federated split learning (FSL) and federated transfer learning (FTL) within heterogeneous environments. Furthermore, given the ever-evolving attack paradigms, more complicated and AI-enhanced advanced persistent threat (APT) models should also be considered for future proactive defense strategies.

ACKNOWLEDGMENT

This work is supported by the National Natural Science Foundation of China (NSFC) under grants No. 62225105, 62394323, and by the Beijing Natural Science Foundation under grant No. 4244084. **G.-M. Muntean wishes to acknowledge the Science Foundation Ireland support via grants 21/FFP-P/10244 (FRADIS) and 12/RC/2289_P2 (INSIGHT).**

REFERENCES

- [1] G. Zhao, P. Zhang, Y. Shen, L. Peng, and X. Jiang, "Passive user authentication utilizing two-dimensional features for iiot systems," *IEEE Transactions on Cloud Computing*, vol. 11, no. 3, pp. 2770–2783, 2023.
- [2] G. Zhao, P. Zhang, Y. Shen, and X. Jiang, "Passive user authentication utilizing behavioral biometrics for iiot systems," *IEEE Internet of Things Journal*, vol. 9, no. 14, pp. 12783–12798, 2021.
- [3] T. M. Ho, K.-K. Nguyen, and M. Cheriet, "Federated deep reinforcement learning for task scheduling in heterogeneous autonomous robotic system," *IEEE Transactions on Automation Science and Engineering*, pp. 1–13, 2022.
- [4] X. Zhou, W. Liang, K. I.-K. Wang, Z. Yan, L. T. Yang, W. Wei, J. Ma, and Q. Jin, "Decentralized p2p federated learning for privacy-preserving and resilient mobile robotic systems," *IEEE Wireless Communications*, vol. 30, no. 2, pp. 82–89, 2023.
- [5] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine learning with adversaries: Byzantine tolerant gradient descent," *Advances in neural information processing systems*, vol. 30, 2017.
- [6] Y. Li, S. Wang, C.-Y. Chi, and T. Q. Quek, "Differentially private federated learning in edge networks: The perspective of noise reduction," *IEEE Network*, vol. 36, no. 5, pp. 167–172, 2022.
- [7] K. Liu, S. Hu, S. Z. Wu, and V. Smith, "On privacy and personalization in cross-silo federated learning," *Advances in Neural Information Processing Systems*, vol. 35, pp. 5925–5940, 2022.
- [8] Z. Li, Y. He, H. Yu, J. Kang, X. Li, Z. Xu, and D. Niyato, "Data heterogeneity-robust federated learning via group client selection in industrial iot," *IEEE Internet of Things Journal*, vol. 9, no. 18, pp. 17844–17857, 2022.
- [9] X. Yuan and P. Li, "On convergence of fedprox: Local dissimilarity invariant bounds, non-smoothness and beyond," *Advances in Neural Information Processing Systems*, vol. 35, pp. 10752–10765, 2022.

- [10] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International conference on machine learning*. PMLR, 2020, pp. 5132–5143.
- [11] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. Quek, and H. V. Poor, "Federated learning with differential privacy: Algorithms and performance analysis," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 3454–3469, 2020.
- [12] C. Xie, S. Koyejo, and I. Gupta, "Zeno: Distributed stochastic gradient descent with suspicion-based fault-tolerance," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6893–6901.
- [13] S. Andreina, G. A. Marson, H. Möllering, and G. Karame, "Baffle: Backdoor detection via feedback-based federated learning," in *2021 IEEE 41st International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2021, pp. 852–863.
- [14] Z. Zhou, C. Xu, M. Wang, X. Kuang, Y. Zhuang, and S. Yu, "A multi-shuffler framework to establish mutual confidence for secure federated learning," *IEEE Transactions on Dependable and Secure Computing*, 2022.
- [15] Z. Zhou, X. Kuang, L. Sun, L. Zhong, and C. Xu, "Endogenous security defense against deductive attack: When artificial intelligence meets active defense for online service," *IEEE Communications Magazine*, vol. 58, no. 6, pp. 58–64, 2020.

Sizhe Huang (swhsz@bupt.edu.cn) is currently pursuing a Bachelor's degree in the School of Computer Science, BUPT. His research interests include network security and active defense.

BIOGRAPHIES

Zan Zhou (zan.zhou@bupt.edu.cn) received the Ph.D. degree from Beijing University of Posts and Telecommunications (BUPT), Beijing, China in 2022. He is currently a Postdoctoral Fellow at the School of Computer Science, BUPT. His research interests include Network Security, Artificial Intelligence Privacy, and Active Defense.

Changqiao Xu (cqxu@bupt.edu.cn) received the Ph.D. degree from the Institute of Software, Chinese Academy of Sciences in Jan. 2009. Currently, he is a Professor with the State Key Laboratory of Networking and Switching Technology, and Director of the Network Architecture Research Center at BUPT. His research interests include Network Security, Mobile Networking, Multimedia Communications, and Future Internet Technology. He is currently serving as the Editor-in-Chief of *Transactions on Emerging Telecommunications Technologies* (Wiley).

Shujie Yang (sjyang@bupt.edu.cn) received the Ph.D. from the Institute of Network Technology, BUPT, Beijing, China, in 2017. He is an associate professor with the State Key Laboratory of Networking and Switching Technology, BUPT. His major research interests include active defense, wireless communications, and wireless networking.

Xiaoyan zhang (xiaoyan@bupt.edu.cn) is a Professor and Ph.D. Supervisor with the School of Computer Science (National Pilot Software Engineering School), BUPT. Her research interests include Network Security, Mobile Networking, Multimedia Communications, and Future Internet Technology.

Hongjing Li (lihongjing@bupt.edu.cn) received the B.E. degree in computer science and technology from the School of Information, Beijing Forestry University, in 2022. She is currently pursuing an M.Eng degree in the School of Computer Science, BUPT. Her research interests include federated learning and privacy computing.

Gabriel-Miro Muntean (gabriel.muntean@dcu.ie) is a Professor with the School of Electronic Engineering, Dublin City University (DCU), Ireland, and the Co-Director of DCU Performance Engineering Laboratory. His research interests include rich media delivery quality, performance, and energy-related issues, technology-enhanced learning, and other data communications in heterogeneous networks. He is an Associate Editor of the *IEEE Transactions on Broadcasting*, and the Area Editor of the *IEEE Communications Surveys and Tutorials*. He coordinated the EU Horizon 2020 project NEWTON and led the DCU team in the EU Horizon 2020 projects **TRACTION** and **HEAT**. He is a Fellow of the IEEE and IEEE Broadcast Technology Society.