

AdRo-FL: Informed and Secure Client Selection for Federated Learning in the Presence of Adversarial Aggregator

Md. Kamrul Hossain*, Walid Aljoby*, Anis Elgabri*, Ahmed M. Abdelmoniem[§], Khaled A. Harras[‡]

*College of Computing and Mathematics, King Fahd University of Petroleum & Minerals

[§]School of Electronic Engineering and Computer Science, Queen Mary University of London

[‡]Department of Computer Science, Carnegie Mellon University

Abstract—Federated Learning (FL) enables collaborative model training without exposing clients' local data. While clients only share model updates with the aggregator, studies reveal that malicious aggregators can infer sensitive information from these updates. Secure Aggregation (SA) protects individual updates during transmission; however, recent work demonstrates a critical vulnerability where adversarial aggregators manipulate client selection to bypass SA protections, constituting a Biased Selection Attack (BSA). Although verifiable random selection prevents BSA, it precludes informed client selection essential for FL performance.

We propose Adversarial Robust Federated Learning (AdRo-FL), which simultaneously enables: (1) informed client selection based on client utility, and (2) robust defense against BSA maintaining privacy-preserving aggregation. AdRo-FL implements two client selection frameworks tailored for distinct settings in FL. The first framework assumes clients are grouped into clusters based on mutual trust, such as different branches of an organization. The second framework handles distributed clients where no trust relationships exist between them. For the cluster-oriented setting, we propose a novel defense against BSA by (1) enforcing a minimum client selection quota from each cluster, supervised by a cluster-head in every round, and (2) introducing a client utility function to prioritize efficient clients. For the distributed setting, we design a two-phase selection protocol: first, the aggregator selects the top 80% of clients based on our utility-driven ranking; then, a verifiable random function (VRF) ensures transparent, auditable, and BSA-resistant final selection from this pool. AdRo-FL also applies quantization to reduce communication overhead and sets strict transmission deadlines to improve energy efficiency. AdRo-FL achieves up to $1.85\times$ faster convergence (time-to-accuracy) and up to $1.06\times$ higher final accuracy compared to insecure baselines across multiple datasets.

Index Terms—Federated learning, Informed Client selection, Biased selection attack, Secure aggregation, Verifiable random function

I. INTRODUCTION

FEDERATED Learning (FL) [1] has gained significant attention as a decentralized paradigm for collaborative machine learning, where multiple clients train a shared model without exchanging raw data. This approach is especially suitable for privacy-sensitive applications such as healthcare, finance, and industrial IoT (IIoT), where data confidentiality is paramount. In FL, clients perform local training on their private data and periodically send model updates to a central aggregator, which aggregates them to refine the global model. While FL enhances privacy by design, it remains susceptible

to adversarial threats, particularly when the central aggregator behaves maliciously.

One of the critical vulnerabilities of FL is its exposure to aggregators, which can exploit model updates to extract sensitive client information [2]–[4]. Several privacy-preserving techniques have been proposed to mitigate this risk, including Differential Privacy (DP) [5], [6] and Secure Aggregation (SA) [7]. DP introduces noise into model updates before transmission, providing strong privacy guarantees. However, this noise reduces model accuracy, making it less suitable for scenarios that require high-performance learning. SA, on the other hand, enables clients to encrypt their updates before transmission, ensuring that only aggregated results are revealed to the aggregator. This makes SA a more practical choice in many FL settings, particularly when privacy must be preserved without significantly compromising model performance.

Despite its advantages, SA itself is not entirely secure. Recent studies have shown that biased selection attacks (BSA) [8], [9] can undermine SA by allowing an adversarial aggregator to manipulate client selection, thereby inferring individual model updates. BSA can manifest in two forms: *non-colluding* attacks, where an adversarial aggregator isolates a specific client to obtain its model update, and *colluding* attacks, where an adversarial aggregator colludes with a subset of clients to uncover the model update of the target client. This vulnerability highlights the need for a resilient FL framework that prevents the exploitation of SA while maintaining computational efficiency as well as model quality.

To defend against BSA, researchers have recently proposed using verifiable random functions (VRF) [8], [9] to allow each client to independently verify whether the selection process was genuinely random. A technical overview of VRFs is provided in the appendix A. While truly random client selection can effectively disrupt BSA, this approach lacks the ability to make informed decisions when selecting clients. Informed client selection [10] remains crucial for achieving faster convergence and improving energy efficiency in FL.

In this paper, to address this problem, we propose Adversarial Robust Federated Learning (AdRo-FL) that allows informed client selection while protecting against biased selection attack by an adversarial aggregator. The development of AdRo-FL is motivated by several key challenges in FL. First, while SA prevents direct leakage of model updates, it remains vulnerable to BSA that allows adversarial aggregator to infer local model update. Second, existing defenses that

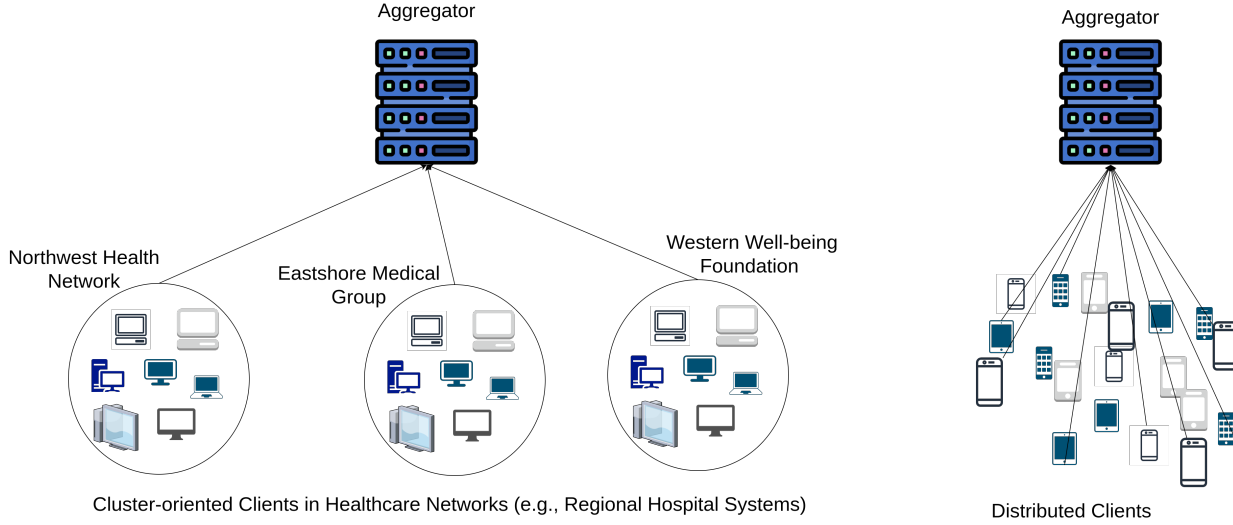


Fig. 1. Federated Learning client settings.

rely on random client selection often degrade performance by disregarding the quality of client contributions. There is a notable lack of solutions that effectively balance resilience to BSA with optimized client selection based on data utility and system efficiency. Lastly, many real-world deployments, such as smart cities, industrial IoT, and finance, exhibit natural client clustering based on organizational and administrative trust, which remains underutilized in current approaches.

AdRo-FL considers two different client settings commonly found in real-world environments. The first setting assumes that clients belong to clusters where there exist mutual trusts within, such as different branches of an organization (Fig. 1, Healthcare Networks). The second setting assumes that the clients are distributed, non-cluster-oriented and that there is no trust between the clients (Fig. 1, Cellphone Users).

For the cluster-oriented setting, we assume a cluster-head per cluster that intermediates communications between the cluster and the aggregator. We enforce that either a minimum number C of clients must be selected from a cluster or no clients will join from that cluster. That is, either the aggregator will select $\geq C$ clients from a cluster or none. The cluster-head will supervise and verify this process in every training round. This process guarantees protection against BSA and is detailed in section V-A6.

As for the distributed, non-cluster-oriented setting, the client selection is completed in two phases. In the first phase, clients compute a utility value and sign it using their private keys. These signed values are sent to the aggregator, which ranks them, compresses the sorted list and signatures, and broadcasts both. Clients in the top 80% proceed to the next selection phase as detailed in section V-B. This two level client selection strategy ensures informed as well as secure client selection. The initial selection retains best clients based on their utility and the final selection ensures verifiable random selection to prevent BSA.

Beyond privacy protection, AdRo-FL enhances FL's efficiency through combining several key approaches:

- **Utility-based Client Selection:** It selects clients based

on a combination of local loss, gradient norm and sample size, prioritizing clients that contribute the most to global model.

- **Deadline-aware Selection:** To ensure timely updates and energy efficiency in resource-constrained IoT environments, it incorporates a deadline-based transmission mechanism, allowing clients to only participate within a specified time constraints.
- **Quantization for Communication Efficiency:** To reduce the communication overhead, it employs quantization, compressing model updates before transmission.

To evaluate the effectiveness of AdRo-FL, we conduct extensive experiments on four benchmark datasets: MNIST, FMNIST, SVHN, and CIFAR-10. We simulate non-IID data distribution using Dirichlet distribution [11], closely mirroring real-world FL scenarios. Our results demonstrate that AdRo-FL significantly improves training accuracy and energy efficiency compared to insecure client selection methods such as Oort [12] which is a state-of-the-art client selection method in FL.

Below we summarize our contributions:

- We propose AdRo-FL that supports both cluster-oriented and non-clustered client settings. In cluster-based environments, AdRo-FL leverages intra-cluster trust to defend against BSA while enabling informed client selection. For non-clustered, distributed setting, it introduces a two-phase client filtering mechanism to ensure secure and informed client selection.
- We propose a client utility measurement function combining gradient norm, local loss, and sample size. Additionally, clients are filtered based on a transmission deadline to improve energy efficiency. We apply 8-bit quantization to reduce communication overhead and enhance scalability in bandwidth-constrained IoT systems without affecting performance.
- We conduct comprehensive experiments on MNIST, FMNIST, SVHN, and CIFAR10 datasets by using Dirichlet

distribution with high heterogeneity, demonstrating superior performance over insecure client selection baselines.

II. THREAT MODEL

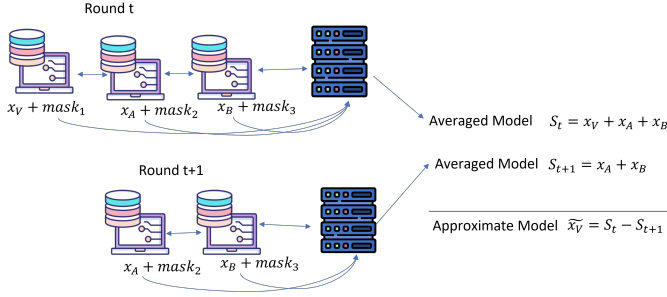


Fig. 2. Non-colluding biased selection attack.

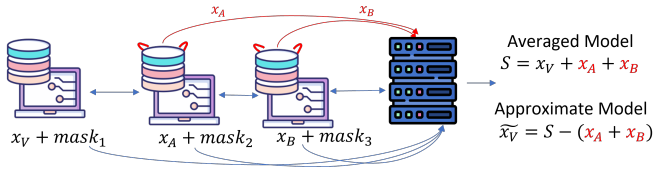


Fig. 3. Colluding biased selection attack.

This section describes how a BSA is executed in SA. In [8], two types of attacks are described. The first, a non-colluding attack (Fig. 2), assumes no collusion between clients and the aggregator, nor the creation of fake clients by the aggregator. Here, the aggregator uses statistical inference by subtracting consecutive aggregated updates. Suppose, in round t , the aggregator selects the victim client V and a subset of other clients (A and B), resulting in an aggregated update $S_t = x_V + x_A + x_B$. Here x represents the local model update. In the following round $t + 1$, the aggregator reselect the same clients, excluding V , obtaining $S_{t+1} = x_A + x_B$. By estimating $S_{t+1} - S_t$, the aggregator approximates V 's model x_V . However, the work [8] mentioned some preconditions to make this attack feasible, such as the aggregator should either send same global model to the selected clients at the rounds t and $t + 1$ or do this attack at the late stage when the global model has already converged. It was also assumed that the clients (other than the victim) do not change the algorithm and local data between the rounds t and $t + 1$.

The second type, a colluding attack (Fig. 3), relaxes these preconditions by assuming that some clients collude with the aggregator by sharing their model updates. Here, the aggregator aggregates the colluding clients' updates (x_A, x_B) alongside victim V 's, obtaining $S = x_V + x_A + x_B$. Since the aggregator knows x_A and x_B through colluding clients, it can approximate x_V using $x_V \approx S - (x_A + x_B)$.

III. RELATED WORK

A. Insecure Informed Client Selection

Client selection plays a crucial role in improving global model convergence and energy efficiency in FL. While random

client selection is conventionally used, researchers have explored more intelligent informed selection strategies. There are two dimensions of this. A group of works [13]–[15] explored ways to improve global model convergence time and accuracy without considering resource efficiency while another group of works explored improving both the convergence and training efficiency in terms of resource consumption [16]–[19].

Most of the works that employ client selection for improved global model convergence leverages metrics such as local loss, gradient norm and training duration. For instance, [20] prioritizes clients with the highest L2 norm of gradients, under the assumption that these clients contribute most to the global model's learning progress. However, relying solely on gradient magnitude may oversimplify client importance, neglecting critical aspects such as data diversity and quality, which are essential for robust FL. Similarly, [21] proposes the Power-of-Choice, where clients with higher local losses are favored to accelerate convergence. While this approach can improve accuracy, it introduces the risk of model bias, as clients with extreme loss values may disproportionately influence learning. Oort [12] also showed that selecting clients with higher training loss can accelerate convergence. FairDPFL-SCS [22] selects different subset of clients in each FL round while prioritizing the selection of new clients. The work [23] proposed an auction-based client selection mechanism for online FL to select better clients while compensating clients for their participation.

In addition to client selection for faster convergence, energy efficiency remains a major concern in FL, especially in resource-constrained environments such as IoT networks and mobile edge computing. Several strategies have been proposed to optimize energy consumption, including message compression, accuracy relaxation, and client dropout mechanisms. In [24], an early stopping mechanism is introduced to terminate training when the global model begins to oscillate among local optima, thereby reducing unnecessary communication rounds. Another approach in [25] utilizes 8-bit approximations instead of traditional 32-bit gradient representations, effectively reducing communication overhead without significantly compromising accuracy. Additionally, PyramidFL [26] employs a trade-off strategy that allows slightly lower accuracy in favor of better communication efficiency by training sub-optimal solutions during local model updates.

Yet another group of work leverages client data and resource homogeneity to cluster them together to facilitate improved client selection and model convergence. For example, the work [27] clustered clients based on their local model weights. [28] does client clustering based on location and data similarity. Another work called 'Auxo' [29], groups clients with statistically similar data distributions, called 'cohorts'. Yet another work [30] clustered clients by training duration. The clusters are subsequently scheduled across training rounds to minimize the total number of iterations needed to reach a desired accuracy level.

While the above mentioned client selection methods are effective in terms of accuracy and energy efficiency, they are insecure against an honest-but-curious aggregator and an adversarial aggregator.

B. Secure Uninformed Client Selection

Various techniques have been developed to allow an aggregator to obtain the sum of clients model updates without access to individual model update. However, many rely on assumptions that limit their practicality, such as an honest-but-curious aggregator or a trusted third party [31], [32]. Secure Aggregation (SA) is one of the few available protocols that were designed to handle a potentially malicious aggregator. SA [7] has been widely adopted as a countermeasure, enabling encrypted aggregation of client updates to obscure individual contributions. However, SA alone is insufficient, as recent research has demonstrated its susceptibility to biased selection attacks (BSA) [8]. These attacks allow a malicious aggregator to manipulate client selection in order to infer sensitive model updates, posing a significant threat to data privacy. Additionally, there is an inherent problem with SA, it can handle only a limited number of clients colluding with the aggregator; beyond this threshold, the adversary could reconstruct a client's update [7]. To strengthen SA, distributed DP [33] adds another layer of privacy to the computed sum by combining DP noise addition with SA, offering robust privacy guarantees. Nevertheless, both SA and distributed DP remain susceptible to BSA when a majority of participants are dishonest. To mitigate BSA, [8] introduces a blockchain-based random client selection strategy, leveraging blockchain's transparency to ensure unbiased client participation. While this approach enhances trustworthiness, it introduces scalability challenges due to blockchain's inherent limitations in transaction speed and data throughput, making it impractical for large-scale FL deployments. Additionally, random client selection overlooks key factors such as client utility which makes informed client selection impossible.

C. Secure and Informed Client Selection

The Lotto [9] approach uses VRF to randomly select clients from a client-pool. The clients in the client-pool were chosen based on some desired client selection criteria where an adversarial aggregator was involved. One of the limitations of this work is that the adversarial aggregator may attempt to manipulate the process of building the client-pool by excluding honest clients which will enable a BSA. Moreover, it fails when the majority of clients are colluding while only the minority are honest. Another limitation of this work is that when the aggregator does a primary selection, it goes without a verification from client side before applying verifiable randomness based final selection. The adversarial aggregator may manipulate this to do a BSA by sending a forged list to the clients. Moreover, the VRF α was chosen to be a deterministic number which gives the adversarial aggregator more control to create Sybils [34] in the selection pool. To the best of our knowledge, this remains the only work specifically addressing BSA in SA, underscoring the need for more research on this area. Our proposed AdRo-FL offers several advantages over Lotto. Firstly, we do not use deterministic α value for the VRF (as Lotto does) which makes it more difficult for the adversarial aggregator to do BSA. Secondly, we include sign verification for the first level selection of clients

which makes the adversarial aggregator unable to forge the utility based selection process. Lotto does not address this vulnerability. Lastly, Lotto did not consider cluster-oriented client setting as a separate problem. We identify the existence of cluster-oriented clients setting and take advantage of the inherent trust of clusters to prevent BSA by an adversarial aggregator.

To summarize, existing works tend to focus on either privacy or model convergence and energy efficiency, but very few address all these challenges simultaneously. In real-world FL deployments, these factors are deeply interconnected, particularly in highly heterogeneous and resource-limited environments.

IV. AdRo-FL OVERVIEW

AdRo-FL aims to defend against BSA attempted by an adversarial aggregator while allowing informed client selection for improved global model and faster convergence. It takes advantage of the client's local loss and weighted gradient norm to determine client's utility. In addition, clients should be able to transmit their payloads within a specific deadline which helps faster transmission as well as less energy consumption. In addition, AdRo-FL incorporates quantization to reduce the payload size, improving communication efficiency in bandwidth-constrained networks. AdRo-FL implements two different frameworks to suit different client settings of real world environment to defend against BSA.

- 1) Informed client selection for cluster-oriented clients.
- 2) Informed client selection for non-cluster-oriented, distributed clients.

In Section V, we discuss the details of the client privacy part, which is our defense approach against the BSA in the two settings mentioned above.

In Section VI, we discuss the details of the optimization part, which includes the client utility that allows informed client selection, deadline-based client selection as well as quantization of the client's payload.

V. AdRo-FL: DEFENSE AGAINST THE BSA

A. Informed client selection for cluster-oriented clients

This section describes a real-world client setting in which clients are grouped under a unit that we call a cluster. There is trust within each cluster but not outside. That is, a client from one cluster does not trust a client from another cluster. One of the advantages we get in this setting is that we don't need to use randomness to prevent BSA from adversarial aggregator. We leverage the intra-cluster trust to prevent BSA. We describe this framework in detail in the following subsections.

1) *Rationale for Clustering*: In many real-world deployments of federated learning, such as within corporate ecosystems, municipal infrastructures, or financial institutions, clients often exhibit natural grouping based on ownership or administrative control. For example, Smart city applications might cluster clients under municipal departments (e.g. traffic, utilities, public safety), Industrial IoT settings could group devices deployed across different factory floors of the same company, and healthcare networks might cluster devices or

data silos within the same hospital system. A real example is the federated learning initiative for breast density classification involving seven clinical institutions worldwide [35]. Each institution (cluster) trained local models on their proprietary data and shared model updates with a central aggregator. This collaborative approach allowed the development of a robust global model without exchanging sensitive patient data, adhering to privacy regulations such as HIPAA and GDPR [36]. These groupings provide a natural trust boundary where clients within the same cluster are assumed to be under a common administrative domain.

2) *Trust Assumptions*: Although clients in the same cluster do not share raw data or model updates with each other, they trust that none of them will collude with the aggregator or at least there will be a minority trusted clients within each cluster. This assumption is vital for defending against BSA under SA. It ensures that the privacy threshold mechanism—requiring at least C clients from each participating cluster—functions effectively which we describe in section V-A6.

3) *Cluster Formation Mechanism*: Cluster formation in our setting is declarative and administrative rather than algorithmic, with clusters defined during deployment based on existing organizational structures. Each client is assigned a signed cluster identifier by a central certificate authority (CA) or a trusted cluster head. This identifier supports cluster membership verification and enforces a privacy threshold that ensures a minimum number of clients participate in each aggregation round. This aligns with real-world workflows where devices are centrally registered, authenticated, and managed.

4) *Example Scenario*: Consider a logistics company with multiple regional warehouses. Each warehouse runs a set of IoT devices monitoring inventory and environmental conditions. These devices are grouped into a cluster per warehouse. Although the devices do not exchange data with one another, they are managed by the same IT team and operate within the same secure network. By enforcing the a minimum client participation rule and monitoring selection behavior via a coordinator (cluster head), the system ensures robust protection against an adversarial aggregator that might try to isolate a single device for BSA.

5) *Cluster and Aggregator Communication*: Fig. 4 presents a high-level overview of the framework, illustrating the interaction flow between the aggregator, cluster heads, and clients. The communication and decision-making process are designed to ensure client privacy, prevent BSA. The process goes through in five main steps, as detailed below:

Step 1) Global Model Broadcast: The aggregator shares the latest global model with all cluster heads. Each cluster head, responsible for coordinating communication within its cluster, distributes this model to the local clients. This ensures that all clients work on a consistent starting point for their local training.

Step 2) Client Metadata Collection: After receiving the global model, each client starts local training. Then it computes a function called client utility function (described in details in section VI-B). The cluster head shares this utility value along with some metadata (client IDs, cluster IDs,

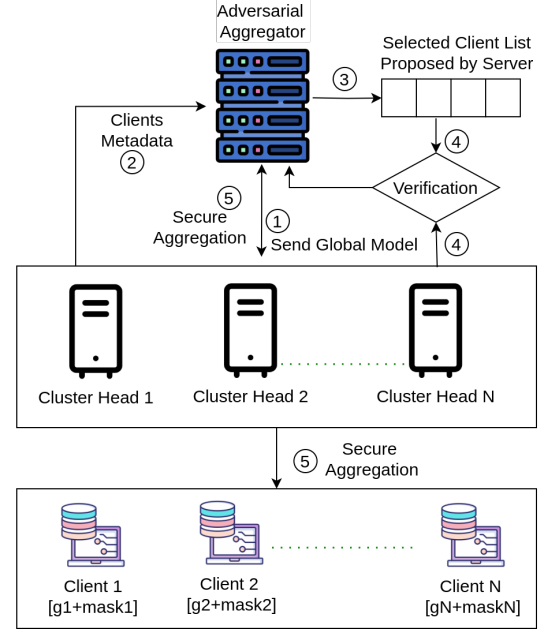


Fig. 4. AdRo-FL for informed selection in cluster-oriented setting.

payload transmission time) with the aggregator. No raw model updates or private data are shared at this stage, maintaining client confidentiality.

Step 3) Client Selection by the aggregator: Using the received utility values, the aggregator performs client selection for the upcoming FL round. Selection can be conducted *cluster-wise* (i.e., locally), or globally (based on the global view of all clients utility values from all clusters). Clients are first filtered based on their ability to meet transmission deadlines (discussed in section VI) and then sorted based on the utility values. The aggregator can limit how many clients can participate at each round and adjust the sorted list accordingly. This sorted list is shared with the cluster heads.

Step 4) Privacy Verification by Cluster Heads: Once the aggregator publishes the list of selected clients, each cluster head reviews the list to ensure compliance with the privacy constraint, i.e., the minimum participation threshold (C clients per cluster, as discussed in section V-A6). If a violation is detected by a cluster head, it withholds its cluster members participation for that round, preventing the aggregator from isolating individual clients to do BSA.

Step 5) Secure Aggregation Participation: Finally, the selected and approved clients participate in the SA protocol. Clients send their masked model updates to the aggregator. This completes one round of FL while maintaining client privacy and ensuring resilience against biased selection attack.

The whole process of client selection for cluster-oriented client setting is formalized in algorithm 1.

Algorithm 1 Secure, Informed Client Selection in Cluster-oriented Client Setting

Input: Global model θ , clusters J , n_j means clients in cluster j where $j \in J$, transmission deadline D , privacy threshold C , client utility function $\mathcal{H}(\omega)$, global sample size S_{global} , local sample size S_{local} , client selection scope E

Output: Participating clients for secure aggregation (SA)

- 1: **Global Model Broadcast:**
- 2: Aggregator sends global model θ to all cluster heads
- 3: **for** each cluster head $j \in J$ **do**
- 4: Distribute θ to all clients $i \in n_j$
- 5: **end for**
- 6: **Client Utility Collection:**
- 7: **for** each client $i \in n_j$ **for** each $j \in J$ **do**
- 8: Compute local training with θ
- 9: Compute transmission time $T_{i,j}$, Client Utility $\mathcal{H}(\omega)$
- 10: Send $M_i = (t_i, \mathcal{H}(\omega))$ to cluster head j
- 11: **end for**
- 12: **for** each cluster head $j \in J$ **do**
- 13: Send $\{M_i\}_{i \in n_j}$ to aggregator
- 14: **end for**
- 15: **Client Selection by Aggregator:**
- 16: Initialize selected clients set $S = \emptyset$
- 17: **for** each cluster $j \in J$ **do**
- 18: Let $F_j = \{i \in n_j \mid T_{i,j} \leq D\}$ ▷ Filter by transmission time
- 19: Sort F_j by $\mathcal{H}(\omega)$ in descending order either locally or globally based on E : let A_j be the sorted list
- 20: Select a subset $S_j \subseteq A_j$ based on ranking (e.g., top- L or all clients meeting a threshold)
- 21: Add S_j to S : $S = S \cup S_j$
- 22: **end for**
- 23: Aggregator publishes selected clients list S
- 24: **Privacy Verification by Cluster Heads:**
- 25: **for** each cluster head $j \in J$ **do**
- 26: **if** $|S \cap n_j| \geq C$ **then**
- 27: Let participating clients from cluster j be $W_j = S \cap n_j$
- 28: **else**
- 29: Let participating clients from cluster j be $W_j = \emptyset$
- 30: ▷ Withhold participation
- 31: **end if**
- 32: **end for**
- 33: Let final participating clients be $W = \bigcup_{j \in J} W_j$
- 34: **SA Participation:**
- 35: **for** each client $i \in W$ **do**
- 36: Client i sends masked model update x_i to aggregator through SA
- 37: **end for**
- 38: Aggregator updates global model θ^{t+1} , then loops to step1

6) *Preventing BSA in Cluster-oriented Setting:* Here we describes how AdRo-FL prevents BSA shown in the threat model.

Although the BSA mentioned in section II did not consider a cluster-oriented client setting, it can very well be applied

on it. To prevent such attack on cluster-oriented client setting, we assume clients are grouped into clusters with trust limited to within each cluster or at least a minority of the clients are trusted in the cluster. We introduce a privacy threshold (C). It imposes that each cluster contributes at least C clients or none. This threshold ensures that a malicious aggregator receives aggregated updates from at least C clients from a cluster. Though setting $C = 2$ ensures that the aggregator cannot obtain a single client's model update, setting $C > 2$ further strengthens security by increasing client resistance to BSA. Without this restriction on clusters, a single client is free to join from a cluster in SA making it vulnerable to BSA. This constraint ensures the privacy of each client by mandating that only clusters with a selected client count greater than or equal to C submit their member client models to the aggregator, while clusters with fewer than C selected clients withhold their participation.

Now we give an example on how AdRo-FL prevents a non-colluding BSA by using similar attack scenario mentioned section II. If $C = 2$, there will be either at least 2 clients from each contributing cluster or none. At round t , the malicious aggregator tries to select a victim with a subset of other clients for SA. Assuming there are four clients from two different clusters (Cluster A and B). With $C = 2$, the aggregator obtains, for instance, $S^t = x_{A1} + x_{A2} + x_{B1} + x_{B2}$ through SA. Here, x_{A1} denotes the model update from client $A1$ and so on. At round $t + 1$, the aggregator attempts to re-select the same subset of clients excluding a victim client (for instance, x_{A1}). Here the malicious aggregator wants to select just one client (x_{A2}) from cluster A . However, since the privacy threshold prevents the cluster A from donating less than $C = 2$ clients, it will withhold from contributing any clients on round $t + 1$. Hence, the aggregator obtains $S^{t+1} = x_{B1} + x_{B2}$ where no clients participate from Cluster A . Therefore, the aggregator fails to uncover x_{A1} by estimating $S^t - S^{t+1}$.

As for colluding BSA, AdRo-FL prevents it in a similar way as discussed above. The aggregator attempts to select a victim client V and a subset of colluding clients. Assume $M3$ and $M4$ are two colluding clients (or Sybils), and $M1$ and $M2$ are two honest clients from a cluster. However, the colluding clients ($M3, M4$) secretly share their model updates with the aggregator. Subsequently, the aggregator, via an SA protocol, obtains $S = x_{M1} + x_{M2} + x_{M3} + x_{M4}$. Since the aggregator knows the local models x_{M3} and x_{M4} of the colluding clients, it attempts to approximate the victim's model as $x_V = S - (x_{M3} + x_{M4})$. This attempt fails because there will be at least C clients aggregated model updates in estimated x_V , not single model update (assuming $C = 2$). Hence the aggregator will obtain $x_V \approx x_{M1} + x_{M2}$.

Furthermore, this scheme protects the privacy of other clusters in case a fake cluster is created by the aggregator. The aggregator may attempt to create colluding clients from the fake cluster but will fail to draw fewer than C clients from other real clusters due to the privacy threshold.

This is to be noted that the defense for cluster-oriented setting is safe against non-colluding attack, as we showed above. As for colluding attacks, there are some nuances. If the aggregator creates a fake cluster or fake clients, AdRo-FL

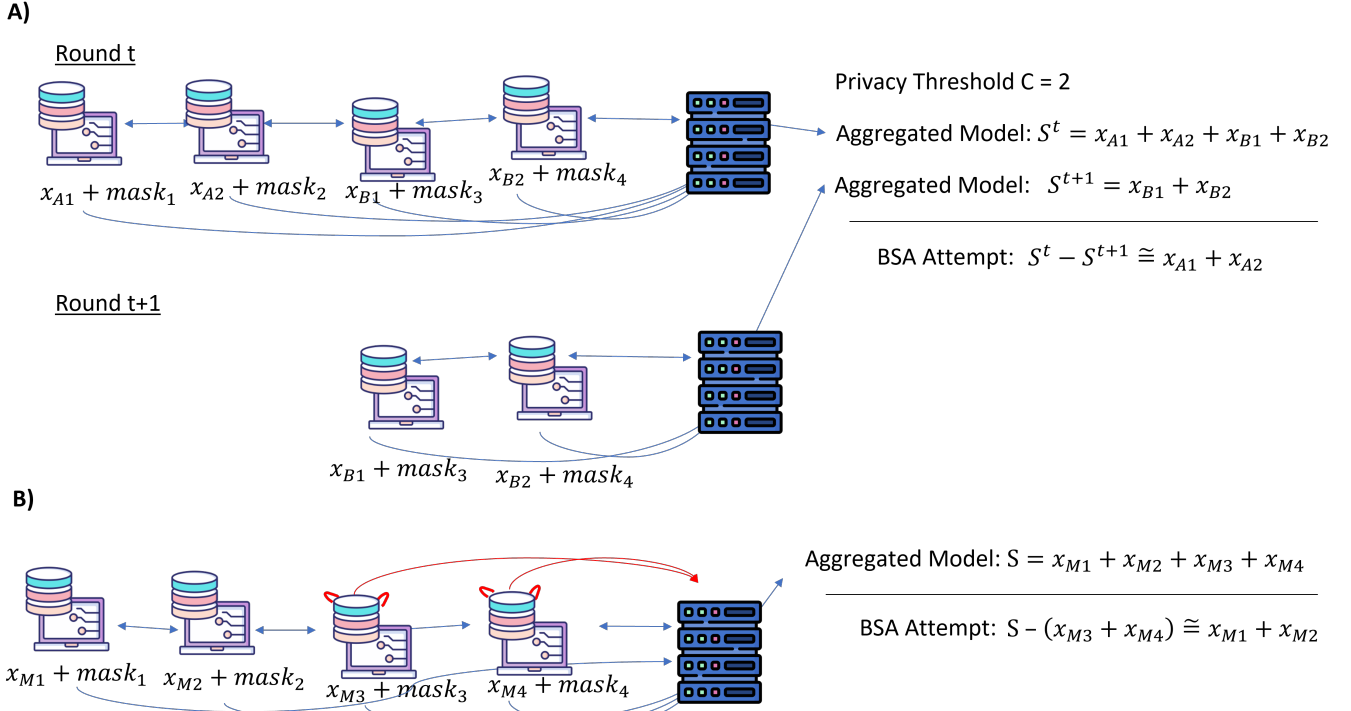


Fig. 5. A) BSA attempt for non-colluding client setting. B) BSA attempt for colluding client setting.

is safe, as the real cluster will still follow the C threshold for client participation which prevents BSA. However, it becomes difficult to prevent BSA on a cluster when there is colluding client(s) within the cluster. For instance, if C is set to 2 and one client from a cluster colludes with the aggregator, then the aggregator will be able to target a client from that cluster. To prevent BSA in such case, we need to set C to at least 3. Similarly, if two clients from a cluster collude with the aggregator, then C should be set to at least 4 and so on. Therefore, in such cases, we suggest to set the value of C to $\geq t + 2$ where t is the number of clients in a cluster that may collude with the aggregator.

However, in many real-world deployments, we do not know the exact number of colluding clients in a cluster. Instead, we assume that each client has an independent probability ϕ of colluding with the aggregator. To protect against BSA, we propose a minimum client selection threshold C such that the number of honest clients among the selected group is at least 2 with high probability. Let $X \sim \text{Binomial}(C, 1 - \phi)$ denote the number of honest clients among the C selected clients. Our goal is to ensure that the probability $P(X < 2) \leq \delta$, where δ is a small risk tolerance parameter (e.g., $\delta = 0.01$). This condition can be satisfied by solving for the smallest C such that $\phi^C + C \cdot \phi^{C-1}(1 - \phi) \leq \delta$. For example, if $\phi = 0.3$ and $\delta = 0.01$, then the required value of C is 12.

Since this calculation may be required frequently in systems with varying levels of assumed collusion or changing security requirements, we recommend precomputing a table of minimum required C values for common combinations of ϕ and δ . This avoids runtime computation and enables fast, adaptive policy enforcement by simply looking up the

appropriate threshold from the table. We provide a sample precomputed lookup table for selecting an appropriate C in Appendix B.

B. Informed client selection for distributed, non-cluster-oriented clients

In this section we describe how informed client selection will be achieved when clients are distributed and there are no trust among clients nor there are any cluster-oriented client setting. Here AdRo-FL does secured informed client selection in two levels. Below we describe the process in details.

1) *First level of selection:* At the first level, after local training, clients are filtered based on a predefined transmission deadline D . The surviving clients compute their local utility value using equation 7. They also produce signatures (like digital signature) for this utility values using their private keys, and submit both the raw utility values and the signatures to the aggregator. The aggregator first ranks these utility values from highest to lowest and saves them into a list. Then it maps the corresponding signatures in another list in the same order and compress it for size reduction. The aggregator broadcasts these two lists for all clients. If a client positions in the top 80% of the sorted utility list, it will participate in the secondary selection process. Any client can verify if the aggregator altered any client's utility value by using the lists and corresponding client's public key. Since every client can explicitly verify the authenticity of the entire selection, all clients can be certain that their selection or non-selection in the first level is completely uninfluenced by the aggregator. In fig 6, we illustrate this as how the primary selection is done.

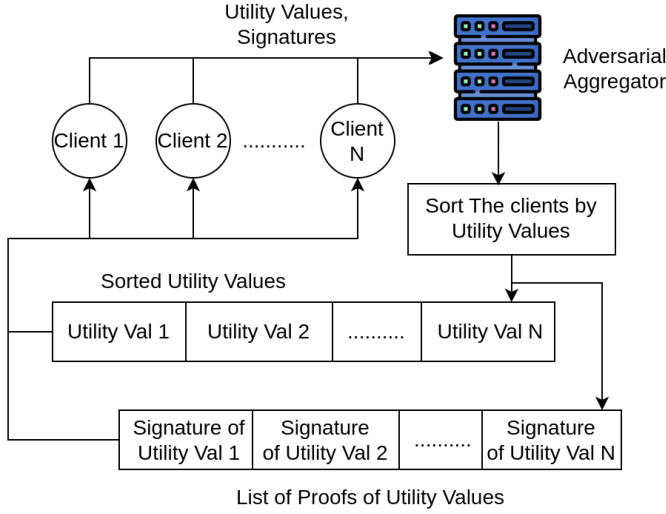


Fig. 6. AdRo-FL: First level of filtering for informed client selection in distributed, non-cluster-oriented setting.

As for how each client generates the signature and how it is utilized for verification, we propose the following method. Each client digitally signs its numeric utility value using its private key (for e.g., using Ed25519 [37]). The aggregator collects the utility values and signatures from all clients, concatenating utility values into a compact binary string and signatures into fixed-sized indexed chunks where each chunk is independently compressed (for e.g., using optimal-level ('zlib') compression [38]). The server then publicly broadcasts these compressed signature chunks along with the concatenated utility values. Upon reception, each client identifies and decompresses only the relevant chunk containing the desired signature, precisely extracts the numeric value and signature based on client indices, and verifies the signature using the corresponding client's public key. This compression and indexing approach significantly reduces payload sizes to enable efficient signature verification suitable for resource-constrained IoT environments. However, more efficient techniques can be used to further reduce the payload size depending on the target IoT devices.

2) *Second level of selection:* To achieve verifiable randomness, we assume that clients use some form of lightweight cryptography to generate private and public keys. In Fig. 7, we illustrate the second-level client filtering process, which consists of three key steps:

Step 1) Generating VRF α : Here, the top 80% values in the sorted utility list from the first level of selection will be concatenated into a string and will be used as the α value for VRF. The reason we don't use a simple deterministic value for VRF α (for e.g., FL round number) is that this will give the aggregator more advantage to create Sybils. For instance, if the aggregator know what will be the α in future rounds, it can carefully generate public key, private key pairs for Sybils that satisfies the VRF threshold-based selection. The aggregator can compute Sybil's hash using $VRF(Sk, \alpha)$, then check if the hash is less than the threshold and repeat until it satisfies.

Step 2) Generation of VRF Hash and VRF Proof: The α

value from the previous step is used as input to the VRF to generate a new value called β . This β is a cryptographic hash, created using secure hash functions like SHA-256 or SHA-512. Along with β , the VRF also produces a proof π . This proof acts as a cryptographic guarantee that β was generated correctly. It is created using either RSA or elliptic curve cryptography (ECC) [39]. Since clients utilize their private keys to independently generate unique hash values and corresponding proofs, both the authenticity and uniqueness of the generated values are ensured.

Step 3) Selection of Winner Clients: Each client compares its VRF hash β against a threshold Z . If the hash value is less than Z , the client is selected for participation in the upcoming FL round. The threshold Z controls how many clients are chosen. Let's assume the hash uses 512-bit. If $Z = 2^{512}$, all clients will be selected; if $Z = 0$, none are selected. This is because the β will be compared against Z . Thus, Z directly influences client participation. Let's assume that in order to eliminate the risk of BSA, we need to select at least K clients out of total clients N . Choosing proper value of K ensures that the chance of a BSA remains negligible. First, we determine the minimum K needed for a desired attack probability P . Then, we calculate the threshold Z based on K and the total number of clients N . A biased selection attack happens when a dishonest aggregator selects mostly colluding clients and isolates a single honest one (discussed in section II). Let's assume $X\%$ of N clients are honest and $Y\%$ are colluding. The aggregator randomly chooses K clients from N (enforced by AdRo-FL using VRF). The probability P of selecting exactly one honest and $K - 1$ colluding clients is given by $P = \frac{\binom{N \cdot X\%}{1} \binom{N \cdot Y\%}{K-1}}{\binom{N}{K}}$. To ensure that this attack probability stays below a safe threshold, we require that $\frac{\binom{N \cdot X\%}{1} \binom{N \cdot Y\%}{K-1}}{\binom{N}{K}} \leq P$.

Solving this inequality gives the value of K needed to satisfy the desired security level. Once K is known, we estimate the threshold Z using $K = N \times \frac{Z}{2^{512}}$, which reflects that the probability of a hash being below Z is proportional to the ratio $Z/2^{512}$. Rearranging gives $Z = \frac{K \times 2^{512}}{N}$. However, due to randomness, fewer than K clients may be selected in some rounds upon applying the threshold-based selection. To address this, we recommend multiplying the term by a conservative factor f whose value can be 2 or 3 depending on the use case. Hence the general formula to find threshold Z becomes: $Z = \frac{f \times K \times 2^{512}}{N}$. Clients whose VRF hash values are lower than Z belong to the winning client set W :

$$\{W\} = \{\text{client}(i) : \text{hash}(\text{client}_i) < Z \text{ for } i = 1 \text{ to } N\}$$

These clients participate in the secure aggregation process. To ensure fairness, any client can verify another client's selection by checking the VRF hash, proof π , and public key. Clients not in $\{W\}$ do not send updates. This two level client selection approach enables informed, unbiased and verifiable random client selection in a distributed setting.

The whole process of client selection for non-cluster-oriented client setting is formalized in algorithm 2. In the following section we discuss the security aspects of this framework.

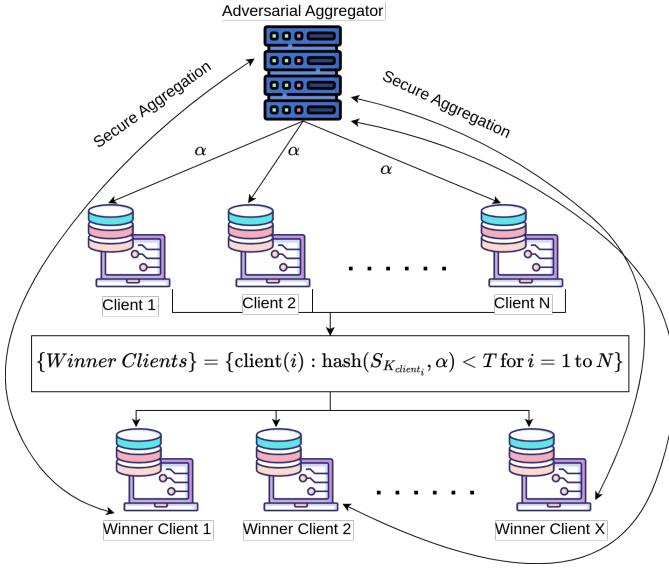


Fig. 7. AdRo-FL: Second level of filtering with verifiable random client selection in distributed, non-cluster-oriented setting.

3) *Security Properties Analysis*: AdRo-FL withstands the following security challenges:

Pool Consistency: Pool consistency ensures all clients observe the same selection outcome, preventing the aggregator from showing different selection results to different clients. In our approach, clients are selected using a VRF hash and a public threshold Z , making the process both random and verifiable. Any client can independently verify the selected set by checking the VRF α , proofs π , hash β , and public keys. This verification ensures transparency and prevents manipulation by the aggregator.

Pool Quality: Pool quality requires that the final selected set includes a minimum number of honest clients. A dishonest aggregator may try to favor colluding clients. To counter this, we ensure the probability P of selecting exactly one honest and $K - 1$ colluding clients remains negligible. This is done by tuning K such that the chance of such biased selection is statistically insignificant. As a result, the final pool $\{W\}$ reliably includes a sufficient number of honest participants.

Anti-Targeting: Anti-targeting prevents the aggregator from selectively including or excluding specific clients. Since the aggregator does not know the VRF α or hash β in advance, it cannot influence the selection outcome. VRF outputs, derived using SHA-512, are uniformly random. Thus, every honest client has an equal chance of being selected, with the probability of selection proportional to $Z/2^{512}$. This guarantees fairness and prevents targeted manipulation.

VI. AdRo-FL: EFFICIENT CLIENT SELECTION

A. Federated Learning Optimization

We adopt Federated Stochastic Gradient Descent (FedSGD) [40] as the foundation for collaborative learning in AdRo-FL. FedSGD is selected over other FL optimization approaches such as FedAvg as it takes less local computational

Algorithm 2 Secure, Informed Client Selection in Non-Cluster-oriented Client Setting

Input: Global model θ , set of all clients N , transmission deadline D , conservative factor f , client utility function $\mathcal{H}(\omega)$, client's secret key sk_i , VRF hash Threshold Z for client selection

Output: Participating clients for secure aggregation (SA)

- 1: **Global Model Broadcast:**
- 2: Aggregator sends global model θ to all clients $i \in N$
- 3: **First-Level Selection:**
- 4: **for** each client $i \in N$ **do**
- 5: Evaluate transmission time T_i
- 6: **if** $T_i \leq D$ **then**
- 7: Compute local client utility score $\mathcal{H}(\omega_i)$
- 8: Generate signature $\sigma_i = \text{Sign}(\mathcal{H}_i)$
- 9: Submit $(\mathcal{H}_i, \sigma_i)$ to aggregator
- 10: **end if**
- 11: **end for**
- 12: Aggregator ranks clients by utility $\mathcal{H}(\omega_i)$, compresses the signature into a list and publishes the ordered lists of utilities and corresponding signatures
- 13: **for** each client i **do**
- 14: **if** i is in top 80% of sorted utility list **then**
- 15: Mark client i eligible for second-level selection
- 16: **end if**
- 17: Verify $(\mathcal{H}_q, \sigma_q)$ for all $q \in N$ using public keys to ensure authenticity of the lists ▷ Optional step
- 18: **end for**
- 19: **Second-Level Selection Using VRF:**
- 20: Concatenate top 80% utility scores into a string to form VRF input α
- 21: **for** each eligible client i **do**
- 22: Generate VRF hash $\beta_i = \text{VRF}(sk_i, \alpha)$
- 23: Generate proof π_i to accompany β_i
- 24: **end for**
- 25: **Winning Client Selection:**
- 26: Initialize $P \leftarrow \emptyset$
- 27: **for** each client i **do**
- 28: **if** $\beta_i < Z$ **then**
- 29: Add client i to P
- 30: **end if**
- 31: **end for**
- 32: **SA Participation:**
- 33: **for** each client $i \in P$ **do**
- 34: Client i sends masked model update x_i to aggregator through SA
- 35: **end for**
- 36: Aggregator updates global model θ^{t+1} , then loops to step 1

time per round, making it particularly suited for scenarios with high computational heterogeneity. Besides, it is typically more memory-efficient as clients process fewer batches per round. FedSGD enables multiple distributed clients to jointly train a shared global model while preserving data privacy. Instead of transmitting raw data, clients compute and transmit model gradients, ensuring that sensitive information remains

TABLE I
SUMMARY OF NOTATION

| Notation | Explanation |
|---------------|--|
| N | Total number of clients |
| J | Total number of clusters |
| t_{max} | Total number of federated learning rounds |
| n_j | Total clients belonging to cluster j |
| $I_{i,j}^k$ | Scheduling decision indicator (1 if client i in cluster j is selected at round t , else 0) |
| $B_{i,j}$ | Channel bandwidth allocated to client i in cluster j |
| $SNR_{i,j}^t$ | Signal-to-noise ratio of client i in cluster j at round t |
| $R_{i,j}^t$ | Achievable data rate (bits/sec) for client i in cluster j at round t |
| PL_q | Quantized payload size (bits) |
| $T_{i,j}^t$ | Transmission time required by client i in cluster j to transmit PL_q |
| D | Deadline, Maximum allowed transmission time per round |
| C | Privacy threshold, Minimum clients needed from each cluster |
| α | Learning rate |

on local devices. In each iteration of FedSGD, the aggregator first distributes the latest global model to all participating clients. Each client then computes the stochastic gradient of its local objective function based on the received global model. These locally computed stochastic gradients are subsequently transmitted to the aggregator, which aggregates them and applies a global SGD update step. This iterative process continues until the model reaches convergence. Each client computes its local gradient at round t : $g_i^t = \nabla f_i(\theta^t)$. The aggregator aggregates these updates and computes the global parameters for round $t+1$: $\theta^{t+1} = \theta^t - \alpha \sum g_i^t$, where θ^t denotes the global model parameters at communication round t , and α represents the learning rate. However, we introduce an additional optimization layer aimed at minimizing overall energy consumption and accelerating convergence at the aggregator. It also considers imposing a minimum client requirement per cluster for cluster-oriented client setting. The optimization formulation for AdRo-FL is defined as follows:

$$\min_{I, \theta} \left[\sum_{j=1}^J \sum_{i \in n_j} (f_{i,j}(\theta) - \sum_{t=1}^{t_{max}} \sum_{j=1}^J \sum_{i \in n_j} I_{i,j}^t \cdot \|\nabla f_{i,j}(\theta^t)\|_2) \right] \quad (1)$$

subject to

$$R_{i,j}^t = B_{i,j} \log_2(1 + SNR_{i,j}^t) \quad (2)$$

$$T_{i,j}^t = I_{i,j}^t \cdot \frac{PL_q}{R_{i,j}^t} \quad (3)$$

$$0 \leq T_{i,j}^t \leq D \quad (4)$$

$$I_{i,j}^t \in \{0, 1\} \quad (5)$$

$$\sum_{i \in n_j} I_{i,j}^t \in \{0, \geq C\}, \forall j \quad (6)$$

The notation used in this formulation is summarized in Table I. The objective function 1 implies that we are minimizing the global loss while ensuring fast convergence and minimum energy consumption. A formal convergence analysis is included in the Appendix C. By selecting the clients based on the loss and gradient (detailed in section VI-B) while they can meet the deadline (constraint (4)) with quantized payload, AdRo-FL speeds up the convergence and avoid transmission failures due to violating the deadline. Constraint (6) is related to the cluster-oriented client setting which enforces that the

number of scheduled clients from each cluster should be $\geq C$ or 0. This constraint C was discussed in details in section V-A6. Equation 1 equally applies to the non-cluster-oriented setting except that the cluster constraints are ignored. For instance, we assume C to be 0 in non-cluster-oriented client setting.

B. Utility-based Client Selection

To accelerate convergence, AdRo-FL considers a client's local loss and weighted gradient norm. In the existing literature, various methods have been used for client selection to improve the global model, notably those based on local loss [21] or gradient norm [20]. Clients with higher gradient norms likely represent data that is harder to fit or provides more informative gradients for global model. On the other hand, by selecting clients with higher local loss, the global model gets updates from clients where it performs worst. However, exclusively relying on gradient norm for measuring a client's contribution towards global model convergence is inefficient as gradient norms are more reflective of how well the global model is currently performing on a client's specific data distribution. If the client's data is mostly from classes that the model is struggling with, its gradient norms may be higher, even if the data is not diverse. It was observed [41], [42] that the gradient norm was significantly higher for clients with imbalanced data than those with comparatively balanced data.

$$\mathcal{H}(\omega) = \omega \mathcal{L}_{\text{local}} + (1 - \omega) \|\nabla f(\theta)\|_2 \cdot \frac{|\mathcal{S}_{\text{local}}|}{|\mathcal{S}_{\text{total}}|} \quad (7)$$

In order to leverage the effectiveness of both gradient norm and local loss, AdRo-FL propose a client utility computation formula in equation 7. The proposed utility function $\mathcal{H}(\omega)$ combines local loss and gradient information to optimize client selection. We empirically determine the value of ω . The first term, $\mathcal{L}_{\text{local}}$, represents the negative log likelihood loss, while $\|\nabla f(\theta)\|_2$ denotes the Euclidean L2 norm of the client's gradients. The ratios $|\mathcal{S}_{\text{local}}|/|\mathcal{S}_{\text{total}}|$ represent the client's proportion of samples. The L2 norm term quantifies the magnitude of model gradient for client's local data, allowing us to prioritize clients with higher gradient norms. This component proves particularly effective for highly imbalanced datasets. Conversely, the $\mathcal{L}_{\text{local}}$ term identifies how well or poorly the model's predictions align with the actual data labels on that specific client and is especially useful for balanced datasets. $\mathcal{H}(\omega)$ achieves a balanced approach, ensuring model generalization across diverse data distributions while maintaining responsiveness to challenging or unique data.

C. Quantization

AdRo-FL uses a fixed quantization for model gradients while ensuring that it does not degrade model performance. The quantization error term E_q is modeled as a function of a fixed quantization level Q : $E_q = \frac{\sigma_q}{Q}$, where σ_q is a scaling factor that controls the quantization error magnitude. The quantization level remains the same throughout all FL rounds, leading to a consistent quantization error.

VII. EXPERIMENTAL RESULTS

A. Dataset

We used MNIST, FMNIST, SVHN, and CIFAR10 datasets to assess AdRo-FL. For all datasets, we used Dirichlet distribution with alpha 0.1 to use high heterogeneity. The privacy threshold C was set to 2 for all the experiments related to cluster-oriented client setting.

B. Baselines

We choose two baseline client selection methods in FL. The first is a conventional random, uninformed client selection approach, commonly used in FL. The second baseline is Oort [12], a state-of-the-art informed client selection technique known for its performance.

C. Environment Setting

1) *Constraints for Energy Efficiency*: For all experiments simulating communication, bandwidth is set at 1 *Mhz*, and power at 0.1 *js*⁻¹. *SNR* is calculated in decibels with a random value in the range of 0 to 30, then converted from dB to linear scale. Energy consumption is computed as (power \times transmission time). These values are chosen arbitrarily and can be modified according to the use case. As for the quantization level, 8 bit quantization was used for all datasets only for AdRo-FL. For the client utility function in VI-B, ω was set to 0.4 based on empirical study. We used deadline $D = 0.5$ seconds for MNIST, $D = 0.11$ seconds for FMNIST, $D = 0.14$ seconds for CIFAR10 and $D = 0.13$ seconds for SVHN dataset. These deadline values are chosen based on the average transmission time found for each dataset with the above setting. All the above values can be tuned to suit the target application.

2) *FL Setting*: We implemented AdRo-FL as well as the baselines in Python using Pytorch. In training, all runs used 3000 rounds. For aggregation, we used FedSGD [43]. For training clients, batch size was set to 64. The learning rate for AdRo-FL and random selection was set to 0.01 while for Oort the learning rate was set to 0.1. These values were chosen after tuning for better performance. The total number of clients were 100 and there were 10 clusters in the cluster-oriented client setting. 20 clients were selected at each FL round. The distribution of clients across clusters was as follows: [12,8,10,11,9,5,15,10,4,16]. Data was distributed across clients using Dirichlet distribution with alpha set to 0.1 for high degree of heterogeneity.

3) *Model Setting*: As for the model, the MNIST dataset uses a feed-forward network with two ReLU-activated hidden layers of 256 and 128 neurons, and a 30% dropout after the first layer to prevent over-fitting. The final layer, using log-softmax, outputs probabilities across 10 classes. For training with FMNIST dataset, we used a multi-layer perceptron with two hidden layers of 64 and 30 units, using ReLU activations and dropout for regularization. For CIFAR10 and SVHN, a CNN was used. It consists two convolutional layers with max pooling, followed by three fully connected layers with dropout regularization and He initialization, culminating in a log softmax output for multi-class classification.

4) *Adversarial Setting*: For cluster-oriented client setting, we set cluster privacy threshold $C = 2$. As for non-cluster-oriented client setting, we set the percentage of honest clients to 90%, the percentage of colluding clients to 10% and set the conservative factor $f = 1$. We set the maximum tolerable probability of biased selection attack to 0.001. All the above values can be tuned to suit the target application.

5) *Hardware*: The experiments utilized a AMD Ryzen Threadripper PRO 5995WX 64-core processor with 500 GB of RAM and an NVIDIA RTX A6000 GPU with 48 GB of VRAM.

D. Results

In this section, we present the performance of AdRo-FL on benchmark datasets. We also compare it with baselines. As mentioned in section IV, we consider two selection approaches for two different client settings. However, for cluster-oriented environment, we further divide it into local, cluster-wise selection and global selection. In the following subsections, we present the experimental outcomes. We report time-to-accuracy as the number of communication rounds required to reach a specified test accuracy threshold.

1) *Cluster-oriented Client Setting*: In this subsection, we report the performance of AdRo-FL. Besides, we present a comparison with the baselines, i.e., random selection and Oort. For each of the methods, we show the results of two approaches of selecting clients in cluster-oriented setting. First is local selection and the second is global selection. In local selection, we choose a fixed number of clients from each cluster while in global selection, clients are selected globally based on the view of all clients utility values. The random selection baseline method selects clients randomly while AdRo-FL and Oort selects client based on respective utility function. We present the results in Fig. 8, 9, 8, 9.

We start by local, cluster-wise selection. Fig. 8 shows the accuracy of AdRo-FL and the baselines against the round numbers. The figures shows that AdRo-FL either achieves better or competitive performance compared to baselines. Fig. 9 shows the loss values of AdRo-FL and the baselines against FL round numbers.

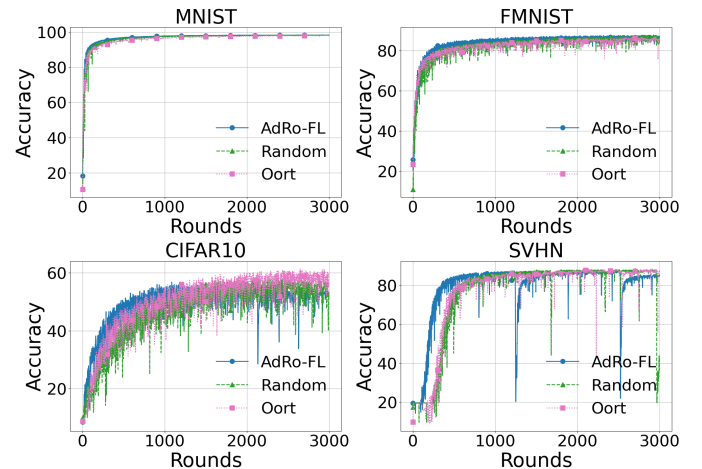


Fig. 8. Accuracy comparison of AdRo-FL with random selection and Oort on four datasets. Here we used local selection (cluster-oriented).

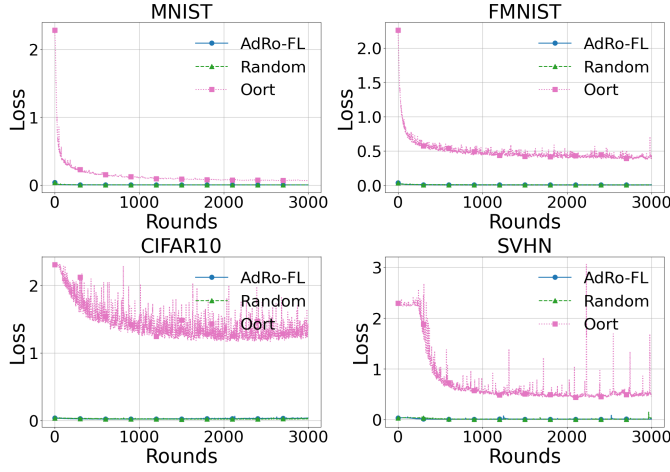


Fig. 9. Loss comparison of AdRo-FL with random selection and Oort on four datasets. Here we used local selection (cluster-oriented).

TABLE II

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE MNIST DATASET (CLUSTER-ORIENTED SETTING). THE BEST VALUE FOR EACH ROW IS HIGHLIGHTED IN BOLD. HERE CLIENTS WERE SELECTED LOCALLY FROM ALL CLUSTERS.

| Method | @ 60 | @ 65 | @ 70 | @ 75 | @ 80 | @ 85 | @ 90 | Best Accuracy (%) |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-------------------|
| AdRo-FL | 12 | 20 | 20 | 21 | 27 | 40 | 63 | 98.23 |
| Random | 16 | 18 | 23 | 23 | 37 | 54 | 86 | 98.20 |
| Oort | 11 | 15 | 15 | 30 | 31 | 48 | 108 | 98.04 |

In table II, we show the number of rounds needed to reach different target accuracies for MNIST dataset. Also we note the final accuracy achieved by each method. On average, compared to random baseline, AdRo-FL achieves $1.22\times$ time-to-accuracy improvement. Compared to Oort, it achieves $1.13\times$ time-to-accuracy improvement on average.

TABLE III

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE FMNIST DATASET (CLUSTER-ORIENTED SETTING). THE BEST VALUE FOR EACH ROW IS HIGHLIGHTED IN BOLD. HERE CLIENTS WERE SELECTED LOCALLY FROM ALL CLUSTERS.

| Method | @ 60 | @ 65 | @ 70 | @ 75 | @ 80 | @ 85 | Best Accuracy (%) |
|---------|-----------|-----------|-----------|------------|------------|------------|-------------------|
| AdRo-FL | 43 | 52 | 56 | 116 | 195 | 639 | 87.75 |
| Random | 50 | 57 | 89 | 156 | 311 | 1127 | 87.38 |
| Oort | 39 | 43 | 69 | 114 | 312 | 1352 | 86.95 |

In table III, we show the number of rounds needed to reach different target accuracies for FMNIST dataset. Also we note the final accuracy achieved by each method (proposed and Oort). On average, compared to random baseline, AdRo-FL achieves $1.43\times$ time-to-accuracy improvement. Compared to Oort, it achieves $1.28\times$ time-to-accuracy improvement on average and $1.01\times$ final accuracy improvement.

TABLE IV

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE CIFAR10 DATASET (CLUSTER-ORIENTED SETTING). THE BEST VALUE FOR EACH ROW IS HIGHLIGHTED IN BOLD. HERE CLIENTS WERE SELECTED LOCALLY FROM ALL CLUSTERS.

| Method | @ 35 | @ 40 | @ 45 | @ 50 | @ 55 | Best Accuracy (%) |
|---------|------------|------------|------------|------------|------------|-------------------|
| AdRo-FL | 155 | 209 | 286 | 469 | 906 | 56.79 |
| Random | 330 | 416 | 548 | 806 | 1334 | 58.15 |
| Oort | 208 | 289 | 393 | 660 | 997 | 61.40 |

In table IV, we show the number of rounds needed to reach different target accuracies for CIFAR10 dataset. Also we note the final accuracy achieved by each method (proposed and Oort). On average, compared to random baseline, AdRo-FL achieves $1.85\times$ time-to-accuracy improvement. Compared to Oort, it achieves $1.32\times$ time-to-accuracy improvement on average.

TABLE V

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE SVHN DATASET (CLUSTER-ORIENTED SETTING). THE BEST VALUE FOR EACH ROW IS HIGHLIGHTED IN BOLD. HERE CLIENTS WERE SELECTED LOCALLY FROM ALL CLUSTERS.

| Method | @ 60 | @ 65 | @ 70 | @ 75 | @ 80 | @ 85 | Best Accuracy (%) |
|---------|------------|------------|------------|------------|------------|------------|-------------------|
| AdRo-FL | 202 | 228 | 238 | 281 | 345 | 561 | 87.33 |
| Random | 388 | 425 | 435 | 477 | 572 | 878 | 88.09 |
| Oort | 354 | 379 | 404 | 438 | 535 | 966 | 88.44 |

In table V, we show the number of rounds needed to reach different target accuracies for SVHN dataset. Also we note the final accuracy achieved by each method (proposed and Oort). On average, compared to random baseline, AdRo-FL achieves $1.76\times$ time-to-accuracy improvement. Compared to Oort, it achieves $1.66\times$ time-to-accuracy improvement on average.

Now we present the results of global selection in cluster-oriented client setting.

While global selection in AdRo-FL, although the aggregator ranks clients globally, each cluster-head enforces the privacy constraint locally by checking if at least C clients from its cluster were selected. If not, it withholds participation, thereby ensuring privacy. Fig. 10 shows the accuracy of AdRo-FL and the baselines against the FL round numbers. The figures shows that AdRo-FL either achieves better or competitive performance compared to baselines. Fig. 11 shows the loss values of AdRo-FL and the baselines against FL round numbers.

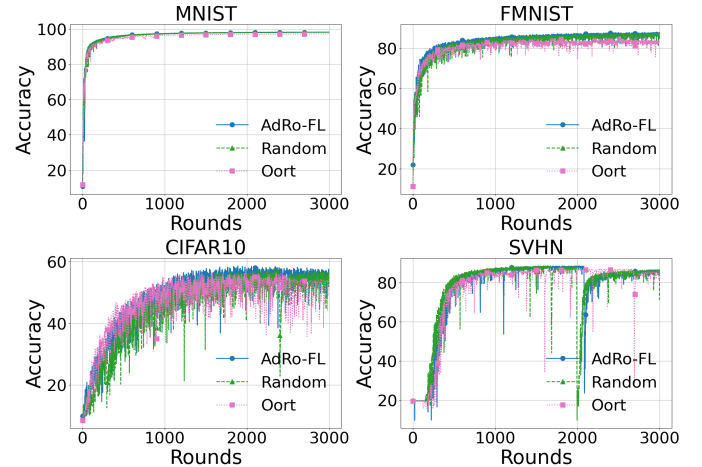


Fig. 10. Accuracy comparison of AdRo-FL with random selection and Oort on four datasets. Here we used global selection (cluster-oriented).

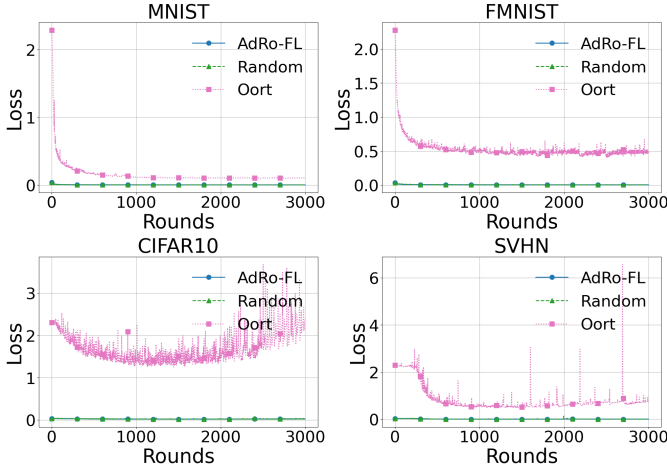


Fig. 11. Loss comparison of AdRo-FL with random selection and Oort on four datasets. Here we used global selection (cluster-oriented).

TABLE VI

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE MNIST DATASET (CLUSTER-ORIENTED SETTING). THE BEST VALUE FOR EACH ROW IS HIGHLIGHTED IN BOLD. HERE CLIENTS WERE SELECTED GLOBALLY FROM ALL CLUSTERS.

| Method | @ 60 | @ 65 | @ 70 | @ 75 | @ 80 | @ 85 | @ 90 | Best Accuracy (%) |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-------------------|
| AdRo-FL | 15 | 15 | 15 | 23 | 37 | 45 | 74 | 98.19 |
| Random | 21 | 21 | 21 | 25 | 36 | 45 | 77 | 98.23 |
| Oort | 18 | 19 | 28 | 28 | 38 | 47 | 108 | 97.25 |

In table VI, we show the number of rounds needed to reach different target accuracies for MNIST dataset. Also we note the final accuracy achieved by each method (proposed and Oort). On average, compared to random baseline, AdRo-FL achieves $1.19\times$ time-to-accuracy improvement. Compared to Oort, it achieves $1.30\times$ time-to-accuracy improvement on average and $1.01\times$ final accuracy improvement.

TABLE VII

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE FMNIST DATASET (CLUSTER-ORIENTED SETTING). THE BEST VALUE FOR EACH ROW IS HIGHLIGHTED IN BOLD. HERE CLIENTS WERE SELECTED GLOBALLY FROM ALL CLUSTERS.

| Method | @ 60 | @ 65 | @ 70 | @ 75 | @ 80 | @ 85 | Best Accuracy (%) |
|---------|-----------|-----------|-----------|------------|------------|------------|-------------------|
| AdRo-FL | 40 | 40 | 67 | 118 | 243 | 795 | 88.06 |
| Random | 50 | 87 | 122 | 167 | 312 | 980 | 87.47 |
| Oort | 37 | 58 | 78 | 157 | 348 | 2195 | 85.82 |

In table VII, we show the number of rounds needed to reach different target accuracies for FMNIST dataset. Also we note the final accuracy achieved by each method (proposed and Oort). On average, compared to random baseline, AdRo-FL achieves $1.53\times$ time-to-accuracy improvement and $1.01\times$ final accuracy improvement. Compared to Oort, it achieves $1.51\times$ time-to-accuracy improvement on average and $1.03\times$ final accuracy improvement.

TABLE VIII

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE CIFAR10 DATASET (CLUSTER-ORIENTED SETTING). THE BEST VALUE FOR EACH ROW IS HIGHLIGHTED IN BOLD. HERE CLIENTS WERE SELECTED GLOBALLY FROM ALL CLUSTERS.

| Method | @ 35 | @ 40 | @ 45 | @ 50 | @ 55 | Best Accuracy (%) |
|---------|------------|------------|------------|------------|------------|-------------------|
| AdRo-FL | 283 | 347 | 496 | 728 | 1108 | 58.68 |
| Random | 408 | 497 | 666 | 974 | 1557 | 57.90 |
| Oort | 210 | 313 | 397 | 655 | 988 | 56.69 |

In table VIII, we show the number of rounds needed to reach different target accuracies for CIFAR10 dataset. Also we note the final accuracy achieved by each method (proposed and Oort). On average, compared to random baseline, AdRo-FL achieves $1.39\times$ time-to-accuracy improvement and $1.01\times$ final accuracy improvement. Compared to Oort, it achieves $1.04\times$ final accuracy improvement.

TABLE IX

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE SVHN DATASET (CLUSTER-ORIENTED SETTING). THE BEST VALUE FOR EACH ROW IS HIGHLIGHTED IN BOLD. HERE CLIENTS WERE SELECTED GLOBALLY FROM ALL CLUSTERS.

| Method | @ 60 | @ 65 | @ 70 | @ 75 | @ 80 | @ 85 | Best Accuracy (%) |
|---------|------------|------------|------------|------------|------------|------------|-------------------|
| AdRo-FL | 344 | 358 | 387 | 418 | 499 | 698 | 88.46 |
| Random | 285 | 300 | 335 | 363 | 441 | 627 | 88.51 |
| Oort | 335 | 351 | 376 | 414 | 512 | 873 | 86.84 |

In table IX, we show the number of rounds needed to reach different target accuracies for SVHN dataset. Also we note the final accuracy achieved by each method (proposed and Oort). Compared to Oort, it achieves $1.03\times$ time-to-accuracy improvement on average and $1.02\times$ final accuracy improvement.

Lastly, in Fig. 12, we show how many clusters became vulnerable to BSA per round because of violating cluster constraints when doing client selection globally in the cluster-oriented setting. We discussed in section II and V-A6 that if the adversarial aggregator can manage to single out a client from a cluster, it can do BSA to know clients private model update. Insecure client selection technique Oort does not consider this important constraint while selecting best clients based on utility. The figure illustrates the privacy risk of Oort if applied on cluster-oriented client setting to select client globally. Compared to Oort, AdRo-FL handles this vulnerability while not compromising on model performance.

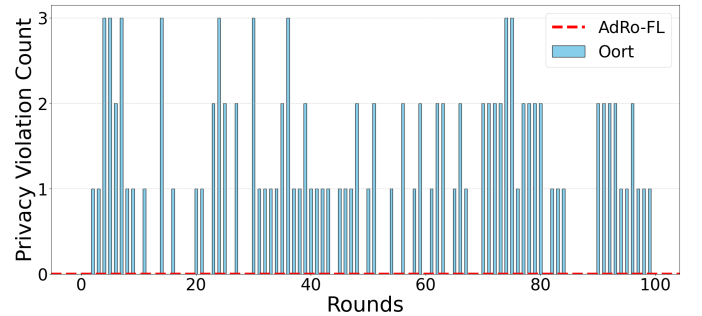


Fig. 12. Number of clusters became vulnerable to BSA per round.

2) *Distributed, Non-cluster-oriented Client Setting*: In this section, we present the results of experiments conducted with

non-cluster-oriented clients. As detailed in section V-B, AdRo-FL uses two level client filtering for informed client selection in distributed, non-cluster-oriented client setting. We compared the result with baselines. In Fig. 13, we show result of 3000 rounds of training on the selected datasets. The plot shows that AdRo-FL achieves better time-to-accuracy as well as target accuracy than Oort. In Fig. 14, we shows the loss reduction across FL rounds.

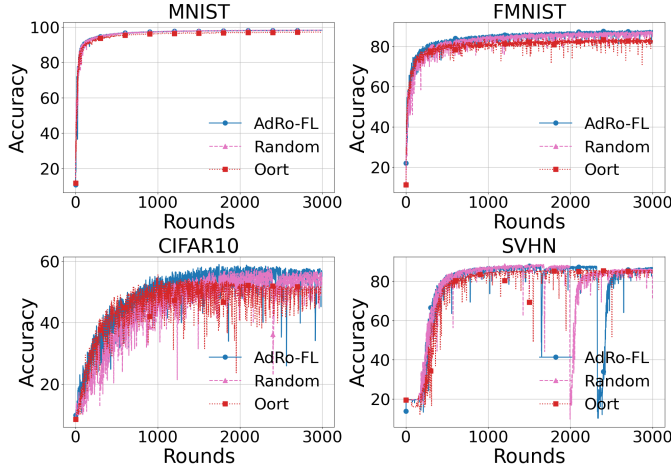


Fig. 13. Accuracy comparison of AdRo-FL (VRF-based) with random selection and Oort on four datasets in non-cluster-oriented setting.

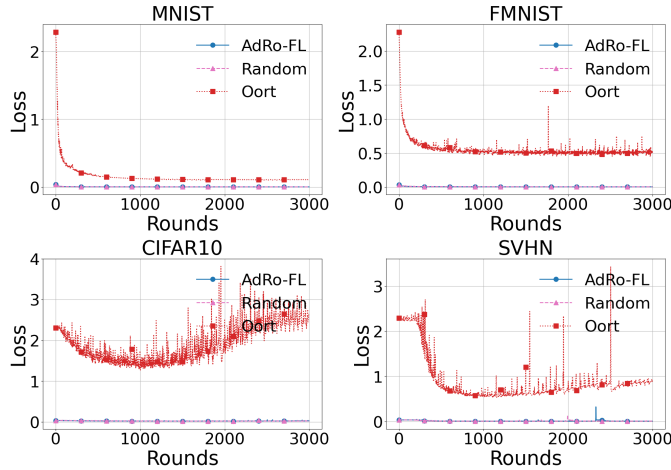


Fig. 14. Loss comparison of AdRo-FL (VRF-based) with random selection and Oort on four datasets in non-cluster-oriented setting.

TABLE X

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE MNIST DATASET (NON-CLUSTER-ORIENTED SETTING). BEST VALUES FOR EACH ROW ARE HIGHLIGHTED IN BOLD.

| Method | @ 60 | @ 65 | @ 70 | @ 75 | @ 80 | @ 85 | @ 90 | Best Accuracy (%) |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-------------------|
| AdRo-FL | 15 | 15 | 15 | 23 | 37 | 45 | 74 | 98.19 |
| Random | 21 | 21 | 21 | 25 | 36 | 45 | 77 | 98.23 |
| Oort | 16 | 21 | 22 | 25 | 36 | 52 | 112 | 97.17 |

In table X, we show the number of rounds needed to reach different target accuracies for MNIST dataset. Also we note the final accuracy achieved by each method. On

average, compared to Oort, AdRo-FL achieves $1.24\times$ time-to-accuracy improvement and $1.01\times$ final accuracy improvement. Compared to random selection, AdRo-FL achieves $1.19\times$ time-to-accuracy improvement on average.

TABLE XI

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE FMNIST DATASET (NON-CLUSTER-ORIENTED SETTING). BEST VALUES FOR EACH ROW ARE HIGHLIGHTED IN BOLD.

| Method | @ 60 | @ 65 | @ 70 | @ 75 | @ 80 | @ 85 | Best Accuracy (%) |
|---------|-----------|-----------|-----------|------------|------------|------------|-------------------|
| AdRo-FL | 40 | 40 | 67 | 118 | 243 | 795 | 88.06 |
| Random | 50 | 87 | 122 | 167 | 312 | 980 | 87.47 |
| Oort | 36 | 55 | 93 | 128 | 345 | N/A | 84.59 |

In table XI, we show the number of rounds needed to reach different target accuracies for FMNIST dataset. Also we note the final accuracy achieved by each method. On average, compared to Oort, AdRo-FL achieves $1.23\times$ time-to-accuracy improvement and $1.04\times$ final accuracy improvement. Compared to random selection, AdRo-FL achieves $1.53\times$ time-to-accuracy improvement and $1.01\times$ final accuracy improvement, on average.

TABLE XII

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE CIFAR-10 DATASET (NON-CLUSTER-ORIENTED SETTING). BEST VALUES FOR EACH ROW ARE HIGHLIGHTED IN BOLD.

| Method | @ 35 | @ 40 | @ 45 | @ 50 | @ 55 | Best Accuracy (%) |
|---------|------------|------------|------------|------------|------------|-------------------|
| AdRo-FL | 219 | 336 | 415 | 739 | 966 | 58.92 |
| Random | 408 | 497 | 666 | 974 | 1557 | 57.90 |
| Oort | 225 | 286 | 428 | 674 | 1003 | 55.53 |

In table XII, we show the number of rounds needed to reach different target accuracies for CIFAR10 dataset. Also we note the final accuracy achieved by each method. On average, compared to Oort, AdRo-FL achieves $1.06\times$ final accuracy improvement. However, compared to random selection, on average, AdRo-FL achieves $1.58\times$ time-to-accuracy and $1.02\times$ final accuracy improvement.

TABLE XIII

TIME-TO-ACCURACY AND FINAL TEST ACCURACY COMPARISON ON THE SVHN DATASET (NON-CLUSTER-ORIENTED SETTING). BEST VALUES FOR EACH ROW ARE HIGHLIGHTED IN BOLD.

| Method | @ 60 | @ 65 | @ 70 | @ 75 | @ 80 | @ 85 | Best Accuracy (%) |
|---------|------------|------------|------------|------------|------------|------------|-------------------|
| AdRo-FL | 272 | 301 | 326 | 361 | 422 | 643 | 87.99 |
| Random | 285 | 300 | 335 | 363 | 441 | 627 | 88.51 |
| Oort | 347 | 355 | 397 | 439 | 552 | 1158 | 85.53 |

In table XIII, we show the number of rounds needed to reach different target accuracies for SVHN dataset. Also we note the final accuracy achieved by each method. On average, compared to Oort, AdRo-FL achieves $1.33\times$ time-to-accuracy improvement and $1.03\times$ final accuracy improvement. Compared to random selection, time-to-accuracy was improved $1.02\times$ on average.

3) *Resource Consumption Comparison:* In Fig. 15, we showed the energy consumption estimations. The figure shows energy consumption estimated in cluster-oriented client setting. However, similar results were observed for non-cluster-oriented client setting. AdRo-FL achieved superior energy efficiency due to deadline based client selection and quantization.

This energy efficiency is achieved without noticeable degradation of performance which indicates AdRo-FL's robustness in delivering better performance with reduced memory consumption.

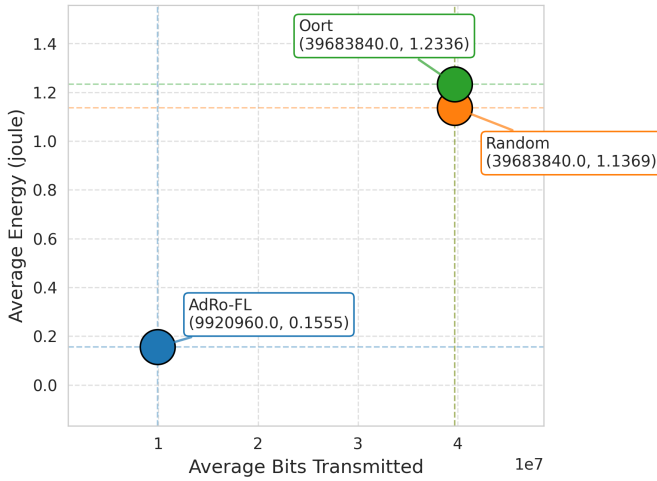


Fig. 15. Average energy consumption (Joules) and bits transmitted per round for SVHN dataset in cluster-oriented setting.

VIII. CONCLUSION

Ensuring faster convergence as well as client privacy remains one of the most pressing challenges in FL, particularly in adversarial settings where the central aggregator may not follow the protocol. Ensuring the aggregator cannot alter the client selection process to do biased selection attack, is still a critical challenge that requires robust mechanisms to guarantee fairness and verifiability in FL systems. In this work we presented AdRo-FL, a practical solution to an underexplored problem of secure and informed client selection in the presence of adversarial aggregator in FL. Unlike existing methods, which either compromise performance through random uninformed client selection or risk privacy by relying on utility-based informed selection, we uniquely cover both that realizes informed as well as secure client selection. AdRo-FL prevents adversarial aggregator from singling out a victim client while ensuring that clients are selected based on their practical utility. Experimental outcomes with several benchmark datasets shows that AdRo-FL securely achieves significant performance improvement over insecure baselines.

REFERENCES

- [1] C. Zhang, Y. Xie, H. Bai, B. Yu, W. Li, and Y. Gao, "A survey on federated learning," *Knowledge-Based Systems*, vol. 216, p. 106775, 2021.
- [2] J. Le, D. Zhang, X. Lei, L. Jiao, K. Zeng, and X. Liao, "Privacy-preserving federated learning with malicious clients and honest-but-curious servers," *IEEE Transactions on Information Forensics and Security*, 2023.
- [3] J. Cai, W. Shen, and J. Qin, "Esvfl: Efficient and secure verifiable federated learning with privacy-preserving," *Information Fusion*, vol. 109, p. 102420, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1566253524001982>
- [4] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawit, Z. Charles, G. Cormode, R. Cummings, R. G. L. D'Oliveira, H. Eichner, S. El Rouayheb, D. Evans, J. Gardner, Z. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, Z. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, S. Koyejo, T. Lepoint, Y. Liu, P. Mittal, M. Mohri, R. Nock, A. Özgür, R. Pagh, H. Qi, D. Ramage, R. Raskar, M. Raykova, D. Song, W. Song, S. U. Stich, Z. Sun, A. Theertha Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao, 2021.
- [5] R. Aziz, S. Banerjee, S. Bouzeffrane, and T. Le Vinh, "Exploring homomorphic encryption and differential privacy techniques towards secure federated learning paradigm," *Future Internet*, vol. 15, no. 9, 2023. [Online]. Available: <https://www.mdpi.com/1999-5903/15/9/310>
- [6] N. Rodríguez-Barroso, G. Stipicich, D. Jiménez-López, J. A. Ruiz-Millán, E. Martínez-Cámara, G. González-Seco, M. V. Luzón, M. A. Veganzones, and F. Herrera, "Federated learning and differential privacy: Software tools analysis, the sherpa.ai fl framework and methodological guidelines for preserving data privacy," *Information Fusion*, vol. 64, pp. 270–292, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1566253520303213>
- [7] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '17. New York, NY, USA: Association for Computing Machinery, 2017, p. 1175–1191. [Online]. Available: <https://doi.org/10.1145/3133956.3133982>
- [8] T. Nguyen, P. Thai, J. Tre'R, T. N. Dinh, and M. T. Thai, "Blockchain-based secure client selection in federated learning," in *2022 IEEE International Conference on Blockchain and Cryptocurrency (ICBC)*. IEEE, 2022, pp. 1–9.
- [9] Z. Jiang, P. Ye, S. He, W. Wang, R. Chen, and B. Li, "Lotto: Secure participant selection against adversarial servers in federated learning," in *33rd USENIX Security Symposium (USENIX Security 24)*. Philadelphia, PA: USENIX Association, Aug. 2024, pp. 343–360. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity24/presentation/jiang-zhifeng>
- [10] J. Li, T. Chen, and S. Teng, "A comprehensive survey on client selection strategies in federated learning," *Computer Networks*, p. 110663, 2024.
- [11] G. D. M. Jimenez, A. Anagnostopoulos, I. Chatzigiannakis, and A. Vitaletti, "Fedartml: A tool to facilitate the generation of non-iid datasets in a controlled way to support federated learning research," *IEEE Access*, 2024.
- [12] F. Lai, X. Zhu, H. V. Madhyastha, and M. Chowdhury, "Oort: Efficient federated learning via guided participant selection," in *15th USENIX Symposium on Operating Systems Design and Implementation (OSDI 21)*. USENIX Association, Jul. 2021, pp. 19–35. [Online]. Available: <https://www.usenix.org/conference/osdi21/presentation/lai>
- [13] Z. Li, Y. Dang, and X. Chen, "Node selection for model quality optimization in hierarchical federated learning based on deep reinforcement learning," *Peer-to-Peer Networking and Applications*, vol. 17, no. 3, pp. 1720–1731, 2024.
- [14] M. Tang, X. Ning, Y. Wang, Y. Wang, and Y. Chen, "Fedgpp: Correlation-based active client selection strategy for heterogeneous federated learning," *arXiv preprint arXiv:2103.13822*, 2021.
- [15] H. Wu and P. Wang, "Node selection toward faster convergence for federated learning on non-iid data," *IEEE Transactions on Network Science and Engineering*, vol. 9, no. 5, pp. 3099–3111, 2022.
- [16] S. Trindade and N. L. S. da Fonseca, "Client selection in hierarchical federated learning," *IEEE Internet of Things Journal*, vol. 11, no. 17, pp. 28 480–28 495, 2024.
- [17] F. Shi, C. Hu, W. Lin, L. Fan, T. Huang, and W. Wu, "Vfedcs: Optimizing client selection for volatile federated learning," *IEEE Internet of Things Journal*, vol. 9, no. 24, pp. 24 995–25 010, 2022.
- [18] X. Chen, X. Zhou, H. Zhang, M. Sun, and H. Vincent Poor, "Client selection for wireless federated learning with data and latency heterogeneity," *IEEE Internet of Things Journal*, vol. 11, no. 19, pp. 32 183–32 196, 2024.
- [19] O. Wehbi, S. Arisdakessian, O. A. Wahab, H. Otrouk, S. Otoum, A. Mourad, and M. Guizani, "Fedmint: Intelligent bilateral client selection in federated learning with newcomer iot devices," *IEEE Internet of Things Journal*, vol. 10, no. 23, pp. 20 884–20 898, 2023.
- [20] O. Marnissi, H. E. Hammouti, and E. H. Bergou, "Client selection in federated learning based on gradients importance," in *AIP Conference Proceedings*, vol. 3034, no. 1. AIP Publishing, 2024.

- [21] Y. J. Cho, J. Wang, and G. Joshi, “Towards understanding biased client selection in federated learning,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 10351–10375.
- [22] F. Sabah, Y. Chen, Z. Yang, A. Raheem, M. Azam, N. Ahmad, and R. Sarwar, “Fairdpfl-scs: Fair dynamic personalized federated learning with strategic client selection for improved accuracy and fairness,” *Information Fusion*, vol. 115, p. 102756, 2025. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1566253524005347>
- [23] J. Guo, L. Su, J. Liu, J. Ding, X. Liu, B. Huang, and L. Li, “Auction-based client selection for online federated learning,” *Information Fusion*, vol. 112, p. 102549, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1566253524003270>
- [24] Z. Niu, H. Dong, A. K. Qin, and T. Gu, “Flrce: Resource-efficient federated learning with early-stopping strategy,” *IEEE Transactions on Mobile Computing*, pp. 1–16, 2024.
- [25] T. Dettmers, “8-bit approximations for parallelism in deep learning,” in *4th International Conference on Learning Representations, ICLR, San Juan, Puerto Rico*, Y. Bengio and Y. LeCun, Eds., 2016.
- [26] C. Li, X. Zeng, M. Zhang, and Z. Cao, “Pyramidfl: A fine-grained client selection framework for efficient federated learning,” in *Proceedings of the 28th Annual International Conference on Mobile Computing And Networking*, 2022, pp. 158–171.
- [27] C. Briggs, Z. Fan, and P. Andras, “Federated learning with hierarchical clustering of local updates to improve training on non-iid data,” in *2020 international joint conference on neural networks (IJCNN)*. IEEE, 2020, pp. 1–9.
- [28] M. Asad, A. Moustafa, F. A. Rabhi, and M. Aslam, “Thf: 3-way hierarchical framework for efficient client selection and resource management in federated learning,” *IEEE Internet of Things Journal*, vol. 9, no. 13, pp. 11085–11097, 2022.
- [29] J. Liu, F. Lai, Y. Dai, A. Akella, H. V. Madhyastha, and M. Chowdhury, “Auxo: Efficient federated learning via scalable client clustering,” in *Proceedings of the 2023 ACM Symposium on Cloud Computing*, ser. SoCC ’23. New York, NY, USA: Association for Computing Machinery, 2023, p. 125–141. [Online]. Available: <https://doi.org/10.1145/3620678.3624651>
- [30] C. Keçeci, M. Shaqfeh, F. Al-Qahtani, M. Ismail, and E. Serpedin, “Clustered scheduling and communication pipelining for efficient resource management of wireless federated learning,” *IEEE Internet of Things Journal*, vol. 10, no. 15, pp. 13303–13316, 2023.
- [31] J. So, C. He, C.-S. Yang, S. Li, Q. Yu, R. E. Ali, B. Guler, and S. Avestimehr, “Lightsecagg: a lightweight and versatile design for secure aggregation in federated learning,” in *Proceedings of Machine Learning and Systems*, D. Marculescu, Y. Chi, and C. Wu, Eds., vol. 4, 2022, pp. 694–720. [Online]. Available: https://proceedings.mlsys.org/paper_files/paper/2022/file/6c44dc73014d66ba49b28d483a8f8b0d-Paper.pdf
- [32] U. Erlingsson, V. Feldman, I. Mironov, A. Raghunathan, K. Talwar, and A. Thakurta, “Amplification by shuffling: from local to central differential privacy via anonymity,” in *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, ser. SODA ’19. USA: Society for Industrial and Applied Mathematics, 2019, p. 2468–2479.
- [33] Z. Jiang, W. Wang, and R. Chen, “Dordis: Efficient federated learning with dropout-resilient differential privacy,” ser. EuroSys ’24. New York, NY, USA: Association for Computing Machinery, 2024, p. 472–488. [Online]. Available: <https://doi.org/10.1145/3627703.3629559>
- [34] C. Fung, C. J. Yoon, and I. Beschastnikh, “The limitations of federated learning in sybil settings,” in *23rd International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2020)*, 2020, pp. 301–316.
- [35] H. R. Roth, K. Chang, P. Singh, N. Neumark, W. Li, V. Gupta, S. Gupta, L. Qu, A. Ihsani, B. C. Bizzo, Y. Wen, V. Buch, M. Shah, F. Kitamura, M. Mendonça, V. Lavor, A. Harouni, C. Compas, J. Tetreault, P. Dogra, Y. Cheng, S. Erdal, R. White, B. Hashemian, T. Schultz, M. Zhang, A. McCarthy, B. M. Yun, E. Sharaf, K. V. Hoebe, J. B. Patel, B. Chen, S. Ko, E. Leibovitz, E. D. Pisano, L. Coombs, D. Xu, K. J. Dreyer, I. Dayan, R. C. Naidu, M. Flores, D. Rubin, and J. Kalpathy-Cramer, “Federated learning for breast density classification: A real-world implementation.” Berlin, Heidelberg: Springer-Verlag, 2020, p. 181–191. [Online]. Available: https://doi.org/10.1007/978-3-030-60548-3_18
- [36] A. Said, A. Yahyaoui, and T. Abdellatif, “Hipa and gdpr compliance in iot healthcare systems,” in *International conference on model and data engineering*. Springer, 2023, pp. 198–209.
- [37] J. Brendel, C. Cremers, D. Jackson, and M. Zhao, “The provable security of ed25519: theory and practice,” in *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2021, pp. 1659–1676.
- [38] A. P. Maulidina, R. A. Wijaya, K. Mazel, and M. S. Astriani, “Comparative study of data compression algorithms: Zstandard, zlib & lz4,” in *International Conference on Science, Engineering Management and Information Technology*. Springer, 2023, pp. 394–406.
- [39] D. Papadopoulos, D. Wessels, S. Huque, M. Naor, J. Včelák, L. Reyzin, and S. Goldberg, “Making nsec5 practical for dnssec,” *Cryptology ePrint Archive*, 2017.
- [40] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [41] Z. Xiao, Z. Chen, S. Liu, H. Wang, Y. Feng, J. Hao, J. T. Zhou, J. Wu, H. H. Yang, and Z. Liu, “Fed-grab: federated long-tailed learning with self-adjusting gradient balancer,” in *Proceedings of the 37th International Conference on Neural Information Processing Systems*, ser. NIPS ’23. Red Hook, NY, USA: Curran Associates Inc., 2023.
- [42] M. Mortaheb, C. Vahapoglu, and S. Ulukus, “Fedgradnorm: Personalized federated gradient-normalized multi-task learning,” in *2022 IEEE 23rd International Workshop on Signal Processing Advances in Wireless Communication (SPAWC)*, 2022, pp. 1–5.
- [43] M. N. Fekri, K. Grolinger, and S. Mir, “Asynchronous adaptive federated learning for distributed load forecasting with smart meter data,” *International Journal of Electrical Power & Energy Systems*, vol. 153, p. 109285, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0142061523003423>

APPENDIX A

VERIFIABLE RANDOM FUNCTION (VRF)

A standard VRF protocol comprises two entities: a prover and a verifier [39]. The prover generates a VRF output $\beta = \text{HASH}(sk, \alpha)$ and an associated proof $\pi = \text{Proof}(sk, \alpha)$, where sk denotes the secret key and α is an input provided by the user. These values, β and π , are sent to the verifier. The verifier initially checks whether β is consistent with the proof by evaluating $\beta = \text{VRF_proof_to_hash}(\pi)$. A mismatch at this stage suggests a potential security breach. To authenticate further, the prover reveals their public key pk , enabling the verifier to perform a formal check using $\text{Verify}(pk, \alpha, \pi)$. This function returns true if the proof π correctly validates β with respect to pk and α ; otherwise, it returns false.

APPENDIX B

LOOKUP TABLE FOR MINIMUM CLIENT SELECTION THRESHOLD IN CLUSTER-ORIENTED SETTING

TABLE XIV

MINIMUM REQUIRED CLIENT SELECTION THRESHOLD C TO ENSURE $P(\text{HONEST CLIENTS} < 2) \leq \delta$, FOR VARIOUS VALUES OF COLLUSION PROBABILITY ϕ .

| $\delta \backslash \phi$ | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 |
|--------------------------|-----|-----|-----|-----|-----|
| 0.05 | 3 | 4 | 5 | 6 | 8 |
| 0.01 | 4 | 5 | 7 | 8 | 11 |
| 0.001 | 5 | 7 | 9 | 11 | 14 |

Each entry in Table XIV represents the minimum required client selection threshold C for a given combination of collusion probability ϕ and risk tolerance δ . Specifically, each cell gives the smallest value of C such that the probability of selecting fewer than 2 honest clients from a group of C selected clients is at most δ . This assumes that each client in the cluster independently colludes with the aggregator with probability ϕ . Practitioners can use this table to choose an appropriate value of C by identifying their system’s estimated collusion probability and desired level of privacy assurance.

APPENDIX C CONVERGENCE ANALYSIS

To thoroughly analyze the convergence of AdRo-FL under static quantization (referring to section VI-A), we begin by introducing following assumptions:

Assumption 1 (Bounded Gradients): We assume that the local gradients for all clients are bounded, which means there exists a constant G such that:

$$\|\nabla f_i(\Theta)\|_2 \leq G, \quad \forall i, \forall \Theta. \quad (8)$$

This assumption prevents excessively large updates, ensuring stability during the optimization process.

Assumption 2 (Bounded Variance of Stochastic Gradients): We assume the variance of gradients computed by the selected subset S_t of clients at round t is bounded as follows:

$$\mathbb{E} \left[\left\| \frac{1}{|S_t|} \sum_{i \in S_t} \nabla f_i(\Theta) - \nabla F(\Theta) \right\|_2^2 \right] \leq \frac{\sigma^2}{|S_t|}. \quad (9)$$

This assumption ensures the updates are not excessively noisy.

Assumption 3 (L-Smoothness of the Global Loss): We assume the global loss function $F(\Theta)$ is L -smooth, which implies:

$$\|\nabla F(\Theta) - \nabla F(\Theta')\|_2 \leq L\|\Theta - \Theta'\|_2, \quad \forall \Theta, \Theta'. \quad (10)$$

This guarantees smooth and gradual changes in the gradient, aiding stable convergence.

Assumption 4 (Selection Bias): We consider the bias R_t introduced by client selection at each round:

$$\mathbb{E} \left[\frac{1}{|S_t|} \sum_{i \in S_t} \nabla f_i(\Theta) \right] = \nabla F(\Theta) + R_t. \quad (11)$$

Here, R_t represents the bias due to the non-uniform selection of clients.

A. Step 1: Update Rule with Static Quantization

The global model update incorporating static quantization is:

$$\Theta^{t+1} = \Theta^t - \eta \frac{1}{|S_t|} \sum_{i \in S_t} \nabla f_i(\Theta^t), \quad (12)$$

where η is the learning rate. The static quantization affects only the transmitted payload size, not explicitly the gradients in this formulation.

B. Step 2: Applying the L-Smoothness Property

Applying the L -smoothness definition to our update step, we obtain:

$$F(\Theta^{t+1}) \leq F(\Theta^t) + \nabla F(\Theta^t)^T (\Theta^{t+1} - \Theta^t) + \frac{L}{2} \|\Theta^{t+1} - \Theta^t\|^2 \quad (13)$$

$$= F(\Theta^t) - \eta \nabla F(\Theta^t)^T \frac{1}{|S_t|} \sum_{i \in S_t} \nabla f_i(\Theta^t) \quad (14)$$

$$+ \frac{L\eta^2}{2} \left\| \frac{1}{|S_t|} \sum_{i \in S_t} \nabla f_i(\Theta^t) \right\|^2. \quad (15)$$

This describes how the loss changes at each iteration, highlighting the effect of the learning rate and gradients.

C. Step 3: Taking Expectation

Taking the expectation conditioned on the current parameters Θ^t , we have:

$$\mathbb{E}[F(\Theta^{t+1})|\Theta^t] \leq F(\Theta^t) - \eta \nabla F(\Theta^t)^T (\nabla F(\Theta^t) + R_t) \quad (16)$$

$$+ \frac{L\eta^2}{2} \mathbb{E} \left[\left\| \frac{1}{|S_t|} \sum_{i \in S_t} \nabla f_i(\Theta^t) \right\|^2 \right]. \quad (17)$$

This accounts explicitly for bias due to client selection.

D. Step 4: Bounding the Gradient Norm Term

Using bounded variance and gradient assumptions, we bound the gradient norm as follows:

$$\mathbb{E} \left[\left\| \frac{1}{|S_t|} \sum_{i \in S_t} \nabla f_i(\Theta^t) \right\|^2 \right] \leq \frac{\sigma^2}{|S_t|} + \|\nabla F(\Theta^t)\|^2. \quad (18)$$

This step ensures stable updates despite the stochastic nature of client sampling.

E. Step 5: Deriving the Final Convergence Bound

Combining the above results, the final expected decrease in the loss becomes:

$$\mathbb{E}[F(\Theta^{t+1})|\Theta^t] \leq F(\Theta^t) - \eta \|\nabla F(\Theta^t)\|^2 + \eta \|\nabla F(\Theta^t)\| \|R_t\| \quad (19)$$

$$+ \frac{L\eta^2}{2} \left(\frac{\sigma^2}{|S_t|} + \|\nabla F(\Theta^t)\|^2 \right). \quad (20)$$

This explicitly captures the effects of bias R_t , variance σ^2 , and learning rate η on the convergence.

Provided that the selection bias R_t diminishes or remains small, and given appropriate choices of η , the global model parameters Θ will converge towards a stationary point.