

---

# Semi-Synchronous Hierarchies and Credibility Management for Robust Federated Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In synchronous Federated Learning (FL) architectures, three major operational  
2 challenges persistently emerge: the stragglers’ effect, which significantly impedes  
3 aggregate computation efficiency; network congestion that compromises commu-  
4 nication efficacy; and the vulnerability to poisoning attacks that endangers model  
5 integrity. In response to these critical issues, this paper introduces a novel FL  
6 framework named Clustered Semi-synchronous Hierarchical Federated Learning  
7 (CSS-HFL). It utilizes edge servers to synchronously train models with their respec-  
8 tive clustered clients, which are clustered based on their computational capabilities  
9 and network conditions. As for the cloud server, a semi-synchronous training  
10 scheme is adopted to defend cloud aggregation against adversarial attacks. To  
11 bolster the robustness of CSS-HFL against poisoning attacks, we propose a new  
12 algorithm, Fusion Credibility (FusCred), which leverages a credibility scoring sys-  
13 tem and a small clean dataset on the cloud server to filter out potentially malicious  
14 updates. We provide a theoretical convergence guarantee and efficiency analysis for  
15 CSS-HFL and extensive experiments on MNIST, FMNIST, and CIFAR-10 datasets  
16 under various attack scenarios to demonstrate its effectiveness. Our results show  
17 that CSS-HFL with FusCred significantly enhances model accuracy and robustness  
18 compared to state-of-the-art FL algorithms. For example, on the Non-IID CIFAR-  
19 10 dataset, FusCred showcased an improvement in accuracy of 17.7%, 17.8% and  
20 10.4%, respectively, over the state-of-the-art algorithm when exposed to three types  
21 of model poisoning attacks in experiments with 40% attackers.

## 22 1 Introduction

23 Federated Learning (FL) [21] is a decentralized machine learning paradigm that enables end devices  
24 to train models locally and share only the parameter updates, thereby alleviating concerns regarding  
25 data privacy [30] and legal compliance [41]. A typical FL framework, such as federated averaging  
26 (FedAvg) [30], trains a global model by iteratively aggregating local updates from many clients  
27 synchronously. This framework is widely adopted in various applications, yet it faces several  
28 challenges that hinder its practical implementation.

29 **Stragglers effect:** firstly, because of the computational capacity and the network constraints, slow  
30 clients (*i.e.*, stragglers) require more time to train local models. This makes normal clients waste  
31 a great deal of time to wait stragglers, which is called *stragglers effect* [4, 25, 33]. To mitigate the  
32 impact of stragglers, [30, 55, 20] aggregate local updates only from a delicately selected subset of  
33 clients. Nevertheless, due to the Non-IID (not identically and independently distributed) distribution,  
34 the absence of excluding clients can greatly reduce the global model performance. Additionally,  
35 Xie *et al.* [51] propose Asynchronous FL framework, where the server can aggregate with the first  
36 received local update without waiting for the lagging devices.

37 **Network congestion:** second, in practical applications, the network condition becomes a bottleneck  
 38 when a significant amount of end devices collaborate to train the global model under Synchronous FL  
 39 framework [11]. Liu *et. al.* [26] introduced Hierarchical Federated Learning (HFL) to relieve the  
 40 congestion on the backbone network. A client-edge-cloud HFL architecture can greatly decrease the  
 41 model training time and the energy consumption of the clients compared to traditional FL [27].

42 **Poisoning attacks:** third, due to its special framework, traditional FL faces some severe security  
 43 problems if some clients are malicious. For instance, malicious clients could upload modified  
 44 parameters (*i.e.*, model poisoning) [39, 56] or dirty training data (*i.e.*, data poisoning) [17, 2]. The  
 45 global model performance would be degraded even though only one single malicious client in  
 46 traditional FL [12]. Actually, some robust algorithms [13, 31, 43, 58] are proposed to protect global  
 47 models against adversarial attacks. Nevertheless, all of them are based on Synchronous FL which  
 48 means they cannot fit perfectly with Asynchronous FL.

49 However, the existing FL frameworks and FL algorithms can only address part of drawbacks of  
 50 typical FL framework. The summarization of the limitation is shown in Table 1. In the real-world FL  
 51 system implementation, we should comprehensively deal with stragglers effect, network constraints,  
 52 and malicious attacks. An urgent need thus arises to propose a new FL framework to simultaneously  
 address the above problems.

Table 1: The limitations of existing FL frameworks.

Limitations	Framework		
	Synchronous FL	Asynchronous FL	HFL
Stragglers effect	✗	✓	✗
Network congestion	✗	✓	✓
Poisoning attacks	✓	✗	✓

53

54 In this paper, we introduce a novel FL framework termed **C**lustered **S**emi-synchronous **H**ierarchical  
 55 **F**ederated **L**earning (CSS-HFL). Within this framework, clients are organized into distinct clusters  
 56 according to their computational capacities and network conditions to mitigate the stragglers effect.  
 57 Then, the edge servers engage in training with their respective clients utilizing Synchronous FL  
 58 methods (*e.g.* Fed-Credit [10], Median [58], GeoMed [13]). Subsequently, the cloud server strategi-  
 59 cally determines the timing for aggregating edge models, ensuring that each edge server has recently  
 60 completed its aggregation. Noticeably there is no requirement for uniform epochs across all edge  
 61 servers. Either the Semi-synchronous FL in the cloud layer or the synchronous FL in the edge layer  
 62 can be accessed to apply robust Synchronous FL algorithms to resist malicious attacks. Additionally,  
 63 under CSS-HFL framework, we propose a robust algorithm named **F**usion **C**redibility (FusCred).  
 64 This algorithm leverages Fed-Credit [10] in the edge layer and maintains a small clean dataset on the  
 65 cloud server. The updated models of each edge cluster are assigned a credit score by comparing it with  
 66 the model trained on the cloud dataset. Subsequently, only the top  $k$  edge parameters are aggregated  
 67 with the cloud parameter. We provide both the convergence guarantee and efficiency analysis of  
 68 CSS-HFL, followed by the efficiency simulation and comprehensive experiments conducted on  
 69 MNIST, FMNIST, and CIFAR-10 datasets. Our experiments encompassed various attack types, ratios  
 70 of malicious clients, and dataset distributions. The empirical findings unequivocally showcase that  
 71 our proposed *FusCred* not only preserves high test accuracies but also exhibits exceptional robustness  
 72 against adversarial attacks. Our main contributions can be summarized as follows:

- 73 • To the best of our knowledge, this is the first work to explore both the robustness and  
 74 efficiency of HFL. By leveraging semi-synchronized aggregation and adaptive clustering,  
 75 CSS-HFL framework is proposed to comprehensively address the limitations of existing FL  
 76 frameworks, when dealing with stragglers effect, network congestion and poisoning attacks.
- 77 • We derive the efficiency analysis of CSS-HFL in comparison with famous FL frameworks  
 78 and prove the convergence guarantee in CSS-HFL framework. We also conducted an  
 79 efficiency simulation to show that our CSS-HFL can significantly enhance efficiency by  
 80 involving few edge servers.
- 81 • Within the CSS-HFL, we design a novel defense algorithm named *FusCred*. *FusCred* utilizes  
 82 Fed-Credit on edge servers and maintains a small clean dataset on the cloud server to assign  
 83 credit scores to edge model updates. This ensures that only top  $k$  edge parameters are  
 84 aggregated with the cloud parameter to mitigate attack effects passed over the edge. Notably,

85 FusCred demonstrates superior performance across various scenarios, outperforming state-  
86 of-the-art algorithms.

- 87 • The extensive comparative experiments between *FusCred* and various prior algorithms  
88 validate its effectiveness. Specifically, on the Non-IID CIFAR-10 dataset, our algorithm  
89 exhibited performance enhancements of 17.7%, 17.8%, and 10.4%, respectively, in compari-  
90 son to the state-of-the-art algorithm when facing three types of model poisoning attacks in  
91 experiments involving 40% attackers.

## 92 2 Observation And Threat Model

93 In this section, we first briefly introduce the observation of extant FL frameworks. Then we will  
94 describe the threat model of Hierarchical Federated Learning (HFL) system considered in this work.

### 95 2.1 Observation

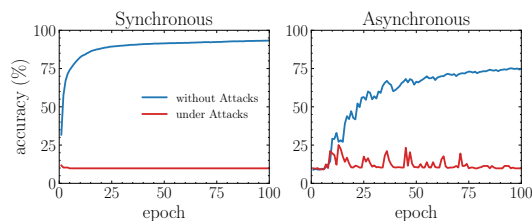


Figure 1: The accuracy of Synchronous FL and Asynchronous FL under no attacks or 20% attacks in Three-tier HFL on the Non-IID Mnist dataset.

96 In Figure 1, we briefly investigate the resilience to malicious attacks in HFL, using both Synchronous  
97 FL and Asynchronous FL. In the experiment, we set the attack as Sign-Flip (SF), in which the  
98 malicious clients upload the local updates by flipping the sign of each number. We assess the accuracy  
99 of Synchronous FedAvg and Asynchronous FedAsync in both attacks-free and 20% attacks scenarios  
100 under the Three-tier FL framework, employing the Non-IID MNIST dataset. A glance at the Figure 1  
101 reveals the same trends between attacks-free and 20% attacks. All of them see a plunge in accuracy  
102 compared to no attacks scenario, which is not acceptable in practical application, when facing attacks.  
103 More comparisons about existing FL framework can be found in Appendix A.

### 104 2.2 Threat Model

105 In this section, we present a comprehensive threat model for poisoning attacks within the context of  
106 CSS-HFL.

107 **Poisoning Attacker’s Goal:** Aligned with numerous prior studies on poisoning attacks [15, 50],  
108 the primary objective of the poisoning attacker in CSS-HFL is to deliberately manipulate the local  
109 training process. Their ultimate aim is to compromise the aggregation process of the global model.

110 **Types of Poisoning Attacks:** The strategies employed in our attacks align with those detailed in the  
111 work of Fed-Credit [10], encompassing data poisoning attacks [17, 2] and model poisoning attacks  
112 [39, 56].

113 **Poisoning Attacker’s Knowledge:** The poisoning attackers are indeed components of CSS-HFL,  
114 possessing specific knowledge within the framework. As clients, they have access to important  
115 information including the training data, model structure, learning algorithms, and the global model.  
116 This knowledge equips them to conduct their attacks effectively within the system.

117 **Poisoning Attacker’s Assumptions:** 1) Poisoning attackers are capable of collaborating with one  
118 another, thereby enabling them to coordinate and execute the same type of attack collectively. 2)  
119 Poisoning attackers are constrained to conducting their operations solely on the client side, implying  
120 that both the edge and cloud components are deemed trustworthy. 3) It is assumed that the number of  
121 malicious clients does not surpass half of the total [15]. 4) We assume that the network communication  
122 in CSS-HFL is reliable.

123 **3 Clustered Semi-synchronous HFL Framework (CSS-HFL)**

124 In this section, we introduce the CSS-HFL framework (Figure 2). CSS-HFL mainly addresses three  
 125 goals: 1) **To mitigate the waiting time of clients.** 2) **To relieve the network congestion.** 3) **To**  
 126 **provide an interface for robust FL algorithms.**

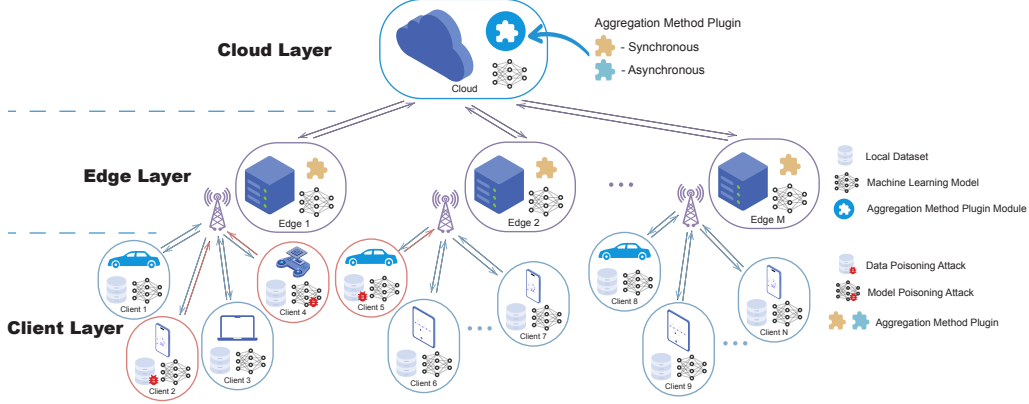


Figure 2: The CSS-HFL Framework.

127 We will begin by describing the components of the CCS-HFL framework. Overall, we adopt  
 128 hierarchical federated learning to enhance communication efficiency and release network congestion  
 129 [26]. 1) **At the client layer**, since we have  $N$  edge servers, we cluster all the clients into  $N$  clusters,  
 130 each belonging to one of  $N$  edges. The cluster criterion is based on the computation capacity and  
 131 network condition [34, 45, 7]. Our objective in this stage is to reduce waiting time of clients within  
 132 respective clusters. 2) **At the edge layer**, each edge server conducts synchronous federated learning  
 133 with its participating clients. During the edge aggregation stage, the edge server can select a robust  
 134 aggregation algorithm to protect the edge model from the attacks of malicious clients. Training at the  
 135 edge server resembles the traditional federated learning. 3) **At the cloud server**, the cloud can choose  
 136 either a different or the same secure aggregation algorithm used by the edges. It is significant for  
 137 cloud server to carefully determine the timing, when each edge server has recently completed an edge  
 138 aggregation, (note: edge servers are not mandated to go through the same number of communication  
 139 rounds with respective clients), to aggregate edge models. It is noteworthy that the semi-synchronous  
 140 aggregation scheme provides interfaces to different robustness algorithms, where the users have the  
 141 flexibility to choose the appropriate algorithm. The overall algorithm of our proposed CSS-HFL  
 142 framework can be found in Algorithm 1.

143 **3.1 Fusion Credibility (FusCred) in CSS-HFL**

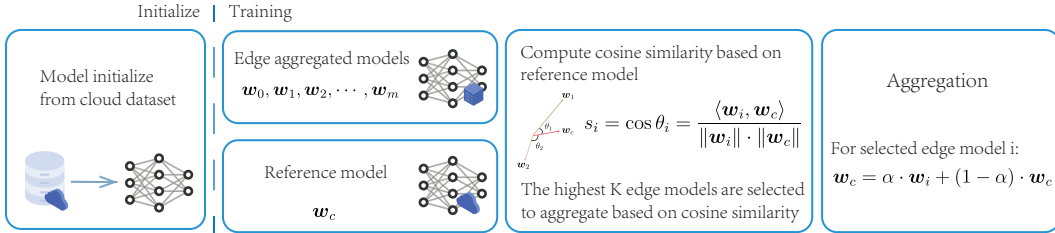


Figure 3: Cloud Aggregation Algorithm

144 Under the CSS-HFL framework, we propose a more robust aggregation algorithm named *FusCred*,  
 145 which comprises both edge aggregation method and cloud aggregation method. This algorithm can  
 146 maintain the efficiency of CSS-HFL while offering a high level of resilience against attacks.

147 In a macroscopic view, we use a non-discriminatory aggregation algorithm at each edge and a  
 148 discriminatory one in the cloud. For the edges, with their narrow perspective limited to the few client  
 149 models they can observe, this non-discriminatory aggregation preserves data diversity and some  
 150 resistance to attack. The cloud, with access to all edge gradients and cloud dataset references, uses a

---

**Algorithm 1: CSS-HFL Training Process**

---

**Input** :  $n$  clients with local training datasets,  $C_1, C_2, C_3, \dots, C_n$ ;  $N$  edge servers,  $\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3, \dots, \mathcal{E}_N$ ; learning rate  $lr$ ; batch size  $B$ ; number of local training iterations  $E$ ; number of cloud communication rounds  $R$ .

**Output** : Convergent cloud model  $w$ .

```
1 The cloud server utilizes the Balanced Clustering Algorithm[45] to form  $N$  clusters by grouping
  clients based on their computation capacities and network conditions. The number of clients in
  the  $\lambda^{th}$  group is denoted as  $N_\lambda$ .
2 The cloud server assigns an edge server to each cluster and defines the communication rounds for
  each edge server denoted as  $E_1, E_2, E_3, \dots, E_N$ .
3 Cloud server excutes:
4    $w \leftarrow$  pre-train model.
5   for  $i_R$  in  $R$  do
6     The cloud sends cloud model  $w$  to  $\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3, \dots, \mathcal{E}_N$ .
7     Receive edge models  $w_1, w_2, w_3, \dots, w_N$ .
8      $w \leftarrow$  Cloud aggregation( $w_1, w_2, w_3, \dots, w_N$ ).
9   end
10 Edge server excutes:
11 for  $\lambda = 1$  to  $N$  parallel do
12   Receive cloud model  $w$ .
13    $w_\lambda \leftarrow w$ .
14   // Semi-synchronous lies in varied edge communication rounds  $E_\lambda$ .
15   for  $i_\lambda$  in  $E_\lambda$  do
16     for  $k = 1$  to  $N_\lambda$  parallel do
17       /* Client executes */
18        $w_{\lambda,k} \leftarrow w_\lambda$ .
19       for  $i_E = 1$  to  $E$  do
20          $w_{\lambda,k} \leftarrow$  SGD( $w_{\lambda,k}$ , local dataset).
21       end
22       The  $k^{th}$  client in the  $\lambda^{th}$  edge server uploads its local model  $w_{\lambda,k}$  to the  $\lambda^{th}$ 
       edge server.
23     end
24      $w_\lambda \leftarrow$  Edge aggregation( $w_{\lambda,1}, w_{\lambda,2}, w_{\lambda,3}, \dots, w_{\lambda,k}$ ).
25   end
26   The  $\lambda^{th}$  edge server uploads its edge model  $w_\lambda$  to the cloud server.
27 end
```

---

151 discriminatory aggregation algorithm to filter out compromised edges. Together, these two methods  
152 allow the global model to converge stably.

153 For edge aggregation algorithm, we employ Fed-Credit [10], the recently proposed robust algorithm  
154 that currently works well at Two-tier FL. Next, we will delve into cloud aggregation algorithm in  
155 detail.

156 As illustrated in Figure 3, the Cloud Aggregation Algorithm is primarily divided into two sections.  
157 In the Initialization phase, the cloud-side dataset is utilized to train the global model for a specified  
158 number of epochs. In the Training phase, the cloud sends the global model to each edge, where  
159 it is used to train the edge-side models,  $w_i$ . Concurrently, the cloud trains for a specified number  
160 of epochs on the cloud dataset based on the global model to obtain a reference model, denoted by  
161  $w_c$ . Thereafter, the cloud calculates the credibility of each edge's updated model according to the  
162 following equation:  $s_i = \cos \theta_i = \frac{\langle w_i, w_c \rangle}{\|w_i\| \cdot \|w_c\|}$ . Subsequently, the most credible edges are selected  
163 for aggregation. Finally, the updates from the selected  $K$  edges are aggregated according to the  
164 following equation:  $w_c = \alpha \cdot w_j + (1 - \alpha) \cdot w_c$ . It should be noted that the aggregation order is  
165 randomised. The training process is repeated until the global model converges or reaches a preset  
166 number of epochs. The pseudo-code for this algorithm can be found in Algorithm 2.

---

**Algorithm 2:** Cloud Aggregation Method

---

**Input :** Cloud dataset, cloud model  $w$ ,  $N$  edge aggregated models  $w_1, w_2, w_3, \dots, w_N$ , aggregate proportion  $p$ , initial epochs  $E_i$ , cloud reference epochs  $E_r$ , cloud communications rounds  $R$ .

**Output :** Convergent cloud model  $w$ .

```
1 for  $epoch = 1$  to  $E_i$  do
2    $w \leftarrow \text{SGD}(w, \text{cloud dataset})$ .
3 end
4 for  $r$  in  $R$  do
5   The cloud sends cloud model  $w$  to all edge.
6   Each edge aggregates their clients' updates and returns new edge models  $w_1,$ 
    $w_2, w_3, \dots, w_N$  to the cloud.
7    $w_c = w$ .
8   for  $epoch = 1$  to  $E_r$  do
9      $w_c \leftarrow \text{SGD}(w_c, \text{cloud dataset})$ .
10  end
11  for  $i = 1$  to  $N$  do
12    Compute cosine similarity  $s_i = \cos \theta_i = \frac{\langle w_i, w_c \rangle}{\|w_i\| \cdot \|w_c\|}$ .
13  end
14  for  $j = 1$  to  $N$  do
15    Compute the rank of  $s_j$  in  $s_1, s_2, s_3, \dots, s_N$ .
16    if  $\text{rank} \geq p \times N$  then
17       $w_c \leftarrow \alpha \cdot w_j + (1 - \alpha) \cdot w_c$ .
18    end
19  end
20   $w \leftarrow w_c$ .
21 end
22 return Cloud model  $w$ .
```

---

### 167 3.2 Efficiency Analysis

168 We introduce a metric called *Average Waiting Time* (AWT), which aims to calculate the average  
169 waiting time of all end devices, to assess the efficiency of the framework. The less value of AWT  
170 indicates higher efficiency of framework. We neglect the time taken for edge aggregation, cloud  
171 aggregation, and communication between the edge server and the cloud server. AWT calculates the  
172 average waiting time across all end devices during one cloud aggregation.

173 Let  $t_{\lambda,k}$  denotes the total training time, including local training and communication overhead, of  $k^{\text{th}}$   
174 client in the  $\lambda^{\text{th}}$  edge. We designate  $T = \max\{t_{\lambda,k}\}$  for all  $1 \leq \lambda \leq N$  and  $1 \leq k \leq N_\lambda$  as the  
175 slowest client among all  $n$  clients, and  $T_\lambda = \max\{T_{\lambda,k}\}$  for all  $1 \leq k \leq N_\lambda$  as the slowest client in  
176 the  $\lambda^{\text{th}}$  edge.  $\Delta_\lambda$  represents the idle time between training of  $\lambda^{\text{th}}$  edge server in FedAT and HiFlash.

Table 2: Average waiting time (AWT) comparisons of various FL frameworks.

Framework	AWT Expression	Framework	AWT Expression
FedAsync [51]	0	FedSync [26]	$\frac{1}{n} \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} (T - t_{\lambda,k})$
FedAT [9], HiFlash [46]	$\frac{1}{n} \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} (T_\lambda - t_{\lambda,k} + \Delta_\lambda)$	CSS-HFL	$\frac{1}{n} \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} (T_\lambda - t_{\lambda,k})$

177 The AWT comparison among various frameworks is shown in Table 2. Since  $T_\lambda \leq T$  and  $\Delta \geq 0$ ,  
178 we have  $AWT_{\text{Asynchronous}} \leq AWTCSS \leq AWTSynchronous$  and  $AWTCSS \leq AWTFedAT, AWTHiFlash$ .  
179 It's important to note that as  $N$  approaches 1, CSS-HFL behaves like Synchronous FL, while as  $N$   
180 approaches  $n$ , it behaves like Asynchronous FL. In summary, the efficiency of CSS-HFL lies between  
181 Synchronous FL and Asynchronous FL, with the choice of  $N$  playing a significant role, and is higher  
182 than FedAT and HiFlash. Experimental results about efficiency can be found in Appendix F.1.

## 183 4 Evaluation

### 184 4.1 Experimental Setup

185 In our experiments, we evaluate CSS-HFL with **FedAvg** [30], **A-Krum** [43], **Median** [58], **GeoMed**  
186 [13], and on MNIST, Fashion-MNIST, and CIFAR-10 datasets under both IID and Non-IID settings.  
187 We utilized Dirichlet distribution to model Non-IID distribution [59]. For each scenario, we take  
188 the average of results of three seeds (2023, 2024, 3047). Experiments are conducted on a server  
189 comprises the AMD EPYC 7742 64-Core Processor and the NVIDIA Tesla A100 40G computing  
190 accelerator.

191 **Datasets and Networks.** A detailed description of the dataset can be found in the Appendix F.2.  
192 For MNIST, we adopt a Multi-Layer Perceptron (MLP) network with two hidden layers and one  
193 output layer to train the model. For Fashion-MNIST, a 7-layer LeNet [22] with convolutional layers is  
194 employed for model training. For CIFAR-10, we opt for a lightweight model Compact Convolutional  
195 Transformers (CCT) [18] due to its compact design and effectiveness, which holds promise for  
196 mitigating the resource constraints in onboard FL end devices.

197 **Attacks.** Our model poisoning attacks are implemented in three distinct forms: Constant Parameter  
198 (CP), where all model parameters remain identical; Normal Parameter (NP), which generates model  
199 parameters following a normal distribution; and Sign-Flip Parameter (SF), producing a model with  
200 parameters opposite to those obtained during training. As for the data poisoning attack, we choose  
201 the based on pairwise (PW) and symmetric (SM) matrices to flip the training labels. Additionally,  
202 20%, 30%, and 40% attack ratios are adopted to evaluate the resilience of algorithms and model  
203 the attackers distribution by Dirichlet ( $G \sim DP(\alpha, G_0)$ ) with three  $\alpha$ s (0.2, 0.5, 0.8). A larger  $\alpha$   
204 corresponds to a distribution closer to uniform, while a smaller  $\alpha$  indicates a more concentrated  
205 distribution.

206 **Evaluation metric.** To assess the performance of the multiple defense algorithms under CSS-HFL  
207 framework, as many prior studies [8, 54], we employ *accuracy* as a key criterion. Higher accuracy  
208 signifies better defense.

### 209 4.2 Results And Analysis

210 We demonstrate the partial accuracy results of our experiments in Table 3 and Table 4. Table 3  
211 presents the average accuracy of different attack types under varying attack ratios. Table 4 displays  
212 the accuracy of various attack types in the presence of 30% attacks.

213 **Impact of Ratio of Malicious Clients.** Firstly, as illustrated in Table 3, A-Krum exhibits lower  
214 accuracy compared to other methods with the absence of attacks. This distinction is particularly  
215 obvious when the dataset distribution is Non-IID. The other methods achieve relatively higher  
216 accuracy. This phenomenon might be because A-Krum tend to heavily rely on few local updates to  
217 update the global update, which lead to the global model cannot fitting the overall dataset well. A  
218 discernible pattern emerges in the results: an increase in the ratio of malicious clients corresponds to  
219 a noticeable decline in accuracy and a growth in bias. Our *FusCred* outperforms other approaches.  
220 Notably, the *FusCred* consistently achieves higher accuracy levels and maintains fewer instances of  
221 extreme variability. This reinforces the assertion that *FusCred* adeptly preserves both accuracy and  
222 stability, even in the presence of an escalating ratio of adversarial entities.

223 **Impact of Attack Types.** Table 4 illustrates the varying effectiveness of different aggregation  
224 approaches against a range of attack techniques on the FMNIST dataset with diverse distributions.  
225 Notably, distinct patterns emerge, particularly in scenarios with high ratios of attackers. For example,  
226 GeoMed performs well under 20% and 30% ratio attacks, similar to *FusCred*, but experiences a  
227 significant decline when facing 40% ratio attacks. Another finding is the compared methods do  
228 not exhibit comprehensive robustness across various attack types. For instance, both FedAvg and  
229 Median perform poorly when facing SF attacks compared to other attack types, while A-Krum only  
230 demonstrates better tolerance for SF attacks. Furthermore, our *FusCred* consistently perform well,  
231 effectively mitigating all types of attacks with higher accuracy compared to alternative methods  
232 considered.

Table 3: Comparing accuracies (%) under various attack ratios. Gold, silver, and bronze respectively denote the top three winners.

Dataset	Attack ratio	FedAvg	Median	GeoMed	A-Krum	FusCred
IID MNIST	0%	98.35±0.04	98.15±0.03	98.31±0.02	93.55±0.24	97.49±0.13
	20%	95.64±2.85	96.35±1.33	97.08±0.86	75.80±32.98	97.41±0.09
	30%	89.14±12.67	93.95±3.56	96.44±1.21	75.92±33.11	97.47±0.15
	40%	82.15±16.81	88.02±10.17	90.92±5.78	63.11±34.34	97.40±0.16
Non-IID MNIST	0%	98.20±0.01	97.98±0.01	98.11±0.01	74.26±0.60	96.86±0.11
	20%	92.14±8.29	95.08±2.94	97.05±0.56	61.78±15.97	96.86±0.28
	30%	89.25±8.52	89.45±9.78	94.60±2.79	49.30±32.25	96.06±1.26
	40%	78.55±16.12	76.33±18.19	75.24±17.68	47.79±31.57	95.81±1.50
IID FMNIST	0%	90.07±0.07	89.70±0.08	90.20±0.06	84.33±0.46	89.08±0.17
	20%	86.27±2.40	86.82±2.64	88.98±0.47	77.94±14.74	88.55±0.29
	30%	79.05±11.18	79.83±14.84	87.97±0.58	73.11±24.00	88.76±0.21
	40%	72.47±16.75	69.66±20.52	79.03±9.32	52.64±29.84	88.53±0.40
Non-IID FMNIST	0%	89.57±0.09	89.18±0.03	89.75±0.10	75.94±0.06	87.04±0.48
	20%	81.51±6.95	82.31±6.76	88.17±0.59	73.47±7.05	86.92±0.84
	30%	73.69±12.96	73.38±17.78	85.87±2.46	56.97±20.62	86.91±0.75
	40%	63.17±19.55	58.91±25.20	72.48±12.05	37.54±27.79	85.95±1.34
IID CIFAR-10	0%	66.07±0.03	65.63±0.02	65.98±0.10	51.79±0.35	62.79±0.20
	20%	48.12±11.74	47.17±13.55	58.95±4.38	51.68±0.96	62.37±0.47
	30%	43.35±12.21	41.56±15.51	47.91±12.52	51.11±1.12	62.15±0.46
	40%	36.89±12.52	34.60±16.37	34.65±18.05	47.94±6.04	61.62±0.46
Non-IID CIFAR-10	0%	66.08±0.07	65.58±0.05	65.87±0.09	51.97±0.01	62.53±0.12
	20%	48.37±11.36	48.60±12.44	60.48±3.96	51.60±0.93	62.73±0.29
	30%	42.61±12.73	41.57±15.97	48.79±11.89	51.57±1.46	62.38±0.40
	40%	36.21±12.24	34.88±16.33	36.33±17.55	47.48±5.89	61.72±0.39

## 233 5 Related Work

234 **Poisoning attacks.** According to Xia *et al.* [47], poisoning attacks can be classified into two  
235 categories: *data poisoning attacks* and *model poisoning attacks*. In data poisoning attacks, malicious  
236 clients have the ability to inject poisoned information into training data or labels. [40, 17, 5, 32, 35, 48]  
237 propose label flipping to attack the models by manipulating labels. Specifically, symmetric flipping  
238 [40] and pairwise flipping [17] are introduced to flip each label to other labels via a specific transition  
239 matrix, significantly enhancing the efficiency of label flipping attacks. [61, 60, 38, 37, 2] focus  
240 on the training data poisoning. They carefully craft the training data or generate fake data with  
241 aim of undermining the performance of the global model. On the other hand, Model poisoning  
242 [3, 56, 39, 49, 19, 23] directly manipulates clients to upload arbitrary or counterfeit local updates  
243 which poses significant threats to the global model.

244 **Robust FL.** In recent years, multiple robust FL algorithms in Synchronous FL have emerged. Broadly,  
245 these algorithms can be categorized into the following two groups. 1) **Discarding rules** detect and  
246 exclude potential attackers when aggregating global models. Krum [6], Multi-Krum [6], Bulyan  
247 [31], A-Krum [43], Trimmed-Mean [58], and MAB-RFL [42] are represent of discarding algorithms.  
248 While discarding algorithms excel in defending against attacks, the removal of partial clients can  
249 be detrimental, especially in cases of non-IID data distribution or when the number of clients is  
250 limited. 2) **Non-discarding rules** aim to leverage all the information of local updates instead of  
251 directly dropping out potential threat clients. Zeno [52], Zeno++ [53], Fed-Credit [10], and FLTrust  
252 [8] assign weights to each candidate local updates. Suspicious clients are assigned lower weights,  
253 while benign clients receive higher weights. GeoMed [13], Median [58], RFA [36] and FoolsGold  
254 [16] utilize statistical characteristics of updates to update the global parameters.

255 **Efficient FL.** Traditional FL is susceptible to the stragglers effect and network congestion. Algorithms  
256 such as FedAsync [51], FedSA [29], ASO-Fed [14], Async-FedED [44], and DP-AFL [28] have  
257 been introduced to mitigate the first issue. These algorithms enable the aggregator to update without  
258 waiting for lagging or lost clients, thereby saving training time. Asynchronous FL, however, can not  
259 effectively address peak network congestion. Additionally, Multi-tier FL (*e.g.* HFL [26], FedAT [9],



Table 4: Comparing robust accuracies (%) under 30% attacks of various types.

Dataset	Attack Type		FedAvg	Median	GeoMed	A-krum	FusCred
IID MNIST	Model poison	CP	95.41±0.12	95.73±0.22	95.93±0.22	92.06±0.52	97.42±0.15
		NP	95.18±0.12	94.40±0.02	94.55±0.27	9.70±0.22	97.47±0.07
		SF	89.23±3.95	87.55±2.18	96.33±0.28	92.78±0.11	97.28±0.09
	Data poison	PW	70.96±18.57	95.21±2.11	97.63±0.10	92.77±0.55	97.56±0.03
		SM	94.91±1.68	96.85±0.10	97.78±0.06	92.28±0.49	97.61±0.09
Non-IID MNIST	Model poison	CP	89.55±0.24	91.93±1.45	91.90±2.60	9.80±0.00	96.73±0.08
		NP	91.70±0.68	91.82±1.98	93.37±2.20	10.16±0.26	96.80±0.07
		SF	74.89±8.93	71.96±8.50	93.07±1.04	75.72±4.21	96.41±0.10
	Data poison	PW	95.43±1.12	95.13±1.46	97.32±0.17	73.55±4.21	94.86±2.16
		SM	94.69±1.07	96.38±0.02	97.37±0.09	77.27±1.77	95.51±0.61
IID FMNIST	Model poison	CP	83.59±0.49	85.46±0.62	88.21±0.07	83.93±0.81	88.85±0.11
		NP	83.73±0.65	85.13±0.61	87.79±0.14	26.26±11.51	88.91±0.10
		SF	64.42±9.73	55.05±18.14	88.51±0.28	86.46±0.09	88.47±0.22
	Data poison	PW	76.32±14.11	86.65±1.11	87.57±0.92	84.38±0.26	88.85±0.08
		SM	87.19±0.53	86.84±0.51	87.78±0.37	84.51±0.56	88.71±0.18
Non-IID FMNIST	Model poison	CP	73.97±2.23	76.65±2.71	85.24±1.15	37.11±9.12	87.21±0.20
		NP	75.40±1.62	75.26±3.14	84.91±1.20	32.20±0.40	87.31±0.24
		SF	51.01±9.66	42.47±16.05	83.70±3.54	78.22±0.57	85.82±0.84
	Data poison	PW	83.18±2.39	86.02±1.92	87.42±1.29	75.60±2.36	87.08±0.36
		SM	84.91±0.74	86.51±0.26	88.08±0.07	61.74±14.38	87.15±0.60
IID CIFAR-10	Model poison	CP	36.78±1.63	35.45±3.01	34.36±3.80	51.73±0.23	62.07±0.52
		NP	35.99±1.63	34.07±2.19	33.38±3.52	51.72±0.22	61.87±0.52
		SF	28.61±2.77	19.98±1.19	49.39±2.19	51.78±0.27	62.71±0.24
	Data poison	PW	56.00±0.45	57.16±0.89	60.46±0.92	50.02±1.34	62.08±0.05
		SM	59.40±0.33	61.13±0.59	61.93±0.42	50.32±1.12	62.01±0.23
Non-IID CIFAR-10	Model poison	CP	35.82±1.12	35.94±1.77	35.37±5.61	51.96±0.02	62.42±0.20
		NP	35.01±0.84	33.92±1.50	37.25±3.49	51.67±0.15	62.42±0.20
		SF	26.99±2.38	18.92±3.61	48.74±6.03	52.62±0.49	62.71±0.11
	Data poison	PW	56.10±0.66	57.65±1.00	60.71±0.52	51.11±1.48	62.41±0.13
		SM	59.13±0.17	61.43±0.29	61.89±0.25	50.51±2.38	61.93±0.61

260 FedEdge [43], HiFlash [46]) combined with cluster algorithms (e.g., FL+HC [7], ClusterFL [34],  
 261 FedCH [45]) have been proposed to address both the stragglers effect and network congestion. To  
 262 our best of knowledge, nevertheless, few works have considered the robustness in efficient FL.

## 263 6 Discussion

264 The *FusCred* consistently outperforms alternative methods across various scenarios. In real-world  
 265 implementation, however, it may be necessary for servers to change aggregation methods for some  
 266 concerns. The edge servers are advised to apply non-discarding robustness algorithms. Given their  
 267 limited visibility, edge servers can only observe their own client models, making it challenging to  
 268 determine whether an outlier client model is due to model diversity or malicious attacks. In contrast,  
 269 the cloud server is recommended to implement discarding robustness algorithms. Since it can perceive  
 270 more information from edge models and it is easier for cloud server to discern whether an outlier  
 271 model indicates a malicious attack.

## 272 7 Conclusion

273 In this paper, we propose CSS-HFL, a novel FL framework that can simultaneously handle with  
 274 stragglers effect, network congestion, and poisoning attacks. The theoretical proofs demonstrating  
 275 the efficiency and convergence guarantee of CSS-HFL are provided. Additionally, within CSS-HFL,  
 276 we design a robust aggregation algorithm, named *FusCred*, outperforming alternative methods in  
 277 defending against adversarial attacks, as exhibited in extensive experiments.

278 **References**

- 279 [1] M. S. H. Abad, E. Ozfatura, D. GÜndÜz, and O. Ercetin. Hierarchical federated learning across  
280 heterogeneous cellular networks. In *ICASSP 2020 - 2020 IEEE International Conference on*  
281 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 8866–8870, 2020.
- 282 [2] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How  
283 to backdoor federated learning. In Silvia Chiappa and Roberto Calandra, editors, *The 23rd*  
284 *International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 August*  
285 *2020, Online [Palermo, Sicily, Italy]*, volume 108 of *Proceedings of Machine Learning Research*,  
286 pages 2938–2948. PMLR, 2020.
- 287 [3] Gilad Baruch, Moran Baruch, and Yoav Goldberg. A little is enough: Circumventing defenses  
288 for distributed learning. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence  
289 d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information*  
290 *Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019,*  
291 *NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 8632–8642, 2019.
- 292 [4] Amir Behrouzi-Far and Emina Soljanin. On the effect of task-to-worker assignment in dis-  
293 tributed computing systems with stragglers. In *2018 56th Annual Allerton Conference on*  
294 *Communication, Control, and Computing (Allerton)*, pages 560–566, 2018.
- 295 [5] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector  
296 machines. In *Proceedings of the 29th International Conference on Machine Learning, ICML*  
297 *2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress, 2012.
- 298 [6] Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning  
299 with adversaries: Byzantine tolerant gradient descent. In Isabelle Guyon, Ulrike von Luxburg,  
300 Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett,  
301 editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural*  
302 *Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages  
303 119–129, 2017.
- 304 [7] Christopher Briggs, Zhong Fan, and Peter Andras. Federated learning with hierarchical cluster-  
305 ing of local updates to improve training on non-iid data. In *2020 International Joint Conference*  
306 *on Neural Networks (IJCNN)*, pages 1–9, 2020.
- 307 [8] Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. Fltrust: Byzantine-robust  
308 federated learning via trust bootstrapping. In *28th Annual Network and Distributed System*  
309 *Security Symposium, NDSS 2021, virtually, February 21-25, 2021*. The Internet Society, 2021.
- 310 [9] Zheng Chai, Yujing Chen, Ali Anwar, Liang Zhao, Yue Cheng, and Huzefa Rangwala. Fedat: a  
311 high-performance and communication-efficient federated learning system with asynchronous  
312 tiers. In Bronis R. de Supinski, Mary W. Hall, and Todd Gamblin, editors, *International*  
313 *Conference for High Performance Computing, Networking, Storage and Analysis, SC 2021, St.*  
314 *Louis, Missouri, USA, November 14-19, 2021*, page 60. ACM, 2021.
- 315 [10] Jiayan Chen, Zhirong Qian, Tianhui Meng, Xitong Gao, Tian Wang, and Weijia Jia. Fed-credit:  
316 Robust federated learning with credibility management. *CoRR*, abs/2405.11758, 2024.
- 317 [11] Mingzhe Chen, Zhaohui Yang, Walid Saad, Changchuan Yin, H. Vincent Poor, and Shuguang  
318 Cui. A joint learning and communications framework for federated learning over wireless  
319 networks. *IEEE Transactions on Wireless Communications*, 20(1):269–283, 2021.
- 320 [12] Yudong Chen, Lili Su, and Jiaming Xu. Distributed statistical machine learning in adversarial  
321 settings: Byzantine gradient descent. In Konstantinos Psounis, Aditya Akella, and Adam  
322 Wierman, editors, *Abstracts of the 2018 ACM International Conference on Measurement and*  
323 *Modeling of Computer Systems, SIGMETRICS 2018, Irvine, CA, USA, June 18-22, 2018*,  
324 page 96. ACM, 2018.
- 325 [13] Yudong Chen, Lili Su, and Jiaming Xu. Distributed statistical machine learning in adversarial  
326 settings: Byzantine gradient descent. In Konstantinos Psounis, Aditya Akella, and Adam  
327 Wierman, editors, *Abstracts of the 2018 ACM International Conference on Measurement and*

- 328 *Modeling of Computer Systems, SIGMETRICS 2018, Irvine, CA, USA, June 18-22, 2018,*  
329 page 96. ACM, 2018.
- 330 [14] Yujing Chen, Yue Ning, Martin Slawski, and Huzefa Rangwala. Asynchronous online federated  
331 learning for edge devices with non-iid data. In Xintao Wu, Chris Jermaine, Li Xiong, Xiaohua  
332 Hu, Olivera Kotevska, Siyuan Lu, Weija Xu, Srinivas Aluru, Chengxiang Zhai, Eyhab Al-Masri,  
333 Zhiyuan Chen, and Jeff Saltz, editors, *2020 IEEE International Conference on Big Data (IEEE  
334 BigData 2020), Atlanta, GA, USA, December 10-13, 2020*, pages 15–24. IEEE, 2020.
- 335 [15] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. Local model poisoning  
336 attacks to byzantine-robust federated learning. In Srdjan Capkun and Franziska Roesner,  
337 editors, *29th USENIX Security Symposium, USENIX Security 2020, August 12-14, 2020*, pages  
338 1605–1622. USENIX Association, 2020.
- 339 [16] Clement Fung, Chris J. M. Yoon, and Ivan Beschastnikh. The limitations of federated learning  
340 in sybil settings. In *23rd International Symposium on Research in Attacks, Intrusions and  
341 Defenses (RAID 2020)*, pages 301–316, San Sebastian, October 2020. USENIX Association.
- 342 [17] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor W. Tsang, and  
343 Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy  
344 labels. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-  
345 Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31:  
346 Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December  
347 3-8, 2018, Montréal, Canada*, pages 8536–8546, 2018.
- 348 [18] Ali Hassani, Steven Walton, Nikhil Shah, Abulikemu Abuduweili, Jiachen Li, and Humphrey  
349 Shi. Escaping the big data paradigm with compact transformers. *CoRR*, abs/2104.05704, 2021.
- 350 [19] Sai Praneeth Karimireddy, Lie He, and Martin Jaggi. Learning from history for byzantine robust  
351 optimization. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International  
352 Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of  
353 *Proceedings of Machine Learning Research*, pages 5311–5319. PMLR, 2021.
- 354 [20] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and  
355 Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learning.  
356 In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference  
357 on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages  
358 5132–5143. PMLR, 13–18 Jul 2020.
- 359 [21] Jakub Konečný, H. Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated  
360 optimization: Distributed machine learning for on-device intelligence. *CoRR*, abs/1610.02527,  
361 2016.
- 362 [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep  
363 convolutional neural networks. In Peter L. Bartlett, Fernando C. N. Pereira, Christopher J. C.  
364 Burges, Léon Bottou, and Kilian Q. Weinberger, editors, *Advances in Neural Information  
365 Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems  
366 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States*,  
367 pages 1106–1114, 2012.
- 368 [23] Liping Li, Wei Xu, Tianyi Chen, Georgios B. Giannakis, and Qing Ling. RSA: byzantine-robust  
369 stochastic aggregation methods for distributed learning from heterogeneous datasets. In *The  
370 Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative  
371 Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on  
372 Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January  
373 27 - February 1, 2019*, pages 1544–1551. AAAI Press, 2019.
- 374 [24] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence  
375 of fedavg on non-iid data. In *8th International Conference on Learning Representations, ICLR  
376 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- 377 [25] Xingyu Li, Zhe Qu, Bo Tang, and Zhuo Lu. Stragglers are not disasters: A hybrid federated  
378 learning framework with delayed gradients. In *2022 21st IEEE International Conference on  
379 Machine Learning and Applications (ICMLA)*, pages 727–732, 2022.

- 380 [26] Lumin Liu, Jun Zhang, Shenghui Song, and Khaled B. Letaief. Client-edge-cloud hierarchical  
381 federated learning. In *2020 IEEE International Conference on Communications, ICC 2020,*  
382 *Dublin, Ireland, June 7-11, 2020*, pages 1–6. IEEE, 2020.
- 383 [27] Lumin Liu, Jun Zhang, Shenghui Song, and Khaled B. Letaief. Hierarchical federated learning  
384 with quantization: Convergence analysis and system design. *IEEE Trans. Wirel. Commun.*,  
385 22(1):2–18, 2023.
- 386 [28] Yunlong Lu, Xiaohong Huang, Yueyue Dai, Sabita Maharjan, and Yan Zhang. Differentially  
387 private asynchronous federated learning for mobile edge computing in urban informatics. *IEEE*  
388 *Trans. Ind. Informatics*, 16(3):2134–2143, 2020.
- 389 [29] Qianpiao Ma, Yang Xu, Hongli Xu, Zhida Jiang, Liusheng Huang, and He Huang. Fedrsa: A  
390 semi-asynchronous federated learning mechanism in heterogeneous edge computing. *IEEE J.*  
391 *Sel. Areas Commun.*, 39(12):3654–3672, 2021.
- 392 [30] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas.  
393 Communication-efficient learning of deep networks from decentralized data. In Aarti Singh  
394 and Xiaojin (Jerry) Zhu, editors, *Proceedings of the 20th International Conference on Artificial*  
395 *Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA,*  
396 volume 54 of *Proceedings of Machine Learning Research*, pages 1273–1282. PMLR, 2017.
- 397 [31] El Mahdi El Mhamdi, Rachid Guerraoui, and Sébastien Rouault. The hidden vulnerability of  
398 distributed learning in byzantium. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings*  
399 *of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan,*  
400 *Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*,  
401 pages 3518–3527. PMLR, 2018.
- 402 [32] Konda Reddy Mopuri, Vaisakh Shaj, and R. Venkatesh Babu. Adversarial fooling beyond  
403 "flipping the label". In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition,*  
404 *CVPR Workshops 2020, Seattle, WA, USA, June 14-19, 2020*, pages 3374–3382. Computer  
405 Vision Foundation / IEEE, 2020.
- 406 [33] Diego Cardoso Nunes, Bruno Loureiro Coelho, Ricardo Parizotto, and Alberto Schaeffer-Filho.  
407 Serene: Handling the effects of stragglers in in-network machine learning aggregation. In  
408 *NOMS 2023-2023 IEEE/IFIP Network Operations and Management Symposium*, pages 1–10,  
409 2023.
- 410 [34] Xiaomin Ouyang, Zhiyuan Xie, Jiayu Zhou, Guoliang Xing, and Jianwei Huang. Clusterfl: A  
411 clustering-based federated learning system for human activity recognition. *ACM Trans. Sen.*  
412 *Netw.*, 19(1), dec 2022.
- 413 [35] Andrea Paudice, Luis Muñoz-González, and Emil C. Lupu. Label sanitization against label  
414 flipping poisoning attacks. In Carlos Alzate, Anna Monreale, Haytham Assem, Albert Bifet,  
415 Teodora Sandra Buda, Bora Caglayan, Brett Drury, Eva García-Martín, Ricard Gavaldà, Stefan  
416 Kramer, Niklas Lavesson, Michael Madden, Ian M. Molloy, Maria-Irina Nicolae, and Math-  
417 ieu Sinn, editors, *ECML PKDD 2018 Workshops - Nemesis 2018, UrbReas 2018, SoGood*  
418 *2018, IWAISe 2018, and Green Data Mining 2018, Dublin, Ireland, September 10-14, 2018,*  
419 *Proceedings*, volume 11329 of *Lecture Notes in Computer Science*, pages 5–15. Springer, 2018.
- 420 [36] Venkata Krishna Pillutla, Sham M. Kakade, and Zaïd Harchaoui. Robust aggregation for  
421 federated learning. *CoRR*, abs/1912.13445, 2019.
- 422 [37] Dazhong Rong, Shuai Ye, Ruoyan Zhao, Hon Ning Yuen, Jianhai Chen, and Qinming He.  
423 Fedrecattack: Model poisoning attack to federated recommendation. In *38th IEEE International*  
424 *Conference on Data Engineering, ICDE 2022, Kuala Lumpur, Malaysia, May 9-12, 2022*, pages  
425 2643–2655. IEEE, 2022.
- 426 [38] Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor  
427 Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural  
428 networks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-  
429 Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31:*  
430 *Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December*  
431 *3-8, 2018, Montréal, Canada*, pages 6106–6116, 2018.

- 432 [39] Virat Shejwalkar and Amir Houmansadr. Manipulating the byzantine: Optimizing model  
433 poisoning attacks and defenses for federated learning. In *28th Annual Network and Distributed*  
434 *System Security Symposium, NDSS 2021, virtually, February 21-25, 2021*. The Internet Society,  
435 2021.
- 436 [40] Brendan van Rooyen, Aditya Krishna Menon, and Robert C. Williamson. Learning with  
437 symmetric label noise: The importance of being unhinged. In Corinna Cortes, Neil D. Lawrence,  
438 Daniel D. Lee, Masashi Sugiyama, and Roman Garnett, editors, *Advances in Neural Information*  
439 *Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015,*  
440 *December 7-12, 2015, Montreal, Quebec, Canada*, pages 10–18, 2015.
- 441 [41] Paul Voigt and Axel Von dem Bussche. The EU general data protection regulation (GDPR). *A*  
442 *Practical Guide, 1st Ed., Cham: Springer International Publishing*, 10:3152676, 2017.
- 443 [42] Wei Wan, Shengshan Hu, Jianrong Lu, Leo Yu Zhang, Hai Jin, and Yuanyuan He. Shielding  
444 federated learning: Robust aggregation with adaptive client selection. In Luc De Raedt, editor,  
445 *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI*  
446 *2022, Vienna, Austria, 23-29 July 2022*, pages 753–760. ijcai.org, 2022.
- 447 [43] Kaibin Wang, Qiang He, Feifei Chen, Hai Jin, and Yun Yang. Fededge: Accelerating edge-  
448 assisted federated learning. In *Proceedings of the ACM Web Conference 2023, WWW '23*, page  
449 2895–2904, New York, NY, USA, 2023. Association for Computing Machinery.
- 450 [44] Qiyuan Wang, Qianqian Yang, Shibo He, Zhiguo Shi, and Jiming Chen. Asyncfed: Asyn-  
451 chronous federated learning with euclidean distance based adaptive weight aggregation. *CoRR*,  
452 abs/2205.13797, 2022.
- 453 [45] Zhiyuan Wang, Hongli Xu, Jianchun Liu, Yang Xu, He Huang, and Yangming Zhao. Accelerat-  
454 ing federated learning with cluster construction and hierarchical aggregation. *IEEE Transactions*  
455 *on Mobile Computing*, 22(7):3805–3822, 2023.
- 456 [46] Qiong Wu, Xu Chen, Tao Ouyang, Zhi Zhou, Xiaoxi Zhang, Shusen Yang, and Junshan Zhang.  
457 Hiflash: Communication-efficient hierarchical federated learning with adaptive staleness control  
458 and heterogeneity-aware client-edge association. *IEEE Transactions on Parallel and Distributed*  
459 *Systems*, 34(5):1560–1579, 2023.
- 460 [47] Geming Xia, Jian Chen, Chaodong Yu, and Jun Ma. Poisoning attacks in federated learning: A  
461 survey. *IEEE Access*, 11:10708–10722, 2023.
- 462 [48] Han Xiao, Huang Xiao, and Claudia Eckert. Adversarial label flips attack on support vector  
463 machines. In Luc De Raedt, Christian Bessiere, Didier Dubois, Patrick Doherty, Paolo Frasconi,  
464 Fredrik Heintz, and Peter J. F. Lucas, editors, *ECAI 2012 - 20th European Conference on*  
465 *Artificial Intelligence. Including Prestigious Applications of Artificial Intelligence (PAIS-2012)*  
466 *System Demonstrations Track, Montpellier, France, August 27-31, 2012*, volume 242 of  
467 *Frontiers in Artificial Intelligence and Applications*, pages 870–875. IOS Press, 2012.
- 468 [49] Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. Fall of empires: Breaking byzantine-  
469 tolerant SGD by inner product manipulation. In Amir Globerson and Ricardo Silva, editors,  
470 *Proceedings of the Thirty-Fifth Conference on Uncertainty in Artificial Intelligence, UAI 2019,*  
471 *Tel Aviv, Israel, July 22-25, 2019*, volume 115 of *Proceedings of Machine Learning Research*,  
472 pages 261–270. AUAI Press, 2019.
- 473 [50] Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. Fall of empires: Breaking byzantine-  
474 tolerant sgd by inner product manipulation. In Ryan P. Adams and Vibhav Gogate, editors,  
475 *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, volume 115 of  
476 *Proceedings of Machine Learning Research*, pages 261–270. PMLR, 22–25 Jul 2020.
- 477 [51] Cong Xie, Sanmi Koyejo, and Indranil Gupta. Asynchronous federated optimization. *CoRR*,  
478 abs/1903.03934, 2019.
- 479 [52] Cong Xie, Sanmi Koyejo, and Indranil Gupta. Zeno: Distributed stochastic gradient descent  
480 with suspicion-based fault-tolerance. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors,  
481 *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June*

- 482 2019, Long Beach, California, USA, volume 97 of *Proceedings of Machine Learning Research*,  
483 pages 6893–6901. PMLR, 2019.
- 484 [53] Cong Xie, Sanmi Koyejo, and Indranil Gupta. Zeno++: Robust fully asynchronous SGD. In  
485 *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18*  
486 *July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages  
487 10495–10503. PMLR, 2020.
- 488 [54] Haonan Yan, Wenjing Zhang, Qian Chen, Xiaoguang Li, Wenhai Sun, Hui Li, and Xiaodong  
489 Lin. RECESS vaccine for federated learning: Proactive defense against model poisoning attacks.  
490 In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine,  
491 editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*  
492 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -*  
493 *16, 2023*, 2023.
- 494 [55] Haibo Yang, Minghong Fang, and Jia Liu. Achieving linear speedup with partial worker partici-  
495 pation in non-IID federated learning. In *International Conference on Learning Representations*,  
496 2021.
- 497 [56] Ming Yang, Hang Cheng, Fei Chen, Ximeng Liu, Meiqing Wang, and Xibin Li. Model poisoning  
498 attack in differential privacy-based federated learning. *Information Sciences*, 630:158–172,  
499 2023.
- 500 [57] Zhengjie Yang, Sen Fu, Wei Bao, Dong Yuan, and Albert Y. Zomaya. Hierarchical federated  
501 learning with momentum acceleration in multi-tier networks. *IEEE Transactions on Parallel*  
502 *and Distributed Systems*, 34(10):2629–2641, 2023.
- 503 [58] Dong Yin, Yudong Chen, Kannan Ramchandran, and Peter L. Bartlett. Byzantine-robust  
504 distributed learning: Towards optimal statistical rates. In Jennifer G. Dy and Andreas Krause,  
505 editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018,*  
506 *Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine*  
507 *Learning Research*, pages 5636–5645. PMLR, 2018.
- 508 [59] Dun Zeng, Siqi Liang, Xiangjing Hu, Hui Wang, and Zenglin Xu. Fedlab: A flexible federated  
509 learning framework. *J. Mach. Learn. Res.*, 24:100:1–100:7, 2023.
- 510 [60] Jiale Zhang, Bing Chen, Xiang Cheng, Huynh Thi Thanh Binh, and Shui Yu. Poisongan:  
511 Generative poisoning attacks against federated learning in edge computing systems. *IEEE*  
512 *Internet Things J.*, 8(5):3310–3322, 2021.
- 513 [61] Jiale Zhang, Junjun Chen, Di Wu, Bing Chen, and Shui Yu. Poisoning attack in federated  
514 learning using generative adversarial nets. In *18th IEEE International Conference On Trust,*  
515 *Security And Privacy In Computing And Communications / 13th IEEE International Conference*  
516 *On Big Data Science And Engineering, TrustCom/BigDataSE 2019, Rotorua, New Zealand,*  
517 *August 5-8, 2019*, pages 374–380. IEEE, 2019.

518 **A Additional Discussions**

519 **A.1 Synchronous v.s. Asynchronous**

520 Synchronous FL is a typical architecture in which a server distributes the global model to a selected  
 521 subset of clients and does not update the global model until it receives all local updates. However the  
 522 server can update global model without waiting for the lagging clients in Asynchronous FL.

523 Due to the synchronous nature of FL, numerous previous robust FL algorithms have been proposed  
 524 which heavily rely on comparisons of local updates. However, waiting for stragglers or offline  
 525 clients can lead to substantial costs. Asynchronous FL significantly enhances convergence efficiency  
 526 compared to Synchronous FL. Nevertheless, defending against malicious attacks becomes challenging  
 527 when updating the global model with just one local update.

528 **A.2 Two-tier v.s. Three-tier**

529 In Two-tier FL, multiple clients are directly connected to a remote server or cloud, which suffers  
 530 from peak network congestion in both Synchronous FL and Asynchronous FL. With the development  
 531 of edge computing, an edge tier is added between the local clients and remote cloud to alleviate the  
 532 strain caused by peak network congestion. In a Three-tier HFL, the clients can first communicate  
 533 with the edge node for edge-level aggregation. Subsequently, the edge nodes communicate with the  
 534 remote cloud for cloud-level aggregation.

535 The Three-tier architecture presents a promising solution for real-world large-scale clients and has  
 536 captivated significant attention from researchers [1, 7, 57]. To the best of our knowledge, however,  
 537 none of these works have focused on security, which poses significant threats to convergence, privacy,  
 538 economics, and even life security.

539 **B Notations**

Table 5: Key Notations For The Clustered Semi-synchronous HFL Framework.

Denote	Description	Denote	Description
$n$	The number of clients	$N$	The number of edges
$N_\lambda$	The number clients of $\lambda^{th}$ edge	$C_i$	The $i^{th}$ client ( $1 \leq i \leq n$ )
$\mathcal{E}_\lambda$	The $\lambda^{th}$ edge ( $1 \leq \lambda \leq N$ )	$lr$	Learning rate
$B$	Batch size	$E$	Number of client training epochs
$E_\lambda$	Number of $\lambda^{th}$ edge training epochs	$R$	Number of cloud training epochs
$w, w_\lambda, w_{\lambda,k}$	The model of cloud, $\lambda^{th}$ edge, $k^{th}$ client in $\lambda^{th}$ edge ( $1 \leq k \leq N_\lambda$ )		

540 **C Convergence Analysis**

541 Our convergence analysis is inspired by [24]. We first make the following assumptions.

542 **Assumption 1.** *The loss functions  $F$  in the cloud server, the edge server, and the client are all*  
 543  *$L$ -smooth:  $\forall v, w, F(w) \leq F(v) + (w - v)^\top \nabla F(v) + \frac{L}{2} \|w - v\|_2^2$ .*

544 **Assumption 2.** *The loss functions  $F$  in the cloud server, the edge server, and the client are all*  
 545  *$\mu$ -strongly convex:  $\forall v, w, F(w) \geq F(v) + (w - v)^\top \nabla F(v) + \frac{\mu}{2} \|w - v\|_2^2$ .*

546 **Assumption 3.** *Let  $\xi_t^{\lambda,k}$  be sampled uniformly at random from local data of the  $k^{th}$  end de-*  
 547 *vice in the  $\lambda^{th}$  edge. The variance of stochastic gradients in each device is bounded as follows:*  
 548  *$\mathbb{E} \|\nabla F_{\lambda,k}(w_t^{\lambda,k}, \xi_t^{\lambda,k}) - \nabla F_{\lambda,k}(w_t^{\lambda,k})\|^2 \leq \sigma_{\lambda,k}^2$ , for  $1 \leq \lambda \leq N$  and  $1 \leq k \leq N_\lambda$ .*

549 **Assumption 4.** *The expected squared norm of stochastic gradients is uniformly bounded, i.e.,*  
 550  *$\mathbb{E} \|\nabla F_{\lambda,k}(w_t^{\lambda,k}, \xi_t^{\lambda,k})\|^2 \leq G^2$ , for  $1 \leq \lambda \leq N$ ,  $1 \leq k \leq N_\lambda$ , and  $1 \leq t \leq T$ .*

551 We define  $F^*$  and  $F_{\lambda,k^*}$  as the minimum value of  $F$  and  $F_{\lambda,k}$  and let  $\Lambda = \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} (F^* - F_{\lambda,k^*})$   
 552 quantify the degree of Non-IID [24]. We assume the cloud server aggregation interval is  $T_c$  and

553 the total number of rounds is  $T$ . Then, under our CSS-HFL framework, we have the following  
 554 convergence guarantee for FedAvg.

555 **Theorem 1.** Let (1) (2) (3) (4) hold and  $L, \mu, \sigma_{\lambda,k}^2, G$  be defined therein. Choose  $\tau = \frac{L}{\mu}$ ,  $\varphi =$   
 556  $\max\{8\tau, T_c\}$  and the learning rate  $\eta_t = \frac{2}{u(\varphi+t)}$ . Then our CSS-HFL framework satisfies

$$\mathbb{E}[F(\mathbf{w}_t)] - F^* \leq \frac{\tau}{\varphi + t - 1} \left( \frac{2\Upsilon}{\mu} + \frac{\mu\varphi}{2} \mathbb{E}\|\mathbf{w}_1 - \mathbf{w}^*\|^2 \right), \quad (1)$$

557 where

$$\Upsilon = \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda^2 p_{\lambda,k}^2 \sigma_{\lambda,k}^2 + 6L\Lambda + 8(T_c - 1)^2 G^2.$$

558 According to the above result, we observe that the right-hand side of the Equation (1) consists of  
 559 two terms. The first term,  $\tau/(\varphi + t - 1)$ , exhibits a decreasing trend concerning  $t$ . As  $t$  grows  
 560 sufficiently large, the constants  $\varphi$  and 1 can be disregarded, leading to an approximate form of  
 561  $\frac{\tau}{t} \left( \frac{2\Upsilon}{\mu} + \frac{\mu\tau}{2} \mathbb{E}\|\mathbf{w}_1 - \mathbf{w}^*\|^2 \right)$ . This implies a convergence rate of  $O(1/t)$ , indicating sub-linear  
 562 convergence.

## 563 D Proof of Convergence

Table 6: Table of Key Notations for Convergence Analysis.

Notation	Description
$N$	The number of edges
$N_\lambda$	The number clients of $\lambda^{th}$ edge
$p_\lambda$	The weight of $\lambda^{th}$ in the cloud aggregation
$p_{\lambda,k}$	The weight of $k^{th}$ client in the $\lambda^{th}$ edge aggregation
$\mathbf{w}^{\lambda,k}$	The model of $k^{th}$ client in $\lambda^{th}$ edge ( $1 \leq k \leq N_\lambda$ )
$\nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k})$	The gradient of $k^{th}$ client in $\lambda^{th}$ edge ( $1 \leq k \leq N_\lambda$ )
$\eta_t$	Learning rate of $t^{th}$ round
$T_c$	The aggregation interval of cloud server

564 We prove the Theorem 1 in this section. The key notations for convergence analysis is presented in  
 565 Table 6.

### 566 D.1 Additional Denotes

567 Let  $T_c, T_\lambda, T_{\lambda,k}$  be the aggregation interval of the cloud server, the  $\lambda^{th}$  edge, and the  $k^{th}$  client in the  
 568  $\lambda^{th}$  edge. Note that for cloud server and edges server, aggregations only occur if the remainder of  $t$   
 569 divided by the interval  $T_c$  or  $T_\lambda$  is 0,  $t$  is the current round. And for the  $k^{th}$  client in the  $\lambda^{th}$  edge, if the  
 570  $t^{th}$  round is not the aggregation round for the client (*i.e.*,  $t \bmod T_{\lambda,k} \neq 0$ ),  $\nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}) = 0$ ;  
 571 otherwise,  $\nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k})$  represent the gradient of the sampled mini-batch local dataset. We  
 572 also adopt the virtual sequence from [24] to represent the immediate result of  $t^{th}$  round. The above  
 573 note can be described as

$$\nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}) = \begin{cases} 0, & \text{if } t \bmod T_{\lambda,k} \neq 0, \\ \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}), & \text{if } t \bmod T_{\lambda,k} = 0. \end{cases} \quad (2)$$

$$\mathbf{v}_{t+1}^{\lambda,k} = \mathbf{w}_t^{\lambda,k} - \eta_t \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}). \quad (3)$$



$$\mathbf{w}_{t+1}^{\lambda,k} = \begin{cases} \mathbf{v}_{t+1}^{\lambda,k}, & \text{if } t+1 \bmod T_\lambda \neq 0 \text{ and } t+1 \bmod T_c \neq 0, \\ \sum_{k=1}^{N_\lambda} p_{\lambda,k} \mathbf{v}_{t+1}^{\lambda,k}, & \text{if } t+1 \bmod T_\lambda = 0 \text{ and } t+1 \bmod T_c \neq 0, \\ \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \mathbf{v}_{t+1}^{\lambda,k}, & \text{if } t+1 \bmod T_c = 0. \end{cases} \quad (4)$$

574 For convenience, we define  $\bar{\mathbf{v}}_t = \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \mathbf{v}_t^{\lambda,k}$ ,  $\bar{\mathbf{w}}_t = \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \mathbf{w}_t^{\lambda,k}$ ,  $\bar{\mathbf{g}}_t =$   
575  $\sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k})$ , and  $\mathbf{g}_t = \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k})$ . Therefore,  $\bar{\mathbf{v}}_{t+1} =$   
576  $\bar{\mathbf{w}}_t - \eta_t \mathbf{g}_t$  and  $\mathbb{E} \mathbf{g}_t = \bar{\mathbf{g}}_t$ .

## 577 D.2 Key Lemmas

578 In this section, we describe and proof the key useful lemmas.

579 **Lemma 1.** *Assume (3) holds, we have*

$$\mathbb{E} \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2 \leq \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda^2 p_{\lambda,k}^2 \sigma_{\lambda,k}^2.$$

580 *Proof.* From (3), the variance of the stochastic gradients in  $k^{\text{th}}$  client device in  $\lambda^{\text{th}}$  edge is bounded  
581 by  $\sigma_{\lambda,k}^2$ , then

$$\begin{aligned} \mathbb{E} \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2 &= \mathbb{E} \left\| \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} (\nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}) - \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k})) \right\|^2 \\ &= \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda^2 p_{\lambda,k}^2 \mathbb{E} \left\| \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}) - \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) \right\|^2 \\ &\leq \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda^2 p_{\lambda,k}^2 \sigma_{\lambda,k}^2. \end{aligned}$$

582

□

583 **Lemma 2.** *Assume (4) holds,  $\eta_t$  is non-increasing, and  $\eta_t \leq 2\eta_{t+T_c}$  for all  $t \geq 0$ . It follows that*

$$\mathbb{E} \left[ \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\| \bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k} \right\|^2 \right] \leq 4\eta_t^2 (T_c - 1)^2 G^2.$$

584 *Proof.*  $\forall t \geq 0, \exists t_0 \leq t$ , such that  $\mathbf{w}_{t_0}^{\lambda,k} = \bar{\mathbf{w}}_{t_0}$ , i.e., round  $t_0$  is the last cloud server aggregation  
585 round. Therefore we can indicate that  $t - t_0 \leq T_c - 1$ . Additionally, we utilize the assumptions that  
586  $\eta_t$  is non-increasing and  $\eta_{t_0} \leq 2\eta_{t_0+T_c} \leq 2\eta_t$ , then

$$\begin{aligned}
\mathbb{E} \left[ \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\| \bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k} \right\|^2 \right] &= \mathbb{E} \left[ \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\| (\mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_{t_0}) - (\bar{\mathbf{w}}_t - \bar{\mathbf{w}}_{t_0}) \right\|^2 \right] \\
&\leq \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \mathbb{E} \left\| \mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_{t_0} \right\|^2 \\
&\leq \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \sum_{t=t_0}^{t-1} (T_c - 1) \eta_t^2 \mathbb{E} \left\| \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}) \right\|^2 \\
&\leq \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \sum_{t=t_0}^{t-1} (T_c - 1) \eta_{t_0}^2 G^2 \\
&\leq \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \eta_{t_0}^2 (T_c - 1)^2 G^2 \\
&\leq 4 \eta_t^2 (T_c - 1)^2 G^2.
\end{aligned}$$

587 **First inequality:** We use the property of variance as follow

$$\mathbb{E} \|X - \mathbb{E}X\|^2 \leq \mathbb{E} \|X\|^2$$

588 where  $X = (\mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_{t_0})$ .

589 **Second inequality:** We use the Cauchy-Schwarz inequality and  $t - t_0 \leq T_c - 1$ .

$$\begin{aligned}
\mathbb{E} \left\| \mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_{t_0} \right\|^2 &= \mathbb{E} \left\| \sum_{t=t_0}^{t-1} \eta_t \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}) \right\|^2 \\
&\leq \sum_{t=t_0}^{t-1} (t - t_0) \eta_t^2 \mathbb{E} \left\| \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}) \right\|^2 \\
&\leq \sum_{t=t_0}^{t-1} (T_c - 1) \eta_t^2 \mathbb{E} \left\| \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}, \xi_t^{\lambda,k}) \right\|^2.
\end{aligned}$$

590 **Third inequality:** We leverage (3) and  $\eta_t$  is non-increasing (i.e.  $\eta_t \leq \eta_{t_0}$  for  $t \geq t_0$ ).

591 **Fourth inequality:** We utilize  $\eta_{t_0} \leq 2\eta_{t_0+T_c} \leq 2\eta_t$ . □

592 **Lemma 3.** We assume the (1), (2), and  $\eta_t = \frac{\alpha}{t+\varphi}$  for some  $\alpha > \frac{1}{\mu}$  and  $\varphi > 0$  such that  $\eta_1 \leq$   
593  $\min\{\frac{1}{\mu}, \frac{1}{4L}\} = \frac{1}{4L}$ , it follows that

$$\begin{aligned}
\mathbb{E} \left\| \bar{\mathbf{v}}_{t+1} - \mathbf{w}^* \right\|^2 &\leq (1 - \eta_t \mu) \mathbb{E} \left\| \bar{\mathbf{w}}_t - \mathbf{w}^* \right\|^2 + \eta_t^2 \mathbb{E} \left\| \mathbf{g}_t - \bar{\mathbf{g}}_t \right\|^2 + 6L \eta_t^2 \Lambda \\
&\quad + 2 \mathbb{E} \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\| \bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k} \right\|^2,
\end{aligned}$$

594 where  $\Lambda = \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} (F^* - F_{\lambda,k}^*) \geq 0$ .

595 *Proof.* We first divide  $\left\| \bar{\mathbf{v}}_{t+1} - \mathbf{w}^* \right\|^2$  into following three parts.

$$\begin{aligned}
\left\| \bar{\mathbf{v}}_{t+1} - \mathbf{w}^* \right\|^2 &= \left\| \bar{\mathbf{w}}_t - \eta_t \mathbf{g}_t - \mathbf{w}^* - \eta_t \bar{\mathbf{g}}_t + \eta_t \bar{\mathbf{g}}_t \right\|^2 \\
&= \underbrace{\left\| \bar{\mathbf{w}}_t - \mathbf{w}^* - \eta_t \bar{\mathbf{g}}_t \right\|^2}_{P_1} + \underbrace{2\eta_t \langle \bar{\mathbf{w}}_t - \mathbf{w}^* - \eta_t \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle}_{P_2} + \eta_t^2 \left\| \bar{\mathbf{g}}_t - \mathbf{g}_t \right\|^2. \quad (5)
\end{aligned}$$

596 Next, we focus on the  $P_1$ :

$$P_1 = \|\bar{\mathbf{w}}_t - \mathbf{w}^* - \eta_t \bar{\mathbf{g}}_t\|^2 = \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 - \underbrace{2\eta_t \langle \bar{\mathbf{w}}_t - \mathbf{w}^*, \bar{\mathbf{g}}_t \rangle}_{Q_1} + \underbrace{\eta_t^2 \|\bar{\mathbf{g}}_t\|^2}_{Q_2}. \quad (6)$$

597 We pay attention to  $Q_1$ :

$$\begin{aligned} Q_1 &= 2\eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\langle \bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k} + \mathbf{w}_t^{\lambda,k} - \mathbf{w}^*, \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) \right\rangle \\ &= 2\eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\langle \bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}, \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) \right\rangle \\ &\quad + 2\eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\langle \mathbf{w}_t^{\lambda,k} - \mathbf{w}^*, \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) \right\rangle \\ &\geq 2\eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \frac{1}{2} \left( -\frac{1}{\eta_t} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 - \eta_t \|\nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k})\|^2 \right) \\ &\quad + 2\eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\langle \mathbf{w}_t^{\lambda,k} - \mathbf{w}^*, \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) \right\rangle \\ &\geq 2\eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \frac{1}{2} \left( -\frac{1}{\eta_t} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 - \eta_t \|\nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k})\|^2 \right) \\ &\quad + 2\eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) - \nabla F_{\lambda,k}(\mathbf{w}^*) + \frac{\mu}{2} \|\mathbf{w}_t^{\lambda,k} - \mathbf{w}^*\|^2 \right). \end{aligned} \quad (7)$$

598 **First inequality:** We derive the first inequality in Equation (7) by Cauchy-Schwarz inequality and  
599 AM-GM inequality.

$$\left\langle \bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}, \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) \right\rangle \geq \frac{1}{2} \left( -\frac{1}{\eta_t} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 - \eta_t \|\nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k})\|^2 \right).$$

600 **Second inequality:** By the  $\mu$ -strong convexity of  $F_{\lambda,k}$ , we have

$$\left\langle \mathbf{w}_t^{\lambda,k} - \mathbf{w}^*, \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) \right\rangle \geq \left( \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) - \nabla F_{\lambda,k}(\mathbf{w}^*) + \frac{\mu}{2} \|\mathbf{w}_t^{\lambda,k} - \mathbf{w}^*\|^2 \right).$$

601 Then, we analyze  $Q_2$ , By the convexity of  $\|\cdot\|^2$  and the  $L$ -smoothness of  $F_{\lambda,k}$ , we have

$$\begin{aligned} Q_2 &= \eta_t^2 \|\bar{\mathbf{g}}_t\|^2 \leq \eta_t^2 \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\| \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) \right\|^2 \\ &\leq 2L\eta_t^2 \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) - F_{\lambda,k}(\mathbf{w}^*) \right) \end{aligned} \quad (8)$$

602 By combining Equation (7) and Equation (8), we have

$$\begin{aligned}
P_1 &\leq \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( \frac{1}{\eta_t} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 + \eta_t \|\nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k})\|^2 \right) \\
&\quad - 2\eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( \nabla F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) - \nabla F_{\lambda,k}(\mathbf{w}^*) + \frac{\mu}{2} \|\mathbf{w}_t^{\lambda,k} - \mathbf{w}^*\|^2 \right) \\
&\quad + 2L\eta_t^2 \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) - F_{\lambda,k}^* \right) \\
&\leq (1 - \mu\eta_t) \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 \\
&\quad + 4L\eta_t^2 \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) - F_{\lambda,k}^* \right) \tag{9} \\
&\quad - 2\eta_t \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) - F_{\lambda,k}(\mathbf{w}^*) \right) \\
&= (1 - \mu\eta_t) \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 \\
&\quad + 4L\eta_t^2 \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( F^* - F_{\lambda,k}^* \right) \\
&\quad + \underbrace{(4L\eta_t^2 - 2\eta_t) \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) - F^* \right)}_S,
\end{aligned}$$

603 where we use the  $L$ -smoothness of  $F_{\lambda,k}(\cdot)$  again and the following inequality for the second inequality,  
604

$$\begin{aligned}
\|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 &= \left\| \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} (\mathbf{w}_t^{\lambda,k} - \mathbf{w}^*) \right\|^2 \\
&\leq \left\| \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} (\sqrt{p_\lambda p_{\lambda,k}}) \right\|^2 \cdot \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} \left( \sqrt{p_\lambda p_{\lambda,k}} (\mathbf{w}_t^{\lambda,k} - \mathbf{w}^*) \right)^2 \\
&= \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \|\mathbf{w}_t^{\lambda,k} - \mathbf{w}^*\|^2
\end{aligned}$$

605 We next focus  $S$ ,

$$\begin{aligned}
S &= \left( \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( F_{\lambda,k}(\mathbf{w}_t^{\lambda,k}) - F_{\lambda,k}(\bar{\mathbf{w}}_t) \right) + \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left( F_{\lambda,k}(\bar{\mathbf{w}}_t) - F^* \right) \right) \\
&\geq \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left\langle \nabla F_{\lambda,k}(\bar{\mathbf{w}}_t), \mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_t \right\rangle + (F(\bar{\mathbf{w}}_t) - F^*) \\
&\geq -\frac{1}{2} \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left[ \eta_t \|\nabla F_{\lambda,k}(\bar{\mathbf{w}}_t)\|^2 + \frac{1}{\eta_t} \|\mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_t\|^2 \right] + (F(\bar{\mathbf{w}}_t) - F^*) \\
&\geq -\sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left[ \eta_t L (F_{\lambda,k}(\bar{\mathbf{w}}_t) - F_{\lambda,k}^*) + \frac{1}{2\eta_t} \|\mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_t\|^2 \right] + (F(\bar{\mathbf{w}}_t) - F^*). \tag{10}
\end{aligned}$$

606 The first inequality arises from the convexity of  $F_{\lambda,k}(\cdot)$ , the second inequality from the AM-GM  
 607 inequality, and the third inequality from the  $L$ -smoothness of  $F_{\lambda,k}$ .

608 By combining Equation (9) and Equation (10), and utilize the notation  $\Lambda = \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} (F^* -$   
 609  $F_{\lambda,k^*})$ , we have

$$\begin{aligned}
 P_1 &\leq (1 - \mu\eta_t) \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 + (4L\eta_t^2)\Lambda \\
 &\quad + (2\eta_t - 4L\eta_t^2) \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \left[ \eta_t L (F_{\lambda,k}(\bar{\mathbf{w}}_t) - F^* + F^* - F_{\lambda,k^*}) \right. \\
 &\quad \left. + \frac{1}{2\eta_t} \|\mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_t\|^2 \right] - (2\eta_t - 4L\eta_t^2) (F(\bar{\mathbf{w}}_t) - F^*) \\
 &= (1 - \mu\eta_t) \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 + (6L\eta_t^2 - 4L^2\eta_t^3)\Lambda \\
 &\quad + \frac{2\eta_t - 4L\eta_t^2}{2\eta_t} \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \|\mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_t\|^2 \\
 &\quad + (2\eta_t - 4L\eta_t^2)(\eta_t L - 1) \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} (F_{\lambda,k}(\bar{\mathbf{w}}_t) - F^*) \\
 &\leq (1 - \mu\eta_t) \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 + 6L\eta_t^2\Lambda \\
 &\quad + \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \|\mathbf{w}_t^{\lambda,k} - \bar{\mathbf{w}}_t\|^2 \\
 &= (1 - \mu\eta_t) \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + 2 \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} \|\bar{\mathbf{w}}_t - \mathbf{w}_t^{\lambda,k}\|^2 + (6L\eta_t^2)\Lambda,
 \end{aligned} \tag{11}$$

610 For the last inequality, we use the following facts:

611 1.  $\Lambda \geq 0$  and  $4L^2\eta_t^3 > 0$ .

612 2.  $\frac{2\eta_t - 4L\eta_t^2}{2\eta_t} \leq 1$ .

613 3.  $\eta_t L - 1 \leq 0$  and  $\sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_\lambda p_{\lambda,k} (F_{\lambda,k}(\bar{\mathbf{w}}_t) - F^*) = F(\bar{\mathbf{w}}_t) - F^* \geq 0$ .

614 Then, back to the Equation (5), notice that  $\mathbb{E}\|\mathbf{g}_t\| = \bar{\mathbf{g}}_t$ , i.e.,

$$\mathbb{E}\|P_2\| = 0. \tag{12}$$

615 Using the Equation (11) and Equation (12), we prove the Lemma 3.  $\square$

### 616 D.3 Proof of Theorem 1

617 *Proof.* It is evident that we always have  $\bar{\mathbf{w}}_t = \bar{\mathbf{v}}_t$ . For a non-increasing learning rate,  $\eta_t = \frac{\alpha}{t+\varphi}$  for  
 618 some  $\alpha > \frac{1}{\mu}$  and  $\varphi > 0$  such that  $\eta_1 \leq \min\{\frac{1}{\mu}, \frac{1}{4L}\} = \frac{1}{4L}$  and  $\eta_t \leq 2\eta_{t+T_c}$  for all  $t \geq 0$ . From  
 619 Lemma 1 Lemma 3, it follows that

$$\mathbb{E}\|\bar{\mathbf{w}}_{t+1} - \mathbf{w}^*\|^2 \leq (1 - \eta_t\mu)\mathbb{E}\|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \eta_t^2\Upsilon, \tag{13}$$

620 where

$$\Upsilon = \sum_{\lambda=1}^N \sum_{k=1}^{N_\lambda} p_{\lambda}^2 p_{\lambda,k}^2 \sigma_{\lambda,k}^2 + 6L\Lambda + 8(T_c - 1)^2 G^2.$$

621 Let  $\Delta_t = \mathbb{E} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2$ ,  $\zeta = \max \left\{ \frac{\alpha^2 \Upsilon}{\alpha\mu - 1}, (\varphi + 1)\Delta_1 \right\}$ . We can easily find that  $\Delta_1 \leq \frac{\zeta}{\varphi + 1}$  holds  
 622 for  $t = 1$ . Next, we prove  $\Delta_t \leq \frac{\zeta}{\varphi + t}$  by induction.

$$\begin{aligned} \Delta_{t+1} &\leq (1 - \eta_t \mu) \Delta_t + \eta_t^2 \Upsilon \\ &\leq \left(1 - \frac{\alpha\mu}{t + \varphi}\right) \frac{\zeta}{t + \varphi} + \frac{\alpha^2 \Upsilon}{(t + \varphi)^2} \\ &= \frac{t + \varphi - 1}{(t + \varphi)^2} \zeta + \left[ \frac{\alpha^2 \Upsilon}{(t + \varphi)^2} - \frac{\alpha\mu - 1}{(t + \varphi)^2} \zeta \right] \\ &\leq \frac{\zeta}{t + \varphi + 1}. \end{aligned}$$

623 If we set  $\alpha = \frac{2}{\mu}$ ,  $\varphi = \max \left\{ \frac{8L}{\mu}, T_c \right\} - 1$ , and define  $\tau = \frac{L}{\mu}$ , then  $\eta_t = \frac{2}{\mu(\varphi + t)}$ . It can be verified  
 624 that this choice of  $\eta_t$  satisfies  $\eta_t \leq 2\eta_{t+T_c}$  for  $t \geq 1$ . Thus, we obtain

$$\zeta = \max \left\{ \frac{\alpha^2 \Upsilon}{\alpha\mu - 1}, (\varphi + 1)\Delta_1 \right\} \leq \frac{\alpha^2 \Upsilon}{\alpha\mu - 1} + (\varphi + 1)\Delta_1 \leq \frac{4\Upsilon}{\mu^2} + (\varphi + 1)\Delta_1,$$

625 and by the  $L$ -smoothness of  $F(\cdot)$ , we have

$$\begin{aligned} \mathbb{E} \|F(\bar{\mathbf{w}}_t)\| - F^* &\leq \frac{L}{2} \Delta_t \leq \frac{L\zeta}{2(\varphi + t)} \leq \frac{\tau}{\varphi + t} \left( \frac{2\Upsilon}{\mu} + \frac{\mu(\varphi + 1)}{2} \Delta_1 \right) \\ &\leq \frac{\tau}{\varphi + t - 1} \left( \frac{2\Upsilon}{\mu} + \frac{\mu\varphi}{2} \mathbb{E} \|\bar{\mathbf{w}}_1 - \mathbf{w}^*\| \right). \end{aligned}$$

626 □

## 627 E Limitations

628 The experiments were conducted on the assumption that the number of attackers would remain  
 629 below 50%. Scenarios involving a higher number of attackers were not considered in the current  
 630 study. Additionally, the framework was derived under the assumption that the computational power  
 631 and communication capabilities among clients would not significantly differ. The performance and  
 632 robustness of the framework in scenarios where there are substantial variations in computational  
 633 power among clients remain areas for future research.

## 634 F Experiments

### 635 F.1 Efficiency Simulation

636 In this section, we evaluate the efficiency of frameworks using the Average Waiting Time (AWT)  
 637 metric.

638 We utilize a mixed normal distribution to model the computational capacity of clients and a normal  
 639 distribution to model the communication conditions between clients and edge servers. Next, we apply  
 640 the *Balanced Cluster Algorithm* [45] to cluster clients into  $N$  groups. Subsequently, we calculate the  
 641 AWT for various values of  $N$ .

642 We plot the simulation results in Figure 4. The scheduler scheme in HiFlash [46] is trained using  
 643 reinforcement learning, which makes its efficiency difficult to simulate. It is evident that the AWT  
 644 value of CSS-HFL decreases with an increase in the number of edge servers. Moreover, the AWT  
 645 decreases rapidly when there are fewer edge servers, exhibiting a trend similar to the elbow pattern

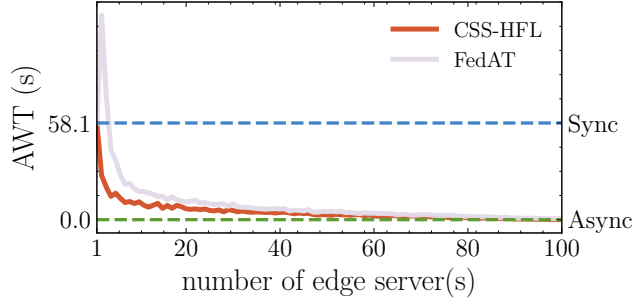


Figure 4: The AWT with different number of edge servers under various frameworks

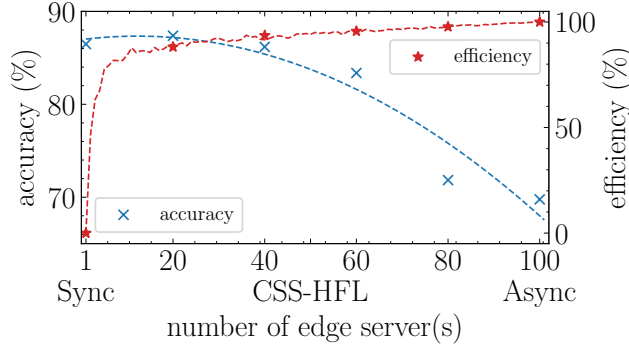


Figure 5: The accuracy and efficiency with different number of edge servers under 30% attacks on FMNIST.

646 often observed in K-means clustering. Another finding, which is in line with our theoretical analysis,  
 647 is that when the number  $N$  of edge servers equals 1, the AWT of CSS-HFL is equivalent to that of  
 648 Synchronous FL. Conversely, when  $N$  equals the number  $n$  of clients, the AWT of CSS-HFL aligns  
 649 with that of Asynchronous FL. Theoretically, AWT follows a strictly decreasing trend. However, our  
 650 figure exhibits small fluctuations. This can be attributed to the fact that clustering task is a NP-hard  
 651 problem. Additionally, we limit the maximum number of iterations in the Balanced Cluster Algorithm  
 652 to 10, which may prevent us from achieving the optimal solution for clustering in each case.

653 To further explore the trade-off relationships of CSS-HFL, we use (14) to define the efficiency  $\text{eff}$ ,  
 654 which converts AWT to an efficiency metric. In this context, the efficiency of Asynchronous and  
 655 Synchronous FL are respectively normalized to 100% and 0%.

$$\text{eff} = \left(1 - \frac{\text{AWT}}{\text{AWT}_{\text{Sync}}}\right) \cdot 100\%. \quad (14)$$

656 We investigated the accuracy and efficiency with various number of edge servers under 30% attacks  
 657 on FMNIST. As shown in Figure 5, a trade-off pattern emerges between accuracy and efficiency. As  
 658 the number of edge servers increases, the accuracy declines while the efficiency improves. Notably,  
 659 it is possible to achieve both robustness and high efficiency by selecting a certain number of edge  
 660 servers (e.g., 20 edge servers in the 100-client scenario).

661 In conclusion, our findings suggest that the inclusion of several edge servers can significantly decrease  
 662 the AWT (*i.e.* enhance efficiency) and maintain robustness against attacks under our CSS-HFL  
 663 framework.

## 664 F.2 Datasets

665 **MNIST:** The MNIST dataset is a well-known collection of handwritten digits widely used in the  
 666 field of machine learning. It consists of 60,000 training examples and 10,000 testing examples. Each  
 667 sample is a 28x28 image of a digit, ranging from 0 to 9. MNIST is a standard benchmark dataset

668 in the machine learning community and is widely employed to assess the performance of various  
 669 algorithms.

670 **Fashion-MNIST:** Fashion-MNIST is a dataset similar in structure to MNIST but comprises images of  
 671 fashion items instead of handwritten digits. The Fashion-MNIST dataset consists of 60,000 training  
 672 samples and 10,000 testing samples, which is consistent with MNIST. Fashion-MNIST stands as  
 673 a benchmark dataset for image classification endeavors, specifically in the domain of fashion and  
 674 clothing recognition. Each image within the dataset is a grayscale 28x28 pixel representation of  
 675 a fashion item, categorized into one of 10 distinct classes, such as shirts, trousers, dresses, and  
 676 shoes. Like MNIST, Fashion-MNIST has become broadly utilized in the machine learning scope for  
 677 evaluating the models.

678 **CIFAR-10:** The CIFAR-10 dataset stands as a widely recognized benchmark in the domain of  
 679 computer vision. It comprises 60,000 color images, each measuring 32x32 pixels, and is categorized  
 680 into 10 distinct classes. Like MNIST and Fashion-MNIST, CIFAR-10 serves as a standard evaluation  
 681 tool for image classification algorithms, facilitating advancements in the field of deep learning.

### 682 F.3 Experimental Results

683 In this section, we present the hyperparameters settings in Table 7 and the overview of experiment  
 684 results in Table 8. Our *FusCred* exhibited superior robustness across various poisoning attack  
 685 scenarios.

Table 7: The Hyperparameters Settings

Parameters	Description	Value	
$n$	Number of clients	100	
$N$	Number of edges	10	
$lr$	Learning rate	0.01	
$B$	Batch size	MNIST	64
		FMNIST	
		CIFAR-10	32
$E$	Number of client training epochs	2	
$R$	Number of cloud training epochs	50	
$ar$	The attack ratio	0%, 20%, 30%, 40%	
$rs$	The random seed	2023, 2024, 3047	
$\alpha$	Parameter of Dirichlet	0.2, 0.5, 0.8	



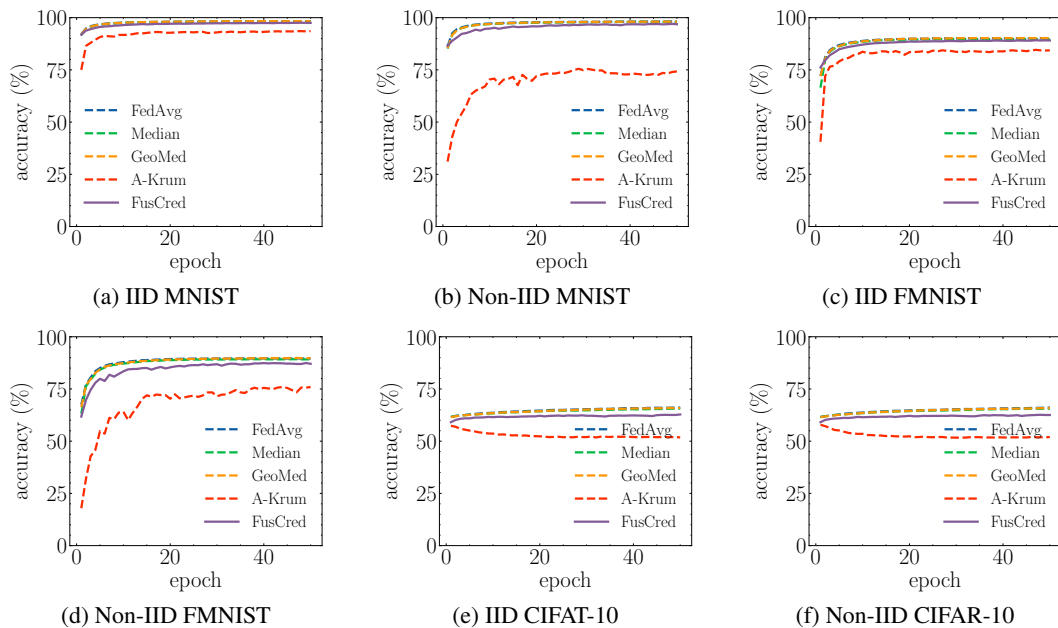


Figure 6: Accuracy without malicious attacks on IID and Non-IID datasets. A-Krum is significantly lower than other methods.

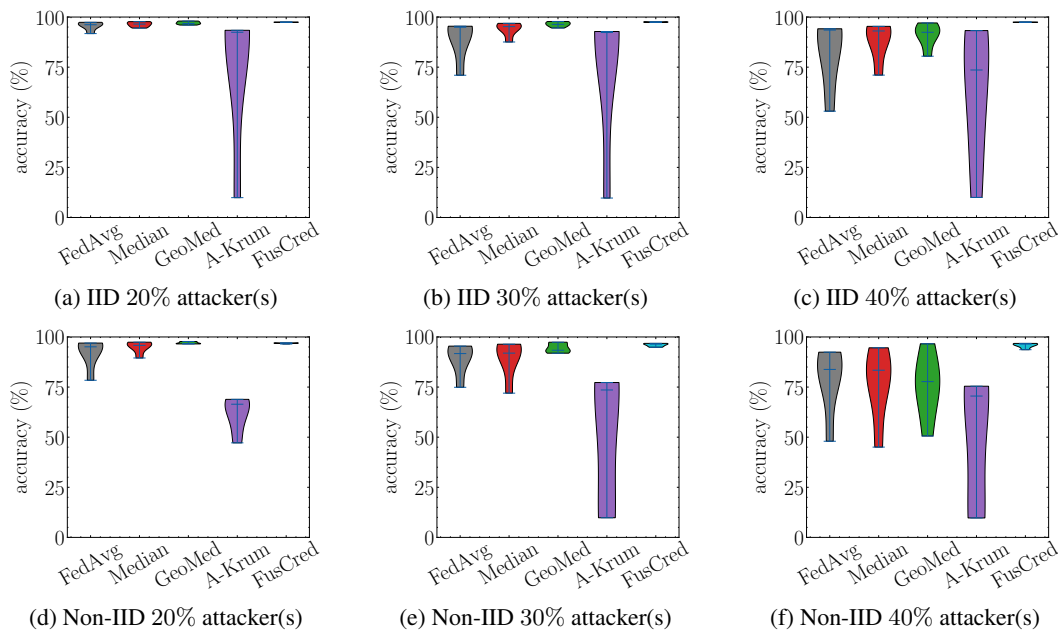


Figure 7: The violin plot of accuracy of various aggregation methods under 20%, 30%, 40% attacker(s) on IID MNIST and Non-IID MNIST.

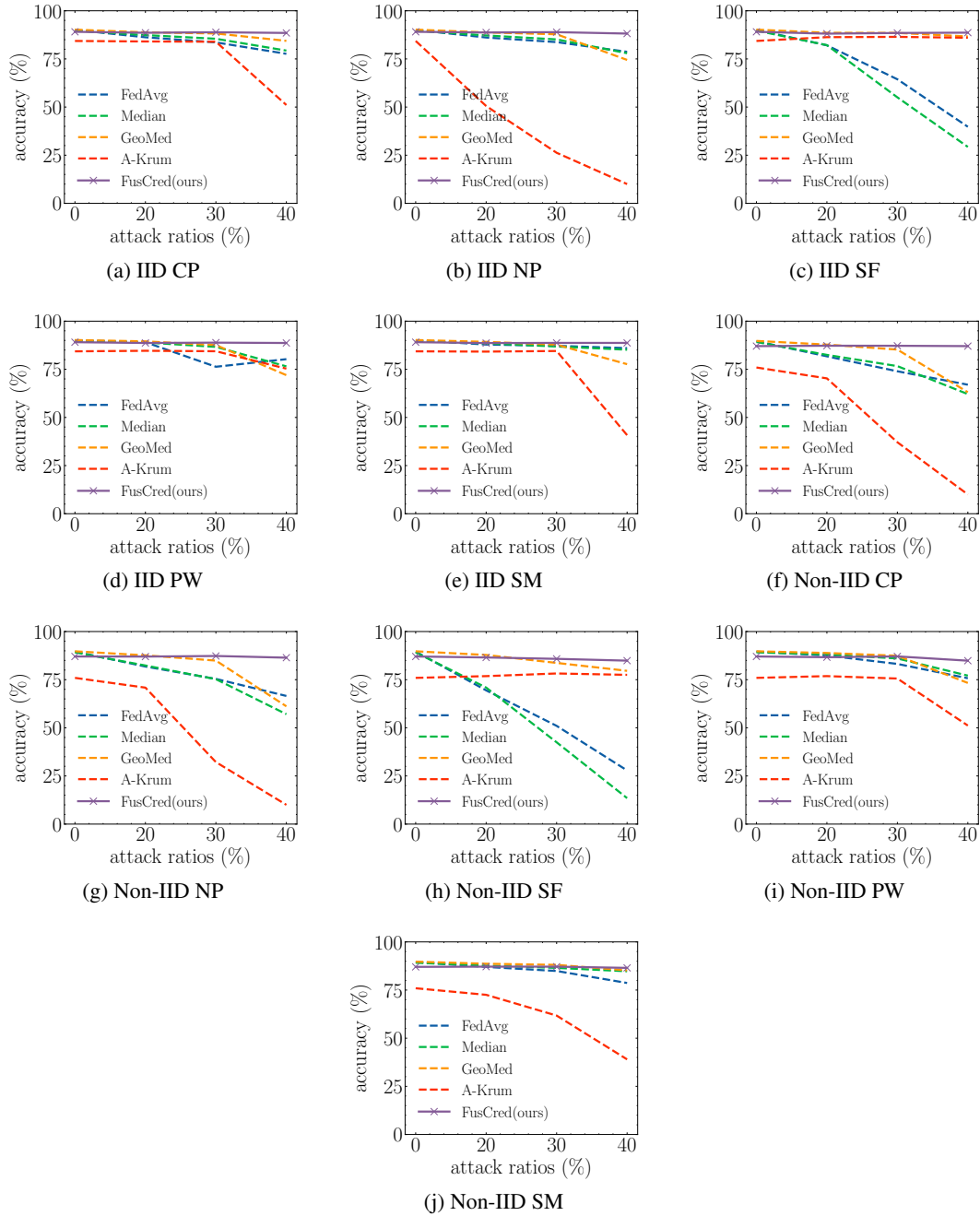


Figure 8: Impact of different attack types on accuracy for IID and Non-IID FMNIST. *FusCred* demonstrates better tolerance against various attack types than other methods.

Table 8: Experiment results overview.

Dataset	Attack ratio	Attack type	FedAvg	Median	GeoMed	A-Krum	FusCred	
IID MNIST	0%	-	98.35±0.04	98.15±0.03	98.31±0.02	93.55±0.24	97.49±0.13	
	20%	Model poison	CP	96.16±0.15	96.48±0.34	96.77±0.28	92.40±0.49	97.33±0.06
			NP	95.76±0.04	95.56±0.39	96.00±0.58	9.90±0.26	97.43±0.01
			SF	91.75±4.42	94.50±1.13	96.76±0.66	92.87±0.75	97.36±0.11
		Data poison	PW	97.17±0.44	97.71±0.06	97.87±0.04	93.39±0.60	97.50±0.01
			SM	97.36±0.13	97.51±0.07	98.00±0.04	90.43±2.20	97.45±0.04
			CP	95.41±0.12	95.73±0.22	95.93±0.22	92.06±0.52	97.42±0.15
	30%	Model poison	NP	95.18±0.12	94.40±0.02	94.55±0.27	9.70±0.22	97.47±0.07
			SF	89.23±3.95	87.55±2.18	96.33±0.28	92.78±0.11	97.28±0.09
			PW	70.96±18.57	95.21±2.11	97.63±0.10	92.77±0.55	97.56±0.03
		Data poison	SM	94.91±1.68	96.85±0.10	97.78±0.06	92.28±0.49	97.61±0.09
			CP	93.86±0.38	94.34±0.43	90.66±1.00	45.87±12.89	97.28±0.12
			NP	94.14±0.47	93.05±0.12	92.39±0.07	9.99±0.18	97.41±0.09
	40%	Model poison	SF	53.03±8.23	71.03±8.01	94.04±0.65	92.87±0.49	97.28±0.08
			PW	76.22±7.13	86.31±6.47	80.46±2.59	93.25±0.11	97.49±0.15
			SM	93.50±0.94	95.35±0.49	97.05±0.14	73.58±26.61	97.55±0.14
Data poison		CP	94.78±0.25	95.81±0.07	96.69±0.29	58.68±18.71	96.92±0.07	
		NP	95.08±0.19	95.57±0.26	96.52±0.19	47.17±22.65	96.95±0.03	
		SF	78.39±10.19	89.51±1.35	96.68±0.27	67.73±3.79	96.49±0.43	
20%	Data poison	PW	95.48±0.76	97.38±0.13	97.63±0.05	68.87±5.36	97.01±0.12	
		SM	96.95±0.29	97.11±0.06	97.74±0.03	66.44±6.18	96.92±0.03	
		CP	89.55±0.24	91.93±1.45	91.90±2.60	9.80±0.00	96.73±0.08	
	Model poison	NP	91.70±0.68	91.82±1.98	93.37±2.20	10.16±0.26	96.80±0.07	
		SF	74.89±8.93	71.96±8.50	93.07±1.04	75.72±4.21	96.41±0.10	
		PW	95.43±1.12	95.13±1.46	97.32±0.17	73.55±4.21	94.86±2.16	
30%	Data poison	SM	94.69±1.07	96.38±0.02	97.37±0.09	77.27±1.77	95.51±0.61	
		CP	83.81±0.37	83.60±0.84	69.97±11.71	9.80±0.00	96.59±0.05	
		NP	87.84±0.89	83.40±1.05	50.59±12.24	9.70±0.31	96.78±0.20	
	Model poison	SF	47.97±1.81	45.05±12.75	77.73±9.88	70.52±11.04	96.44±0.15	
		PW	80.75±6.89	75.03±8.58	81.33±6.51	75.42±2.15	93.65±2.07	
		SM	92.41±1.21	94.59±0.44	96.60±0.06	73.49±4.64	95.60±0.51	
40%	Data poison	CP	86.24±0.05	87.38±0.50	88.79±0.34	84.06±0.46	88.59±0.14	
		NP	86.12±0.20	87.20±0.47	88.76±0.34	50.55±12.01	88.73±0.21	
		SF	82.07±0.85	82.16±2.21	88.59±0.50	86.30±0.57	88.13±0.35	
	Model poison	PW	89.06±0.18	88.92±0.08	89.51±0.07	84.61±0.17	88.67±0.09	
		SM	87.85±0.17	88.46±0.28	89.26±0.07	84.19±0.24	88.62±0.09	
		CP	83.59±0.49	85.46±0.62	88.21±0.07	83.93±0.81	88.85±0.11	
30%	Model poison	NP	83.73±0.65	85.13±0.61	87.79±0.14	26.26±11.51	88.91±0.10	
		SF	64.42±9.73	55.05±18.14	88.51±0.28	86.46±0.09	88.47±0.22	
		PW	76.32±14.11	86.65±1.11	87.57±0.92	84.38±0.26	88.85±0.08	
	Data poison	SM	87.19±0.53	86.84±0.51	87.78±0.37	84.51±0.56	88.71±0.18	
		CP	77.65±1.36	79.32±1.69	84.38±1.95	51.07±12.81	88.49±0.25	
		NP	78.69±1.08	77.96±1.67	74.48±10.05	10.00±0.00	88.16±0.46	
40%	Model poison	SF	39.76±3.38	29.31±3.70	86.66±0.51	85.99±0.34	88.64±0.24	
		PW	80.26±3.05	76.52±2.85	71.96±6.26	75.34±12.93	88.67±0.32	
		SM	86.00±0.40	85.17±0.23	77.66±11.42	40.80±23.10	88.70±0.40	
	Data poison	CP	90.07±0.07	89.70±0.08	90.20±0.06	84.33±0.46	89.08±0.17	
		NP	86.24±0.05	87.38±0.50	88.79±0.34	84.06±0.46	88.59±0.14	
		SF	82.07±0.85	82.16±2.21	88.59±0.50	86.30±0.57	88.13±0.35	
20%	Data poison	PW	89.06±0.18	88.92±0.08	89.51±0.07	84.61±0.17	88.67±0.09	
		SM	87.85±0.17	88.46±0.28	89.26±0.07	84.19±0.24	88.62±0.09	
		CP	83.59±0.49	85.46±0.62	88.21±0.07	83.93±0.81	88.85±0.11	
	Model poison	NP	83.73±0.65	85.13±0.61	87.79±0.14	26.26±11.51	88.91±0.10	
		SF	64.42±9.73	55.05±18.14	88.51±0.28	86.46±0.09	88.47±0.22	
		PW	76.32±14.11	86.65±1.11	87.57±0.92	84.38±0.26	88.85±0.08	
40%	Data poison	SM	87.19±0.53	86.84±0.51	87.78±0.37	84.51±0.56	88.71±0.18	
		CP	77.65±1.36	79.32±1.69	84.38±1.95	51.07±12.81	88.49±0.25	
		NP	78.69±1.08	77.96±1.67	74.48±10.05	10.00±0.00	88.16±0.46	
	Model poison	SF	39.76±3.38	29.31±3.70	86.66±0.51	85.99±0.34	88.64±0.24	
		PW	80.26±3.05	76.52±2.85	71.96±6.26	75.34±12.93	88.67±0.32	
		SM	86.00±0.40	85.17±0.23	77.66±11.42	40.80±23.10	88.70±0.40	

Gold, silver, and bronze respectively denote the top three winners.

Table 8 continued from previous page

Dataset	Attack ratio	Attack type	FedAvg	Median	GeoMed	A-Krum	FusCred	
Non-IID FMNIST	0%	-	89.57±0.09	89.18±0.03	89.75±0.10	75.94±0.06	87.04±0.48	
	20%	Model poison	CP	81.74±1.03	82.49±1.42	87.83±0.32	70.29±9.57	87.24±0.46
			NP	81.82±1.32	82.36±2.12	87.64±0.35	70.86±10.19	87.00±0.79
			SF	69.45±5.15	70.89±5.39	87.83±0.41	76.84±2.85	86.54±1.06
	Data poison	PW	87.54±0.44	88.12±0.41	88.87±0.15	76.85±0.79	86.67±0.99	
		SM	87.00±0.30	87.68±0.15	88.67±0.25	72.53±1.97	87.15±0.46	
		CP	73.97±2.23	76.65±2.71	85.24±1.15	37.11±9.12	87.21±0.20	
	30%	Model poison	NP	75.40±1.62	75.26±3.14	84.91±1.20	32.20±0.40	87.31±0.24
			SF	51.01±9.66	42.47±16.05	83.70±3.54	78.22±0.57	85.82±0.84
			PW	83.18±2.39	86.02±1.92	87.42±1.29	75.60±2.36	87.08±0.36
	Data poison	SM	84.91±0.74	86.51±0.26	88.08±0.07	61.74±14.38	87.15±0.60	
		CP	67.00±1.56	62.12±5.21	63.11±10.94	10.00±0.00	87.04±0.31	
		NP	66.55±1.76	57.08±5.55	61.17±9.92	10.00±0.00	86.43±0.27	
	40%	Model poison	SF	27.87±15.25	13.48±4.92	79.59±2.53	77.51±0.33	84.88±0.60
			PW	75.75±0.74	77.10±4.16	73.31±8.46	51.17±20.62	84.86±2.04
			SM	78.68±1.29	84.75±0.15	85.24±0.40	39.00±11.60	86.51±0.46
IID CIFAR-10	0%	-	66.07±0.03	65.63±0.02	65.98±0.10	51.79±0.35	62.79±0.20	
	20%	Model poison	CP	42.68±0.99	42.39±0.66	54.79±0.33	51.68±0.11	62.13±0.58
			NP	41.63±0.78	41.28±0.70	53.73±0.93	51.88±0.24	62.08±0.47
			SF	32.65±0.89	27.35±0.36	58.43±0.90	51.85±0.10	62.80±0.20
	Data poison	PW	61.21±0.32	61.52±0.32	63.80±0.17	50.66±1.68	62.35±0.32	
		SM	62.41±0.29	63.32±0.08	64.03±0.13	52.33±0.46	62.47±0.27	
		CP	36.78±1.63	35.45±3.01	34.36±3.80	51.73±0.23	62.07±0.52	
	30%	Model poison	NP	35.99±1.63	34.07±2.19	33.38±3.52	51.72±0.22	61.87±0.52
			SF	28.61±2.77	19.98±1.19	49.39±2.19	51.78±0.27	62.71±0.24
			PW	56.00±0.45	57.16±0.89	60.46±0.92	50.02±1.34	62.08±0.05
	Data poison	SM	59.40±0.33	61.13±0.59	61.93±0.42	50.32±1.12	62.01±0.23	
		CP	29.96±2.27	28.11±1.35	15.09±0.42	42.74±6.72	61.46±0.39	
		NP	30.00±2.05	25.87±0.51	15.06±5.00	42.82±6.58	61.37±0.20	
	40%	Model poison	SF	21.83±3.19	12.69±2.02	33.82±3.61	52.14±0.34	62.34±0.07
			PW	48.05±1.63	48.89±2.62	52.58±3.31	52.02±0.31	61.59±0.20
			SM	54.61±0.62	57.46±1.50	56.67±2.58	49.98±1.35	61.31±0.32
Non-IID CIFAR-10	0%	-	66.08±0.07	65.58±0.05	65.87±0.09	51.97±0.01	62.53±0.12	
	20%	Model poison	CP	43.39±0.94	44.62±2.77	57.63±3.15	51.56±0.10	62.67±0.22
			NP	42.21±1.34	43.55±3.38	57.19±3.80	51.90±0.21	62.78±0.17
			SF	33.37±3.37	30.50±3.93	59.81±3.34	51.85±0.23	62.90±0.44
	Data poison	PW	60.87±0.54	61.36±0.48	63.80±0.20	51.93±0.20	62.60±0.26	
		SM	62.03±0.60	62.99±0.26	63.98±0.12	50.74±1.77	62.70±0.15	
		CP	35.82±1.12	35.94±1.77	35.37±5.61	51.96±0.02	62.42±0.20	
	30%	Model poison	NP	35.01±0.84	33.92±1.50	37.25±3.49	51.67±0.15	62.42±0.20
			SF	26.99±2.38	18.92±3.61	48.74±6.03	52.62±0.49	62.71±0.11
			PW	56.10±0.66	57.65±1.00	60.71±0.52	51.11±1.48	62.41±0.13
	Data poison	SM	59.13±0.17	61.43±0.29	61.89±0.25	50.51±2.38	61.93±0.61	
		CP	29.42±0.37	29.52±1.56	15.26±1.53	42.71±7.14	61.42±0.39	
		NP	29.36±0.50	26.43±1.21	18.23±3.54	43.69±6.29	61.51±0.36	
	40%	Model poison	SF	21.58±3.02	12.17±1.23	37.00±1.44	51.76±0.78	62.15±0.27
			PW	46.59±2.04	48.53±0.96	53.48±0.68	49.72±2.82	61.82±0.04
			SM	54.08±1.20	57.76±1.25	57.68±0.93	49.52±3.13	61.71±0.25

Gold, silver, and bronze respectively denote the top three winners.

686 **NeurIPS Paper Checklist**

687 **1. Claims**

688 Question: Do the main claims made in the abstract and introduction accurately reflect the  
689 paper’s contributions and scope?

690 Answer: [Yes]

691 Justification: Our work introduces a novel FL framework, Clustered Semi-synchronous  
692 Hierarchical Federated Learning (CSS-HFL), to address the limitations of existing FL  
693 frameworks, namely the stragglers effect, network congestion, and robustness against  
694 poisoning attacks. We also propose a robust algorithm, Fusion Credibility (FusCred), to  
695 enhance the robustness of the CSS-HFL framework.

696 **2. Limitations**

697 Question: Does the paper discuss the limitations of the work performed by the authors?

698 Answer: [Yes]

699 Justification: Please refer to Appendix E for the limitations of existing FL frameworks.

700 **3. Theory Assumptions and Proofs**

701 Question: For each theoretical result, does the paper provide the full set of assumptions and  
702 a complete (and correct) proof?

703 Answer: [Yes]

704 Justification: Please refer to Appendix D.

705 **4. Experimental Result Reproducibility**

706 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
707 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
708 of the paper (regardless of whether the code and data are provided or not)?

709 Answer: [Yes]

710 Justification: Yes, we provide detailed information on the experimental settings necessary to  
711 reproduce the main experimental results in Section 4.1.

712 **5. Open access to data and code**

713 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
714 tions to faithfully reproduce the main experimental results, as described in supplemental  
715 material?

716 Answer: [Yes]

717 Justification: We provide open access to the data and code, along with detailed instructions  
718 to reproduce the main experimental results in the supplemental material.

719 **6. Experimental Setting/Details**

720 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
721 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
722 results?

723 Answer: [Yes]

724 Justification: Please refer to Section 4.1 for the detailed experimental settings.

725 **7. Experiment Statistical Significance**

726 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
727 information about the statistical significance of the experiments?

728 Answer: [Yes]

729 Justification: We report error bars in the experimental results to show the statistical signifi-  
730 cance of the experiments.

731 **8. Experiments Compute Resources**

732 Question: For each experiment, does the paper provide sufficient information on the com-  
733 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
734 the experiments?

735 Answer: [Yes]  
736 Justification: Section 4.1 provides detailed information on the computer resources used in  
737 the experiments.

738 **9. Code Of Ethics**

739 Question: Does the research conducted in the paper conform, in every respect, with the  
740 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

741 Answer: [Yes]  
742 Justification: Our research conforms to the NeurIPS Code of Ethics.

743 **10. Broader Impacts**

744 Question: Does the paper discuss both potential positive societal impacts and negative  
745 societal impacts of the work performed?

746 Answer: [Yes]  
747 Justification: Our work addresses the limitations of existing FL frameworks, namely the  
748 stragglers effect, network congestion, and robustness against poisoning attacks, which are  
749 crucial for the real-world FL system implementation. This work can potentially provide  
750 positive societal impacts by improving the efficiency and robustness of FL systems.

751 **11. Safeguards**

752 Question: Does the paper describe safeguards that have been put in place for responsible  
753 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
754 image generators, or scraped datasets)?

755 Answer: [NA]  
756 Justification: Our work does not involve data or models that have a high risk for misuse.

757 **12. Licenses for existing assets**

758 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
759 the paper, properly credited and are the license and terms of use explicitly mentioned and  
760 properly respected?

761 Answer: [Yes]  
762 Justification: We include the license and terms of use for all assets used in the paper in the  
763 supplemental material.

764 **13. New Assets**

765 Question: Are new assets introduced in the paper well documented and is the documentation  
766 provided alongside the assets?

767 Answer: [NA]  
768 Justification: Our work does not introduce new assets.

769 **14. Crowdsourcing and Research with Human Subjects**

770 Question: For crowdsourcing experiments and research with human subjects, does the paper  
771 include the full text of instructions given to participants and screenshots, if applicable, as  
772 well as details about compensation (if any)?

773 Answer: [NA]  
774 Justification: Our work does not involve crowdsourcing experiments or research with human  
775 subjects.

776 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
777 Subjects**

778 Question: Does the paper describe potential risks incurred by study participants, whether  
779 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
780 approvals (or an equivalent approval/review based on the requirements of your country or  
781 institution) were obtained?

782 Answer: [NA]  
783 Justification: Our work does not involve research with human subjects.