ASAP: ADAPTIVE SLIDING AGNOSTIC POISONING ATTACK ON FEDERATED LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

The primary risk in the federated learning (FL) framework arises from the potential for manipulating local training data and updates, known as a poisoning attack. Among various attack strategies, agnostic attacks have emerged as a significant category that attempts to operate without explicit knowledge of the server's aggregation rules (AGRs). However, existing AGR-agnostic attacks still suffer from a critical dependency: they rely heavily on staying inside the natural per-coordinate variance of honest client updates. These attacks typically operate by analyzing benign clients' gradient patterns, statistical properties, and behavioral characteristics to strategically position their malicious updates. Therefore, to overcome these fundamental limitations of current AGR-agnostic attacks, this work presents the Adaptive Sliding Agnostic Poisoning Attack (ASAP) on FL, which can adaptively, robustly and precisely manipulate the degree of poisoning without the knowledge of AGRs algorithm of the server.

Instead of relying on benign client patterns, ASAP incorporates Adaptive Sliding Model Control (ASMC) theory — a sophisticated robust nonlinear control framework that enables adaptive attack. We implement our attack through comprehensive experiments on state-of-the-art (SOTA) Byzantine-robust federated learning methods using real-world datasets. These evaluations reveal that ASAP significantly outperforms all existing agnostic attacks while maintaining complete independence from benign client information, representing a fundamental advancement in FL attack strategies.

1 Introduction

The distributed diagram of Federated Learning (FL) ensures training models among clients devices without sharing local data but only sending model updates to a central server (Li et al., 2021). The central server initially sends the global model to selected clients, and each client then trains the model locally using its own data. The locally updated models are then transmitted to the central server, which applies a specified aggregation rule (AGR) to compute the next global model.

However, distributed systems are susceptible to poisoning attacks including both data poisoning attacks (Fang et al., 2020) and model poisoning attacks (Panda et al., 2022) due to its natural mechanism. Most poisoning attacks are designed relying on knowledge of the server's aggregation rules, which is typically difficult to obtain in practical scenarios. Therefore, the development of AGRagnostic attacks enables attack deployment without aggregation rule awareness or specification. Current AGR-agnostic attack methods, such as LIE (Baruch et al., 2019), depend on estimating statistical properties of benign client updates, particularly coordinate-wise mean μ and standard deviation σ , to generate small noises in malicious gradient updates to prevent the optimal convergence. Furthermore, Min-Max and Min-Sum (Shejwalkar & Houmansadr, 2021) attacks constrain the malicious update to lie inside the benign cluster, using a max-distance or sum-of-distances bound, while pushing in an adversarial direction. The fundamental constraint of those AGR-agnostic attacks is their need to remain indistinguishable from benign updates. Moreover, these methods aim to maximize deviation of the global model from optimal convergence, resulting in dynamics that converge to biased equilibrium. Consequently, the estimation of benign clients updates or their statistical properties remains a mandatory requirement for local malicious devices, even though aggregation rule knowledge is no longer required.

To overcome these limitations, we propose Adaptive Sliding Poisoning Attack (ASAP) on FL, a novel FL attack framework that operates without prior knowledge of server aggregation rules or benign client information. The method leverages a combination of Adaptive Sliding Mode Control (ASMC) theory and Fourier series approximation (Young et al., 1999; Ge et al., 1997; Huang & Kuo, 2001). Moreover, ASAP provides precise attack control through adjustable convergence rates and flexible attack objectives capabilities.

To achieve both AGR-agnosticism and precise attack control, we consider the entire FL process as a dynamical system and introduce an adaptive law to estimate the unknown information from malicious clients, along with a control law that guides the global model towards a specified poisoned reference. In particular, we employ Adaptive Sliding Mode Control (ASMC)—an adaptive robust control framework designed for nonlinear systems with uncertain dynamics—which exhibits strong resilience to parameter variations and external disturbances. This eliminates the reliance on explicit knowledge of the AGRs or benign client updates. By employing a Fourier series approximation, the unknown AGR behaviors and benign update patterns are treated as system uncertainties and approximated using a finite number of orthonormal basis functions. ASMC then ensures that the system state converges to and remains on a predefined sliding manifold, thereby enforcing the alignment of the global model with the attack objective. Rather than mimicking benign updates or exploiting AGR structures, the proposed ASAP attack observes the uncertainty from local malicious clients and directly manipulates the malicious gradients to achieve the desired poisoning effect without any prior access to AGR algorithms or benign statistics. Furthermore, ASAP provides fine-grained control over the attack convergence rate, enabling persistent, adaptive, and target-driven manipulation throughout the poisoning process.

Our key contribution can be concluded as below:

- We introduce ASAP, a novel adaptive AGR-agnostic controllable attack that dynamically achieves attack objectives without requiring prior knowledge of aggregation mechanisms.
- ASAP operates without knowledge of server aggregation rules or benign client statistics by leveraging Adaptive Sliding Mode Control (ASMC), distinguishing it from existing agnostic attacks that rely on coordinate-wise estimations and distance constraints.
- We provide theoretical analysis proving that ASAP achieves precise control of attack objectives and converges to predefined targets within finite time at controlled speeds, regardless of the underlying AGR algorithm. Extensive experiments on benchmark datasets against multiple robust AGRs demonstrate consistent superiority over current state-of-the-art (SOTA) methods.

2 Related Work

2.1 Federated Learning

McMahan et al. (2017) firstly demonstrated the algorithm of federated learning (FL). introduced the federated learning (FL) paradigm. In a typical FL system, a central server coordinates N clients, where client i ($i \in [1, N]$) holds its local private dataset D_i drawn from an underlying distribution D, and the datasets can be independently and identically distributed (IID) or statistically heterogeneous (Non-IID). Let $g_t \in \mathbb{R}^r$ denote the global model at iteration t, and $g_{\{t,i\}} \in \mathbb{R}^r$ he corresponding local model on each client i. The objective of FL is to address the optimization problem:

$$g^* = \operatorname*{arg\,min}_{g_{\{t,i\}}} \frac{1}{N} \sum_{i \in [1,N]} L(g_{\{t,i\}}, D_i), \tag{1}$$

where L denotes the loss function, $g^* \in \mathbb{R}^r$ is the optimal global model, and r represents dimensionality of the parameter space encompassing all network weights and biases. At iteration t, the central server broadcasts the current global parameters g_t to a participating subset of clients. Each participating client initializes its local model over g_t and trains on its private dataset D_i , producing a local model $g_{\{t,i\}}$. The client then returns the update $\nabla_{\{t,i\}} = g_{\{t,i\}} - g_t$ to the central server. After that, the central server aggregates the collected client updates according to a predefined aggregation rule $F_{\text{AGR}}(\cdot)$, yielding the aggregated update $\nabla_t = F_{\text{AGR}}(\{\nabla_{\{t,i\}}|i\in N\})$. A new global model is subsequently updated as $g_{t+1} = g_t - \eta \nabla_t$, where η is the global learning rate. This iterative procedure is repeated until the convergence criterion of global model is satisfied.

2.2 Poisoning Attacks on FL

Poisoning attacks in FL can generally be divided into two main types: data poisoning attacks and model poisoning attacks. Data poisoning attacks (Jagielski et al., 2018; Muñoz-González et al., 2017) involve adversaries contaminating the training datasets on their devices, while model poisoning attacks (Bagdasaryan et al., 2020; Baruch et al., 2019; Fang et al., 2020; Mhamdi et al., 2018; Xie et al., 2020) entail the direct manipulation of local model gradients by malicious participants, who then transmit these altered gradients to the server during the learning process. A notable study by Shejwalkar & Houmansadr (2021) mention two agnostic attacks in FL named Min-Max and Min-Sum. Min-Max aims to minimize the distance between benign and malicious clients, while Min-Sum looks for the minimized sum distances between benign and malicious clients. Without knowing the knowledge of the aggregation rules of the server, the attacks can compare the distances between themselves and benign clients instead, iteratively searching for an optimal parameter to update the malicious model, thereby achieving the performance of the attack. However, this strategy of looking for the minimum distances between benign and malicious is not realistic in the real-world FL scenarios. Additionally, the lack of control over the speed of the attack and the requirement for a significant proportion of malicious clients can lead to easy detection by robust AGRs, and low attack efficiency. Furthermore, once initiated, the predetermined attack objective in traditional attacks cannot be modified, which poses a limitation on the flexibility of the attack strategy.

2.3 Existing Byzantine-robust Defenses

Current defenses against poisoning attacks in FL are categorized based on the detection and mitigation strategies servers employ to handle suspicious models. These strategies are generally grouped into three types: statistics-based, distance-based, and performance-based approaches (Shen et al., 2022). Statistics-based defenses, such as Median (Yin et al., 2021) and Trimmed Mean (Yin et al., 2021), use statistical features like mean or median to aggregate input gradients on each dimension, mitigating outlier impacts. Distance-based defenses, including methods like Krum (Blanchard et al., 2017), Mkrum (Blanchard et al., 2017), and Bulyan (Mhamdi et al., 2018), assess distances such as Euclidean distance or Cosine similarity between local updates to pinpoint statistical outliers. Performance-based defenses, exemplified by Fang (Fang et al., 2021), rely on a validation dataset to evaluate the performance of uploaded models, removing those that diverge significantly from expected outcomes. These varied approaches reflect the complexity of securing FL systems from sophisticated attacks aimed at compromising the collaborative learning process.

2.4 Adaptive Sliding Mode Control

Adaptive Sliding Mode Control (ASMC) (Huang & Kuo, 2001) is a robust control technique that combines sliding mode control's insensitivity to matched disturbances and parameter uncertainties with adaptive mechanisms that can estimate unknown system parameters or disturbance bounds in real-time. The proposed approach, Foriers approximation technique (), systematically aggregates all uncertain parameters inherent in the controller synthesis process and represents these uncertainties through finite linear combinations of orthonormal basis functions. Consider a first order nonlinear system with the dynamic model $\dot{g}_t = u_t + d_t$, where \dot{g}_t is the derivative of g_t with respect to time t, and d_t is the disturbance which is an unknown function of time. Conventionally, the sliding surface can be specified as:

$$\dot{s}_t = \dot{e}_t + \lambda e_t,\tag{2}$$

where $e_t = \tilde{g} - g_t$ is the loss function of the system state g_t and the desired state \tilde{g} at iteration t, and λ is a hyperparameter to govern the convergence rate of e_t . If the controller can ensure that $s_t = 0$, i.e., $\dot{e}_t = -ke_t$, solving this first order differential equation will result in $e_t = e_0 e^{-kt}$, which converges exponentially to 0 as t increases.

The next step establishes the design of control law u_t to ensure $\lim_{t\to T} s_t = 0$ alongside real-time estimation of time-varying uncertainties d_t to achieve desired tracking performance. The controller design follows a two-stage methodology: (i) formulation of control law u_t to reach the sliding surface, and (ii) estimation of d_t by finite-term Fourier series approximation technique(Huang & Kuo, 2001). The control law is chosen as:

$$u_t = -\lambda e_t - \hat{d}_t + \alpha \cdot \operatorname{sign}(s_t), \tag{3}$$

where α is a positive constant selected to force the trajectory of the system to reach the sliding

 $\text{mode surface, and sign } (\cdot) \text{ is defined as follows: } \operatorname{sign}(s_t) = \begin{cases} +1 & \text{if } s_t > 0, \\ 0 & \text{if } s_t = 0,. \ \hat{d}_t \text{ is the estimation} \\ -1 & \text{if } s_t < 0. \end{cases}$

of the unknown disturbance d_t , both of them can be approximated using finite-term Fourier series approximation as:

$$d_t = w_d^T \cdot z_t, \quad \hat{d} = \hat{w}_d^T \cdot z_t \tag{4}$$

where

$$w_d = [w_0, w_1, w_2, \dots, w_{2nd}]^T$$
 (5)

$$\hat{w}_d = [\hat{w}_0, \ \hat{w}_1, \ \hat{w}_2, \ \dots , \hat{w}_{2nd}]^T \tag{6}$$

$$z_d = [1, \cos\omega_1 t, \sin\omega_1 t, \cos\omega_2, t \sin\omega_2 t, \dots, \cos\omega_{nd} t, \sin\omega_{nd} t]^T$$
(7)

and the error is $\tilde{w}_d = w_d - \hat{w}_d$. By defining the energy function

$$V_t = \frac{1}{2}s_t^2 + \frac{1}{2}\tilde{w}_d^2 \tag{8}$$

the convergence of the s_t can be ensured when $\dot{V}_t \leq 0$ (i.e., $V_{t+1} < V_t$ for $V_t \neq 0$). Substituting Eq. equation 2 and Eq. equation 3 into V_t , we can get the time derivative of V_t as:

$$\dot{V}_t = s_t \cdot \dot{s}_t + \tilde{w}_d \cdot \hat{w}_d = s_t \cdot (\tilde{w}_d z_t - \alpha \operatorname{sign}(s_t))$$
(9)

Define the adaptive law \hat{w}_d as

$$\hat{w}_d = s \cdot z_t \tag{10}$$

 $\hat{w}_d = s \cdot z_t \label{eq:wd}$ then the time derivative of V_t can be expressed as

$$\dot{V}_t = -\alpha |s_t| = -\sqrt{2}\alpha V_t^{1/2}.$$
 (11)

Therefore, the convergence of the system can be guaranteed.

3 ASAP OVERVIEW

In this section, the detailed workflow of ASAP will be demonstrated. According to the example demonstrated in Sec. 2.4, the FL training process can be treated as a nonlinear system. Considering the practical implementation, we consider the attacker can only control those compromised clients for a limited number of communication rounds. In order to formulate the attack scenario, we firstly introduce the threat model of the attack.

3.1 THREAT MODEL

Adversary's Goal The goal of the adversary is to control the malicious clients updates therefore when the malicious gradients are uploaded to the central server, the accuracy of the global model can adaptively reduce to a target accuracy without the knowledge of AGRs.

Adversary's Capability We assume the adversary controls m malicious clients of total n clients, and (m/n) < 0.5. The agnostic adversary can access global parameters and directly manipulate the malicious clients gradients to the server. Moreover, we assume that the adversary does not know any knowledge of AGRs of central server or gradients of benign clients. In FL, malicious clients naturally have access to the global model.

Comparing our attacks LIE attacks (Baruch et al., 2019) estimate coordinate-wise mean and standard deviation of all client updates to generate statistically similar malicious perturbations. Min-Max and Min-Sum attacks (Shejwalkar & Houmansadr, 2021) constrain malicious updates within benign clusters using maximum distance or sum-of-distances bounds while pushing in adversarial directions. In contrast, FMPA (Zhang et al., 2023) uses predictive reference models from historical data and subsequently fine-tunes them through gradient-based optimization to achieve desired accuracy levels with precise control. However, as demonstrated in Fig. 1, our proposed attack fundamentally differs by seeking updates that are closest to the global optima rather than diverging from it, thereby maintaining consistent effectiveness across different training phases and defensive measures without requiring statistical estimation, distance constraints, or iterative fine-tuning processes.

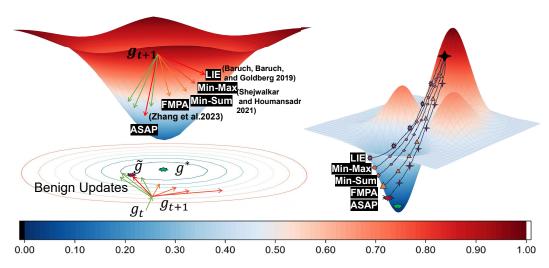


Figure 1: The comparison of existing attacks and our attack. ASAP can directly manipulate the malicious model updates and then force the global model g_{t+1} to reach the desired attack objective \tilde{g} , which is chosen as the closest point to the global optima g^* . The attack effect is illustrated via loss contours—blue area indicates low loss and red area indicates high loss.

3.2 ASAP'S ALGORITHM

We treat the overall FL global model

$$g_t = F_{AGR}\{g_{\{t,1\}}, g_{\{t,2\}}, ..., g_{\{t,m\}}, ..., g_{\{t,N\}}\}$$
(12)

as a nonlinear system, and in particular, the malicious local models are chosen as

$$\dot{g}'_{\{t,i\in m\}} = u_t. \tag{13}$$

The goal of ASAP is to design the control law u_t , and design the adaptive law \hat{w}_{Φ} by applying the function approximation technique using Fourier Series to transform the uncertainties into a finite combination of orthonormal basis functions—thus to ensure that the global model g_t will slide along the surface $s_t=0$ to achieve:

$$e_t(\tilde{g}, g_t) = -C/k \tag{14}$$

exponentially fast, where $C \in \mathbb{R}$ is a constant to adjust the convergence status of e_t , and $k \in \mathbb{R}$ $(k \neq 0)$ is a parameter to adjust the convergence speed of e_t . To achieve the adversary's goal, we design the error function as

$$e_t = g_t - \tilde{g},\tag{15}$$

To realize this new error, we design the sliding surface as

$$s_t = \int (\dot{e}_t + ke_t + \Phi_t + C)dt + C_1,$$
 (16)

where $C_1 \in \mathbb{R}$ is the initial value of the sliding surface s_t , which can be any constant, and Φ_t is the unknown disturbance.

After selecting the sliding surface s_t , the control law u_t is designed based on the FL system, the dynamic model in Eq. (13), the error function in Eq. (15) and sliding surface s_t in Eq. (16), as follows:

$$u_t = \left[\frac{dg_t}{dg'_{\{t,i\}}}\right]^{-1} \left[-ke_t + \eta \operatorname{sign}(s_t) - \hat{\Phi}_t + C\right],\tag{17}$$

where $\eta > 0$ is a positive constant selected to force the system trajectory to reach the sliding mode surface. Here, $dg_t/dg'_{\{t,i\}}$ is the derivative of g_t with respective to $g'_{\{t,i\}}$ and $\hat{\Phi}_t$ is the estimation function of Φ_t . Using the Fourier Series and approximation technique to estimate Φ_t , it can be represented as:

$$\Phi_t = \mathbf{w_{\Phi}}^T z_{\Phi} + \epsilon_{\Phi}, \quad \hat{\Phi}_t = \hat{w}_{\Phi}^T z_{\Phi}$$
 (18)

where

$$w_{\Phi} = [w_0, \ w_1, \ w_2, \ \dots , w_{2n\Phi}]^T \tag{19}$$

$$w_{\Phi} = [\hat{w}_0, \ \hat{w}_1, \ \hat{w}_2, \ \dots , \hat{w}_{2n\Phi}]^T$$
 (20)

$$z_{\Phi} = [1, \cos\omega_1 t, \sin\omega_1 t, \cos\omega_2 t, \sin\omega_2 t, \dots, \cos\omega_{n\Phi} t, \sin\omega_{n\Phi} t]^T$$
 (21)

 $w_\Phi \in \mathbb{R}$ is the weighting parameter, $\hat{w}_\Phi \in \mathbb{R}$ is the estimated weighting parameter, $z_\Phi \in \mathbb{R}^r$ is the vector of orthonormal basis function, and ϵ_Φ is the approximation error. Note that the number of n_Φ needs to be chosen bigger enough to ensure the performance of approximation of Φ_t . The adaptive law of w_Φ is defined as:

$$\dot{\hat{w}}_{\Phi} = -sz^T \tag{22}$$

The workflow of ASAP to compromise a client is demonstrated in Algorithm 1 in Appendix A.2.

3.3 ASAP CONVERGENCE ANALYSIS

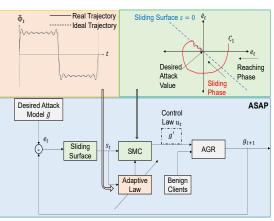


Figure 2: The block diagram of ASAP. The function of adaptive law is to automatically adjust the weight of the estimator in order to track the unknown function. In the SMC block, the system state is forced to slide along the sliding surface which means global model is forced to the desired attack objectives.

The convergence analysis is illustrated in Theorem 3.1 as below.

Theorem 3.1. Consider a FL system characterized by the dynamics in Eq. equation 13, with error function specified in Eq. equation 15 and a sliding manifold defined by Eq. equation 16. Given the control law u_t formulated in Eq. equation 26 with parameters k > 0, $\eta > 0$, $C \in \mathbb{R}$, and the derivative of the aggregation function F_{AGR} with respect to the malicious model $g'_{\{t,i\}}$ is continuous. Then the ASMC framework guarantees: (i) Fourier series approximation of the unknown uncertainty Φ_t ; (ii) finite-time convergence of the sliding surface s_t to zero with subsequent invariance; (iii) exponential convergence of the error $e_t = \tilde{w} - w_t$ to -C/k.

Note that the theoretical proof provided addresses scalar dynamics rather than vector dynamics. Since a vector is composed of multiple scalars, proving the property for each individual scalar inherently establishes the same

property for the entire vector. Thus, demonstrating the desired property at the scalar level is sufficient to confirm the corresponding property for the vector as a whole.

Due to space limit, the proof of Theorem 3.1 is delayed to Appendix A.1. Below, we highlight significant remarks on the new features of ASAP.

Remark 1: AGR-Agnostic Operation. Unlike existing AGR-agnostic attacks (LIE, Min-Max, Min-Sum) that still require statistical estimation of benign client updates, ASAP achieves complete independence from both aggregation rules and benign client information. The ASMC framework treats unknown aggregation effects as system disturbances Φ_t , which are estimated in real-time through Fourier series approximation without requiring any prior knowledge of $F_{\rm AGR}$ or benign gradient statistics.

Remark 2: Convergence Speed. The parameter k serves as a convergence rate controller, enabling precise manipulation of e_t . On the sliding surface where $s_t = \dot{s}_t = 0$, solving the differential equation $\dot{e}_t = -ke_t - C$, produces $e_t = 1/k \cdot e_0^{-kt} - C/k$. The analytical solution reveals that k determines the exponential convergence characteristics: larger values of k correspond to faster exponential convergence rates. This mathematical property enables ASAP to offer flexible convergence speed modulation capabilities.

Remark 3: Adjustable Objectives. The adversary can dynamically modify attack objectives throughout ASAP execution by appropriately selecting parameter C in e_t as evaluated in Eq. (16).

Table 1: The comparison of the accuracy of the global model between different attacks on CIFAR10, MNIST, and Tiny ImageNet against different AGRs. More experimental results against different AGRs under various attack objectives are demonstrated in Appendix A.3.4.

Dataset	AGRs	No	Test Acc. (Difference to the Targeted Acc. δ (%))						
(Model)	AGKS	Attack(%)	LIE	Min-Max	Min-Sum	FMPA	ASAP		
	Target Acc 60%								
	FedAvg	66.42	53.28 (-11.20)	32.75 (-45.42)	51.06 (-14.90)	64.33 (7.22)	61.58 (2.63)		
	Median	64.28	33.40 (-44.33)	28.08 (-53.20)	33.73 (-43.78)	63.57 (5.95)	56.53 (-5.78)		
	Trmean	66.23	46.43 (-23.78)	30.95 (-48.42)	41.19 (-31.52)	55.44 (7.60)	61.87 (3.12)		
	NB	66.73	51.95 (-84.05)	45.64 (-24.07)	55.51 (-7.48)	64.29 (7.15)	61.33 (2.22)		
CIFAR10	Bulyan	66.07	36.91 (-38.48)	25.95 (-56.75)	23.52 (-60.80)	62.55 (4.25)	61.87 (3.12)		
(AlexNet)	Mkrum	66.79	45.03 (-82.45)	52.29 (-12.85)	31.74 (-47.11)	63.26 (6.05)	60.65 (0.92)		
	Fltrust	66.59	31.53 (-47.42)	50.79 (-15.18)	52.56 (-12.4)	65.52 (9.20)	61.94 (3.23)		
	CC	66.62	63.53 (5.88)	10.53 (-82.45)	14.94 (-75.10)	67.22 (12.03)	62.13 (3.55)		
	DNC	66.55	62.92 (4.87)	63.94 (6.57)	58.26 (-2.90)	65.01 (8.35)	61.25 (2.08)		
	1	Target Acc 90%							
	FedAvg	97.98	94.12 (4.58)	91.67 (1.85)	92.84 (3.16)	95.28 (5.80)	91.04 (1.16)		
	Median	97.81	90.99 (1.10)	91.15 (1.28)	92.84 (3.16)	43.79 (-50.88)	88.22 (-1.98)		
	Trmean	97.42	91.80 (2.00)	91.30 (1.44)	92.43 (2.70)	97.26 (8.10)	90.69 (0.77)		
	NB	97.96	92.82 (3.13)	91.88 (2.09)	93.02 (3.36)	60.20 (-33.56)	90.95 (1.06)		
MNIST	Bulyan	97.97	88.92 (-1.20)	91.96 (2.18)	92.29 (2.54)	45.28 (-50.84)	89.22 (-0.87)		
(MLP)	Mkrum	97.94	92.33 (2.59)	96.14 (7.93)	95.39 (5.99)	93.41 (3.57)	92.19 (2.43)		
	Fltrust	97.96	87.89 (-2.34)	73.49 (-18.34)	93.12 (3.47)	95.01 (5.38)	92.46 (2.73)		
	CC	97.96	95.35 (5.94)	94.61 (5.12)	94.54 (5.04)	96.99 (7.86)	93.54 (3.93)		
	DNC	97.95	93.08 (3.42)	92.90 (3.22)	93.36 (3.73)	93.22 (3.58)	92.46 (2.73)		
	Target Acc 45%								
	FedAvg	57.49	51.63 (14.73)	38.37 (-23.26)	53.20 (6.40)	54.64 (21.42)	48.07 (6.82)		
	Median	53.47	22.14 (-55.72)	54.08 (8.16)	34.24 (-31.52)	42.94 (-4.58)	46.93 (4.29)		
	Trmean	54.78	51.60 (14.67)	54.59 (9.18)	39.82 (-20.36)	55.90 (24.22)	44.94 (0.13)		
Tiny	NB	58.62	52.98 (17.73)	52.95 (5.90)	53.09 (6.18)	56.12 (24.71)	45.57 (1.27)		
ImageNet	Bulyan	54.93	24.93 (-44.60)	48.01 (-3.98)	33.51 (-32.98)	5.15 (-88.56)	44.98 (0.01)		
(ResNet50)	Mkrum	54.96	27.02 (-39.96)	49.68 (-10.08)	26.39 (-47.22)	36.06 (18.44)	45.46 (1.02)		
(Resinet50)	Fltrust	54.35	33.57 (-25.40)	47.04 (4.53)	53.45 (6.90)	55.48 (23.29)	45.31 (0.69)		
	CC	54.31	29.13 (-35.27)	32.26 (-35.48)	30.99 (-31.13)	47.88 (6.40)	44.13 (-1.93)		
	DNC	55.97	68.12 (51.37)	69.66 (39.32)	54.29 (20.64)	46.98 (4.40)	44.36 (-1.42)		

When the system reaches equilibrium on the sliding manifold where both $\dot{s}_t = 0$ and $\dot{e}_t = 0$, the constraint $\dot{s}_t = \dot{e}_t + ke_t + C$ results in the equilibrium relationship $e_t = -C/k$ or $g_t = \tilde{g} + C/k$.

4 PERFORMANCE EVALUATION

4.1 Experiment Settings

Datasets and Models Our experimental evaluation of ASAP encompasses diverse architectures and benchmark datasets. We deploy AlexNet following Yang (Yang et al., 2017) for CIFAR10 experiments, utilize a fully connected (FC) neural network architecture for MNIST (Deng, 2012), and employ ResNet50 for Tiny ImageNet (Le & Yang, 2015) evaluation. The experimental framework incorporates both Independent and Identically Distributed (IID) and Non-Independent and Identically Distributed (Non-IID) data partitioning schemes. For Non-IID configurations, we leverage the Dirichlet distribution parameterized by concentration values {0.1, 0.3, 0.5, 0.7, 0.9} to systematically vary data heterogeneity levels. Smaller concentration parameters (e.g., 0.1) generate severely imbalanced client datasets with pronounced class skewness, while larger values approach uniform class distributions across participating clients. The experimental configurations are tailored to optimize performance across different architecture-dataset combinations. Comprehensive details of each dataset are provided in Appendix A.3.2.

Attack Settings The experimental setup involves a federated network of 50 clients with a 10% malicious participation rate, consistent with established benchmarks in adversarial federated learning research (Zhang et al., 2023; Shejwalkar & Houmansadr, 2021; Baruch et al., 2019). Under our threat model, adversaries gain control over compromised client devices, enabling strategic manipulation of local parameter updates to achieve precise global model accuracy targets. The attack targets are stratified across datasets: CIFAR10 targets at 60% (reference), 55%, 50%, and 10% through C parameter tuning. MNIST configurations target 90% (reference), 85%, 80%, and 10% accuracies via C adjustment. Tiny ImageNet targets of 45% (reference), 40%, 35%, and 0.5% through C modulation. The lower bounds (10% for CIFAR10/MNIST, 0.5% for Tiny ImageNet) represent random guess performance baselines. We compare our attack with existing methods including AGRagnostic approaches LIE (Baruch et al., 2019), Min-Max (Shejwalkar & Houmansadr, 2021), and Min-Sum (Shejwalkar & Houmansadr, 2021), as well as FMPA (Zhang et al., 2023) which provides

precise control capabilities but requires AGR knowledge. The details of each attack are introduced in Appendix A.3.2.

Evaluation Defenses In the experiments, various defenses are considered such as FedAvg (McMahan et al., 2017), Median (Yin et al., 2021), Trmean (Yin et al., 2021), Norm-Bounding (NB) (Sun et al., 2019) Bulyan (Mhamdi et al., 2018), Mkrum (Blanchard et al., 2017), Fltrust (Cao et al., 2022), CC (Karimireddy et al., 2021), and DNC (Shejwalkar & Houmansadr, 2021). The details of each defense are demonstrated in Appendix A.3.3.

Evaluation Metric Define I_T and I_0 as the target and achieved attack accuracies, respectively. The normalized deviation $\varsigma = ((I_T - I_0)/I_0) \times 100\%$ measures the relative distance between attack objectives and actual results. Attack method comparison employs the absolute metric $|\varsigma|$, where smaller values denote better objective fulfillment and higher attack quality.

4.2 EXPERIMENTS RESULTS

Experimental results presented in Table 1 and Figure 3 demonstrate the comparative performance of attack methods against different AGRs using CIFAR10/AlexNet, MNIST/MLP, and Tiny ImageNet/ResNet50 benchmarks. More experimental results under different scenarios are demonstrated in Appendix A.3.4. Overall, ASAP achieves the minimal $|\delta|$ values and consistently outperforms all baseline attacks.

As shown in Fig. 3, ASAP achieves robust convergence to attack objectives without triggering AGR detection, requiring fewer communication rounds than competing methods. In contrast to AGR-agnostic attacks including LIE, Min-Max and Min-Sum, which fail to achieve precise control

and demand increased communication resources, and unlike FMPA, which encounters detection by AGRs under various conditions, causing the test accuracy to converge near the optimal performance achieved without any attack presence.

Table 2: Time Complexity and Effective Communication Rounds comparisons.

Comparison	LIEN	/lin-Ma	x M	in-Suı	n FMPA	ASAP
Time (hrs)	0.8	0.9		0.9	1.0	1.1
Rounds (epochs	s) 781	778		767	34	19

The comprehensive evaluation demonstrates

ASAP's consistent performance compared to existing SOTA AGR-agnostic attack methods, coupled with fine-grained controllability for precise attack execution. The subsequent discussion examines the findings across three key dimensions.

Time Complexity The computational cost analysis, detailed in Table 2, reveals that ASAP demonstrates the highest execution time among evaluated methods, primarily due to the computational demands of its underlying mathematical framework. Nevertheless, the increased computational cost compared to competing methods remains feasible for practical deployment.

Effective Communication Rounds To maintain evaluation consistency, we utilize Effective Communication Rounds (ECR) as the standardized communication efficiency metric. Table 2 presents the average convergence performance on CIFAR10 dataset, establishing that ASAP requires the minimum number of communication rounds to achieve attack objectives compared to existing approaches.

Precise Control Table 3 in Appendix A.3.4 presents comprehensive evaluation results across multiple AGRs under diverse attack objectives. FedSA consistently exhibits the lowest $|\delta|$ scores while surpassing all comparative attacks, validating its capability for accurate objective targeting with minimal loss variance. The CIFAR10 results show average $|\delta|$ values of 2.18%, 2.61%, and 1.62% for attack objectives of 60%, 55%, and 50% respectively.

4.3 ABLATION STUDY

In this section, we conduct extensive sensitivity analysis to evaluate ASAP's robustness under varying experimental conditions, including the impact of attack speed, the impact of percentage of attackers, the impact of number of clients, the impact of clients sampling rate and the impact of

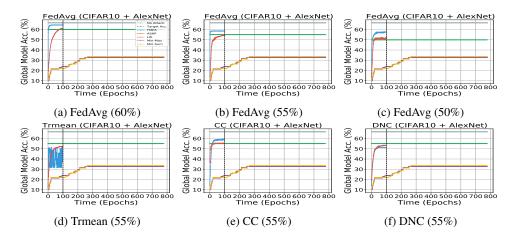


Figure 3: Comparison of each attack against various AGRs with different attack objectives on CI-FAR10 with AlexNet under different attack objectives and different attacks under the same target accuracy. Comparison figures on MNIST and Tiny ImageNet are given in Appendix A.3.4.

Non-IID degrees. The outcomes of ablation study and additional ablation studies are provided are demonstrated below and the detailed statements are illustrated in Appendix A.3.5.

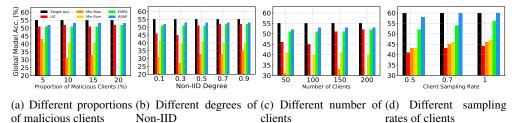


Figure 4: Ablation study results against Trmean on CIFAR10 with AlexNet to target accuracy at 55%.

5 CONCLUSION

In this paper, we introduced ASAP, a novel AGR-agnostic model poisoning attack on Federated Learning, inspired by Adaptive Sliding Mode Control theory. Unlike prior agnostic attacks that rely on heuristic distance-based strategies or require partial knowledge of benign updates, ASAP formulates the poisoning process as a controllable nonlinear system. By leveraging a Fourier series-based estimator, ASAP precisely tracks the global model trajectory and adaptively adjusts the direction and magnitude of malicious updates toward a predefined target. This enables both fine-grained control over convergence speed and resilience against diverse aggregation rules.

Our theoretical analysis guarantees convergence to the attack objective under finite time, without requiring knowledge of the server's aggregation strategy or benign client behavior. Extensive experiments on CIFAR-10, MNIST, and Tiny ImageNet across various robust AGRs which demonstrate that ASAP consistently outperforms SOTA AGR-agnostic attacks in both convergence efficiency and target alignment.

ASAP opens a new attack surface in FL by enabling precise, stealthy, and adaptive poisoning. To counteract this threat, future research should explore dynamic defense mechanisms. In particular, we propose leveraging system identification techniques to model and detect abnormal update dynamics introduced by adaptive attackers. By identifying deviations from expected system behavior, such defenses could adaptively reject suspicious updates in real time.

ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. In this study, no human subjects or animal experimentation was involved. All datasets used, including CIFAR10, MNIST and Tiny ImageNet, were sourced in compliance with relevant usage guidelines, ensuring no violation of privacy. We have taken care to avoid any biases or discriminatory outcomes in our research process. No personally identifiable information was used, and no experiments were conducted that could raise privacy or security concerns. We are committed to maintaining transparency and integrity throughout the research process.

REPRODUCIBILITY STATEMENT

We have made every effort to ensure that the results presented in this paper are reproducible. All code and datasets have been made publicly available in an anonymous repository to facilitate replication and verification. The experimental setup, including training steps, model configurations, and experimental setting details, is described in detail in the paper.

Additionally, the public benchmark datasets used in this paper, such as CIFAR10, MNIST and Tiny ImageNet, are publicly available, ensuring consistent and reproducible evaluation results.

We believe these measures will enable other researchers to reproduce our work and further advance the field.

REFERENCES

- Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How To Backdoor Federated Learning. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, pp. 2938–2948. PMLR, June 2020. URL https://proceedings.mlr.press/v108/bagdasaryan20a.html.
- Gilad Baruch, Moran Baruch, and Yoav Goldberg. A Little Is Enough: Circumventing Defenses For Distributed Learning. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/hash/ec1c59141046cd1866bbbcdfb6ae31d4-Abstract.html.
- Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine Learning with Adversaries: Byzantine Tolerant Gradient Descent. In Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/hash/f4b9ec30ad9f68f89b29639786cb62ef-Abstract.html.
- Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. FLTrust: Byzantine-robust Federated Learning via Trust Bootstrapping, April 2022. URL http://arxiv.org/abs/2012.13995. arXiv:2012.13995.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. IMAGENET: A large-scale hierarchical image database. In 2009 IEEE Conference on Computer Vision and Pattern Recognition, pp. 248–255. Ieee, 2009.
- Li Deng. The MNIST Database of Handwritten Digit Images for Machine Learning Research [Best of the Web]. *IEEE Signal Processing Magazine*, 29(6):141–142, November 2012. ISSN 1558-0792. doi: 10.1109/MSP.2012.2211477. URL https://ieeexplore.ieee.org/abstract/document/6296535. Conference Name: IEEE Signal Processing Magazine.
- Liandi Fang, Shihong Ding, Ju H. Park, and Li Ma. Adaptive Fuzzy Control for Stochastic High-Order Nonlinear Systems With Output Constraints. *IEEE Transactions on Fuzzy Systems*, 29(9): 2635–2646, September 2021. doi: 10.1109/tfuzz.2020.3005350.
- Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. Local Model Poisoning Attacks to Byzantine-Robust Federated Learning. pp. 1605-1622, 2020. ISBN 978-1-939133-17-5. URL https://www.usenix.org/conference/usenixsecurity20/presentation/fang.

- Shuzhi S. Ge, C.C. Hang, and L.C. Woon. Adaptive neural network control of robot manipulators in task space. *IEEE Transactions on Industrial Electronics*, 44(6):746–752, 1997. doi: 10.1109/41.649934.
 - An-Chyau Huang and Yeu-Shun Kuo. Sliding control of non-linear systems containing time-varying uncertainties with unknown bounds. *International Journal of Control*, 74(3):252–264, 2001.
 - Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, and Bo Li. Manipulating Machine Learning: Poisoning Attacks and Countermeasures for Regression Learning. In 2018 IEEE Symposium on Security and Privacy (SP), pp. 19–35, May 2018. doi: 10. 1109/SP.2018.00057. URL https://ieeexplore.ieee.org/abstract/document/8418594. ISSN: 2375-1207.
 - Sai Praneeth Karimireddy, Lie He, and Martin Jaggi. Learning from History for Byzantine Robust Optimization, June 2021. URL http://arxiv.org/abs/2012.10333.arXiv:2012.13995.
 - Suiyang Khoo, Lihua Xie, and Zhihong Man. Robust finite-time consensus tracking algorithm for multirobot systems. *IEEE/ASME Transactions on Mechatronics*, 14(2):219–228, April 2009. ISSN 1941-014X. doi: 10.1109/TMECH.2009.2014057.
 - A. Krizhevsky. Learning Multiple Layers of Features from Tiny Images. 2009. URL https://www.semanticscholar.org/paper/Learning-Multiple-Layers-of-Features-from-Tiny-Krizhevsky/5d90f06bb70a0a3dced62413346235c02b1aa086.
 - Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. CS 231N, 7(7):3, 2015.
 - Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. FedBN: Federated Learning on Non-IID Features via Local Batch Normalization, May 2021. URL http://arxiv.org/abs/2102.07623. ArXiv:2102.07623 arXiv: 2102.07623.
 - Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, pp. 1273–1282, April 2017. URL https://proceedings.mlr.press/v54/mcmahan17a.html.
 - El Mahdi El Mhamdi, Rachid Guerraoui, and Sébastien Rouault. The Hidden Vulnerability of Distributed Learning in Byzantium, July 2018. URL http://arxiv.org/abs/1802.07927. arXiv:1802.07927.
 - Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wongrassamee, Emil C. Lupu, and Fabio Roli. Towards Poisoning of Deep Learning Algorithms with Backgradient Optimization. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, pp. 27–38, 2017. ISBN 978-1-4503-5202-4. doi: 10.1145/3128572.3140451. URL https://dl.acm.org/doi/10.1145/3128572.3140451.
 - Ashwinee Panda, Saeed Mahloujifar, Arjun Nitin Bhagoji, Supriyo Chakraborty, and Prateek Mittal. Sparsefed: Mitigating model poisoning attacks in federated learning with sparsification. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.), *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 7587–7624. PMLR, 28–30 Mar 2022. URL https://proceedings.mlr.press/v151/panda22a.html.
 - Virat Shejwalkar and Amir Houmansadr. Manipulating the Byzantine: Optimizing Model Poisoning Attacks and Defenses for Federated Learning. In *Proceedings 2021 Network and Distributed System Security Symposium*, 2021. ISBN 978-1-891562-66-2. doi: 10.14722/ndss. 2021.24498. URL https://www.ndss-symposium.org/wp-content/uploads/ndss2021_6C-3_24498_paper.pdf.
 - Liyue Shen, Yanjun Zhang, Jingwei Wang, and Guangdong Bai. Better Together: Attaining the Triad of Byzantine-robust Federated Learning via Local Update Amplification. In *Proceedings of the 38th Annual Computer Security Applications Conference*, pp. 201–213, Austin TX

USA, December 2022. ACM. ISBN 978-1-4503-9759-9. doi: 10.1145/3564625.3564658. URL https://dl.acm.org/doi/10.1145/3564625.3564658.

Ziteng Sun, Peter Kairouz, Ananda Theertha Suresh, and H. Brendan McMahan. Can You Really Backdoor Federated Learning?, December 2019. URL http://arxiv.org/abs/1911.07963. arXiv:1911.07963.

- Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. Fall of Empires: Breaking Byzantine-tolerant SGD by Inner Product Manipulation. In *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, pp. 261–270. PMLR, August 2020. URL https://proceedings.mlr.press/v115/xie20a.html.
- Wei Yang, Shuang Li, Wanli Ouyang, Hongsheng Li, and Xiaogang Wang. Learning feature pyramids for human pose estimation. In 2017 IEEE International Conference on Computer Vision (ICCV). IEEE, October 2017. doi: 10.1109/iccv.2017.144.
- Dong Yin, Yudong Chen, Kannan Ramchandran, and Peter Bartlett. Byzantine-Robust Distributed Learning: Towards Optimal Statistical Rates, February 2021. URL http://arxiv.org/abs/1803.01498. arXiv:1803.01498.
- K.D. Young, V.I. Utkin, and U. Ozguner. A control engineer's guide to sliding mode control. *IEEE Transactions on Control Systems Technology*, 7(3):328–342, May 1999. doi: 10.1109/87.761053. URL http://ieeexplore.ieee.org/document/761053/.
- Hangtao Zhang, Zeming Yao, Leo Zhang, Shengshan Hu, Chao Chen, Alan Liew, and Zhetao Li. Denial-of-service or fine-grained control: Towards flexible model poisoning attacks on federated learning. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence (IJCAI-23)*, 2023.

A APPENDIX

A.1 PROOF OF THEOREM 1

Proof. To design the update law for \hat{w}_{Φ} , defining $\tilde{w}_{\Phi} = w_{\Phi} - \hat{w}_{\Phi}$ and a Lyapunov function (or energy function)

$$V_t = \frac{1}{2}s_t^2 + \frac{1}{2}\tilde{w}_{\Phi}^2 \tag{23}$$

and differentiating V_t with respect to time, we have

$$\dot{V}_t = s_t \dot{s}_t - \tilde{w}_{\Phi} \dot{\hat{w}}_{\Phi} \tag{24}$$

$$= s_t \left(-\frac{dg_t}{dg_{\{t,i\}}} [\dot{g}_{\{t,i\}} - \hat{\Phi}_t + ke_t + C] \right). \tag{25}$$

Using control law

$$u_t = \left[\frac{dg_t}{dg'_{\{t,i\}}}\right]^{-1} \left[ke_t + \eta \operatorname{sign}(s_t) - \hat{\Phi}_t + C\right],\tag{26}$$

we get

$$\dot{V}_t = s_t[-ke_t - \eta \operatorname{sign}(s_t) - C + \Phi_t - \hat{\Phi}_t + ke_t + C] - \tilde{w}_{\Phi} s z^T$$
(27)

$$\leq s_t[-\eta_1 \mathrm{sign}(s_t)] \tag{28}$$

$$= -\eta_1 |s_t| = -\sqrt{2}\eta_1 V_t^{1/2} \tag{29}$$

where $\eta = \eta_1 + \delta$, $\delta > 0$, $\dot{w}_{\Phi} = -sz^T$. By the finite time stability theorem proved in the study (Khoo et al., 2009), V_t will converge to zero in a finite time, and hence results in $s_t = \dot{s}_t = 0$ in a finite time.

A.2 ALGORITHM

648

649 650

651

652

653 654

655

656

657

658

659

660

661

662

663

665

666

667

668

669

670671672673

674

675

676

677

678

679

680

682

683 684

685

686

687

688 689

690

691

692

693

696

697

699

700

In this section, the algorithm of ASAP is demonstrated. We firstly initialize the value of g_t , s_t as g_0 , s_0 respectively (line 2), and malicious clients need to initialize weight of the estimator \hat{w}_{Φ} . At the t-th communication round, the client is selected by the server and receives the current global model q_t .

Algorithm 1 The workflow of ASAP to compromise a client

```
Require: Global model g_t, desired poisoning model \tilde{g}.
Ensure: malicious model update g'_t.
 1: if t=0 then
       g_t \leftarrow g_0, \tilde{g}
 2:
 3:
       Initialize \hat{w}_{\Phi}
 4: else
       for malicious client i = 1 to m do
 5:
         Update the adaptive law in Eq. (22)
 6:
         Calculate \Phi_t in Eq. (18)
 7:
         Calculate e_t of g_t and g_t^* in Eq. (15)
 8:
 9:
         Calculate s_t in Eq. (16)
10:
       end for
11:
       calculate g'_t from Eq. (26)
                                                                                           12:
       Output g'_t
13: end if
14: Update the malicious client model g'_t on FL
```

A.3 DATASETS, ATTACKS AND DEFENSES

In this section, we give details of our experiments settings. For CIFAR10 experiments with AlexNet, we establish a global learning rate of 0.02, a global batch size of 128, and conduct training over 100 global rounds, with local client updates using a batch size 10 across 5 local epochs. MNIST experiments employing MLP utilize a global learning rate of 0.01, a global batch size of 128, and 100 training rounds, while local training proceeds with a batch size 5 over 3 epochs. The Tiny ImageNet-ResNet50 configuration employs a global learning rate of 0.001, maintains a batch size 128, and executes 20 global rounds, with local updates using a batch size 10 for 3 epochs. These hyperparameter selections reflect architecture-specific optimization requirements and dataset complexity considerations.

A.3.1 DATASETS

- **CIFAR10** (Krizhevsky, 2009). It is an image database with 60,000 colour images of 32 * 32 size in 10 classes equally, and it is divided into training dataset with 50,000 images and test dataset with 10,000 images.
- MNIST (Deng, 2012). It is a dataset with 70,000 hand-written digital images in 28 * 28 size with 10 classes equally, and it is divided into training dataset with 60,000 images and test dataset with 10,000 images.
- Tiny ImageNet (Le & Yang, 2015) It is a subset of ILSVRC (ImageNet challenge) (Deng et al., 2009), which is one of the most famous benchmarks for image classification. As a subset, Tiny ImageNet only has 200 different classes. In addition, each class contains 500 training images, 50 validation images, and 50 test images totally. Moreover, the size of the images is revised to 64 * 64 pixels instead of 224 * 224 pixels in standard ImageNet.

A.3.2 ATTACKS

• LIE (Baruch et al., 2019). It inserts an appropriate amounts of noise which are large for the adversary to impact the global model while small to avoid attention by Byzantine-robust AGRs to each dimension of the average of the benign gradients.

- 702 704
- 705 706 708
- 710 711 712

709

718

- 719 720 721 722
- 723 724 725 726
- 727 728 729 730
- 731 732 733 734
- 735 737 738
- 739 740 741 742
- 743 744 745
- 746 747 748
- 749

750 751

752

753 754

755

- Min-Max (Shejwalkar & Houmansadr, 2021). They minimize the distance of malicious clients to benign clients, and then ensure the poisoned updates lie closely to the clique of benign gradients.
- Min-Sum (Shejwalkar & Houmansadr, 2021). They minimize the sum of the squared distance of malicious clients to benign clients, and then ensure the poisoned updates lie closely to the clique of benign gradients.
- FMPA (Zhang et al., 2023). It generates an estimator to predict the global model in the next iteration as a benign reference model to fine-turn the global model to the desired poisoned model by collecting the historical information.

A.3.3 DEFENSES

- FedAvg (McMahan et al., 2017). It is a basic algorithm on FL without defense. It collects all the local updates from the clients and computes the average of them as the output of aggregation.
- Median (Yin et al., 2021). It computes the median of the values from each dimension of gradients as a new global gradient.
- Trmean (Trimmed-mean) (Yin et al., 2021). It drops the specific number of maximum and minimum values from the local updates from the clients, and use the average value of the remaining updates as the aggregation output.
- Norm-bounding (Sun et al., 2019). It will scale the local update of the clients if the l_2 norm of it is bigger than the fixed threshold. Then it will average the scaled local updates as it's aggregation.
- Bulyan (Mhamdi et al., 2018). It uses Mkrum to select the updates as a selection set and then use Trmean (Yin et al., 2021) to aggregate the gradients. Trmean averages the gradients after removing the m largest and smallest values from the updates, m is usually set as the number of malicious clients.
- Mkrum (Blanchard et al., 2017). It was modified by krum (Blanchard et al., 2017) to aggregate the information provided from the clients effectively. Krum selects the single gradient which is closest to (N-m-2) neighboring gradients, where N and m are the number of all clients and malicious clients respectively. Mkrum select multi gradients using krum to obtain a selection set and then average the gradients.
- Fltrust (Cao et al., 2022). It assigns a trust score to each clients based on the updates from them to the global update direction, the lower trust score the client get, the more the direction deviates. Then Fltrust normalizes the gradients of local model updates by the trust cores, and then average the updates as a global model.
- CC (Centered Clipping) (Karimireddy et al., 2021). It clips all the gradients to the bad vector ρ to ensure the error is less than a specific value. Then it averages the normalized local updates with the weight of the trust score to generate a new global model.
- DNC (Shejwalkar & Houmansadr, 2021). Singular value decomposition (SVD) is employed for Divide-and-conquer (DnC) to extract the common features. The projection of a subsampled gradients generated from a selection of a sorted set of indices is computed, and then the gradients with highest scores of outlier vector will be removed. DnC averages the gradients after repeating this process.

EXPERIMENTAL RESULTS A.3.4

A.3.5 ABLATION STUDY

Impact of percentage of attackers The impact of malicious client proportion on FL is analyzed by incrementally increasing the adversarial ratio from 5% to 20%. As illustrated in Figure 4a, ASAP exhibits consistent performance advantages compared to competing attack strategies across all evaluated ratios.

Table 3: The comparison of the accuracy of the global model between different attacks on CIFAR10, MNIST, and Tiny ImageNet against different AGRs.

Dataset	Lagn	l No	No Test Acc. (Difference to the Targeted Acc. δ (%))						
(Model)	AGRs	Attack(%)	LIE	Min-Max	Min-Sum	FMPA	ASAP		
	Target Acc 55%								
CIFAR10	FedAvg Median Trmean NB Bulyan Mkrum	66.42 64.28 66.23 66.73 66.07 66.79	53.28 (-11.20) (-82.82) 46.43 (-23.78) (-84.05) 36.91 (-38.48) (-82.45)	32.75 (-45.42) 28.08 (-53.20) 30.95 (-48.42) 45.64 (-24.07) 25.95 (-56.75) 52.29 (-12.85)	51.06 (-14.90) 33.73 (-43.78) 41.19 (-31.52) 55.51 (-7.48) 23.52 (-60.80) 31.74 (-47.11)	58.44 (6.40) 51.05 (-6.87) 58.22 (-21.60) 58.07 (5.98) 48.71 (-11.29) 51.10 (-7.15)	56.37 (2.49) 51.47 (-6.42) 52.50 (-4.55) 56.42 (2.58) 53.71 (-2.35) 54.92 (-0.15)		
	Fltrust CC DNC	66.59 66.62 66.55	31.53 (-47.42) (-72.65) 62.92 (4.87)	50.79 (-15.18) 10.53 (-82.45) 63.94 (6.57)	52.56 (-12.4) 14.94 (-75.10) 58.26 (-2.90)	53.62 (-37.27) 58.98 (7.95) 51.43 (-6.51)	55.16 (0.29) 55.71 (1.29) 53.14 (-3.38)		
(AlexNet)	Target Acc 50%								
	FedAvg Median Trmean NB Bulyan Mkrum Fltrust CC DNC	66.42 64.28 66.23 66.73 66.07 66.59 66.62 66.55	53.28 (-11.20) (-82.82) 46.43 (-23.78) (-84.05) 36.91 (-38.48) (-82.45) 31.53 (-47.42) (-72.65) 62.92 (4.87)	32.75 (-45.42) 28.08 (-53.20) 30.95 (-48.42) 45.64 (-24.07) 25.95 (-56.75) 52.29 (-12.85) 50.79 (-15.18) 10.53 (-82.45) 63.94 (6.57)	51.06 (-14.90) 33.73 (-43.78) 41.19 (-31.52) 55.51 (-7.48) 23.52 (-60.80) 31.74 (-47.11) 52.56 (-12.4) 14.94 (-75.10) 58.26 (-2.90)	57.84 (15.40) 48.98 (-2.70) 52.78 (6.34) 55.79 (12.04) 58.94 (18.06) 62.56 (25.04) 43.86 (4.52) 43.99 (12.02) 51.36 (4.88)	50.87 (1.74) 50.99 (1.98) 50.71 (1.42) 50.26 (0.52) 49.83 (-0.34) 50.53 (1.06) 51.87 (3.74) 50.42 (0.84) 51.46 (2.92)		
	1	Target Acc 85%							
	FedAvg Median Trmean NB Bulyan Mkrum Fltrust CC DNC	97.98 97.81 97.42 97.96 97.97 97.94 97.96 97.96 97.95	94.12 (4.58) 90.99 (6.46) 91.80 (7.99) 88.92 (4.61) 92.33 (8.62) 95.19 (11.88) 87.89 (3.40) 95.35 (12.18) 93.08 (9.51)	91.67 (1.85) 91.15 (7.24) 91.30 (7.41) 91.88 (8.09) 91.96 (8.19) 95.21 (12.01) 73.49 (-13.54) 89.61 (5.42) 92.90 (9.29)	92.84 (3.16) 92.84 (9.22) 92.43 (2.70) 92.29 (2.54) 92.29 (2.54) 95.39 (12.22) 93.12 (3.47) 94.54 (5.04) 93.36 (9.84)	83.21 (-2.11) 51.74 (39.13) 95.84 (12.75) 88.35 (3.94) 98.68 (16.09) 86.71 (2.01) 93.00 (9.41) 94.86 (11.60) 99.47 (17.02)	85.70 (0.82) 88.37 (3.96) 84.45 (-0.65) 86.15 (1.35) 87.94 (3.46) 85.40 (0.47) 87.94 (3.46) 83.72 (-1.51) 86.38 (1.62)		
MNIST	Target Acc 80%								
(MLP)	FedAvg Median Trmean NB Bulyan Mkrum Filtrust CC DNC	97.98 97.81 97.42 97.96 97.97 97.94 97.96 97.96 97.95	94.12 (4.58) 90.99 (1.10) 91.80 (2.00) 88.92 (-1.20) 88.92 (-1.20) 95.19 (18.99) 87.89 (-2.34) 94.48 (18.10) 93.08 (3.42)	91.67 (1.85) 91.15 (1.28) 91.30 (1.44) 91.96 (2.18) 91.96 (2.18) 95.21 (19.01) 73.49 (-18.34) 94.66 (18.32) 92.90 (2.11)	92.84 (3.16) 92.84 (3.16) 92.43 (2.70) 92.29 (2.54) 92.29 (2.54) 95.39 (19.24) 93.12 (3.47) 94.54 (18.18) 93.36 (3.73)	92.39 (15.49) 35.52 (55.60) 97.44 (21.80) 46.36 (-42.05) 68.69 (-14.14) 25.52 (-68.10) 95.11 (18.89) 92.39 (15.49) 92.62 (15.77)	80.68 (0.85) 69.70 (-12.88) 78.75 (-1.56) 79.87 (-0.16) 77.84 (-2.70) 80.02 (0.03) 77.47 (-3.16) 76.85 (-3.94) 82.54 (3.18)		
	Target Acc 40%								
Tiny	FedAvg Median Trmean NB Bulyan Mkrum Fitrust CC DNC	56.46 52.87 56.02 57.23 56.03 54.98 55.63 53.23 53.13	51.60 (3.20) 22.14 (-55.72) 51.60 (3.20) 52.98 (17.73) 24.93 (-44.60) 27.02 (-39.96) 33.57 (-25.40) 29.13 (-35.27) 68.12 (51.37)	38.37 (-23.26) 54.08 (8.16) 54.59 (9.18) 52.95 (5.90) 48.01 (-3.98) 49.68 (-10.08) 47.04 (4.53) 32.26 (-35.48) 69.66 (39.32)	53.20 (6.40) 34.24 (-31.52) 39.82 (-20.36) 53.09 (6.18) 33.51 (-32.98) 26.39 (-47.22) 53.45 (6.90) 30.99 (-31.13) 54.29 (20.64)	38.79 (-3.03) 43.46 (8.65) 46.95 (17.38) 34.09 (-14.78) 38.9 (-2.75) 35.61 (-10.98) 35.97 (-10.08) 41.76 (4.40) 40.64 (1.60)	40.42 (1.05) 39.45 (-1.37) 40.12 (0.30) 38.74 (-3.15) 39.52 (-1.20) 40.82 (2.05) 37.86 (-5.35) 40.91 (2.27) 39.31 (-1.73)		
ImageNet	Target Acc 35%								
(ResNet50)	FedAvg Median Trmean NB Bulyan Mkrum Fltrust CC DNC	56.46 52.87 56.02 57.23 56.03 54.98 55.63 53.23 53.13	51.60 (3.20) 22.14 (-55.72) 51.60 (3.20) 52.98 (17.73) 24.93 (-44.60) 27.02 (-39.96) 33.57 (-25.40) 29.13 (-35.27) 68.12 (51.37)	38.37 (-23.26) 54.08 (8.16) 54.59 (9.18) 52.95 (5.90) 48.01 (-3.98) 47.04 (4.53) 32.26 (-35.48) 69.66 (39.32)	53.20 (6.40) 34.24 (-31.52) 39.82 (-20.36) 53.09 (6.18) 33.51 (-32.98) 26.39 (-47.22) 53.45 (6.90) 30.99 (-31.13) 54.29 (20.64)	48.06 (45.09) 33.88 (-6.86) 48.47 (36.57) 50.12 (39.37) 4.27 (-91.20) 37.30 (9.31) 49.27 (30.40) 41.70 (24.00) 48.34 (37.49)	33.73 (-0.03) 34.45 (-3.20) 34.35 (2.09) 32.73 (1.46) 37.75 (0.86) 34.51 (1.11) 35.34 (2.40) 37.16 (4.09) 35.49 (0.40)		

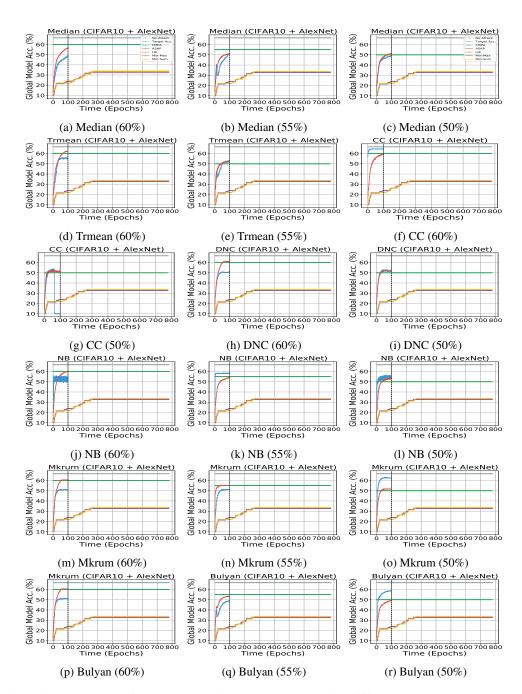


Figure 5: Comparison of each attack against various AGRs with different attack objectives on CI-FAR10 with AlexNet under different attack objectives and different attacks under the same target accuracy.

Impact of Non-IID degrees The impact of data heterogeneity on attack efficacy is assessed using CIFAR10 with Dirichlet concentration parameters spanning {0.1, 0.3, 0.5, 0.7, 0.9}, targeting 55% accuracy under Trmean aggregation. As presented in Figure 4b, ASAP successfully accomplishes the attack objectives while demonstrating robust outperform of existing attack strategies regardless of statistical heterogeneity intensity.

Impact of number of clients While our baseline experiments employ 50 total clients, we extend the evaluation to assess ASAP's scalability under larger federation sizes of 100, 150, and 200 par-

ticipants using CIFAR10 with a 55% target accuracy. Figure 4c demonstrates that ASAP maintains consistent superiority over competing attack methods across all federation scales.

Impact of clients sampling rates The impact of client sampling rate variations on attack performance is examined in Figure 4d. Experimental findings indicate that ASAP exhibits enhanced consistency and reduced performance variance relative to competing attack approaches across all sampling configurations.

A.4 THE USE OF LARGE LANGUAGE MODELS

We used Large Language Models (LLMs) to aid and polish writing in this paper. Specifically, LLMs were used to help improve the clarity, grammar, and flow of certain sections of the manuscript, and to assist in refining the presentation of ideas and ensuring consistent writing style throughout the paper. All core concepts, methodological contributions, experimental designs, results, and conclusions represent our original work. The LLMs did not contribute to the research ideation, experimental methodology, data analysis, or the generation of novel scientific insights. All content assisted by LLMs was thoroughly reviewed, fact-checked, and edited by the human authors to ensure accuracy and alignment with our intended contributions. The authors take full responsibility for all claims, results, and content presented in this work.