

---

# Learning to Attack Federated Learning: A Model-based Reinforcement Learning Attack Framework

---

Henger Li\*, Xiaolin Sun\*, and Zizhan Zheng

Department of Computer Science

Tulane University

New Orleans, LA 70118

{hli30, xsun12, zzheng3}@tulane.edu

## Abstract

We propose a model-based reinforcement learning framework to derive untargeted poisoning attacks against federated learning (FL) systems. Our framework first approximates the distribution of the clients' aggregated data using model updates from the server. The learned distribution is then used to build a simulator of the FL environment, which is utilized to learn an adaptive attack policy through reinforcement learning. Our framework is capable of learning strong attacks automatically even when the server adopts a robust aggregation rule. We further derive an upper bound on the attacker's performance loss due to inaccurate distribution estimation. Experimental results on real-world datasets demonstrate that the proposed attack framework significantly outperforms state-of-the-art poisoning attacks. This indicates the importance of developing adaptive defenses for FL systems.

## 1 Introduction

Federated learning (FL) is a promising machine learning framework that allows multiple devices with private data to jointly train a learning model (coordinated by a server) without sharing their local data. It has recently been applied to consumer digital products [34], credit risk prediction [1], drug discovery [35], and digital health [40]. However, federated learning systems are vulnerable to adversarial attacks [32] such as model poisoning attacks [7, 52, 15, 4], data poisoning attacks [5, 18, 20], and inference attacks [36, 22, 60]. To this end, various robust aggregation rules such as coordinate-wise median [54], trimmed mean [54], Krum [8], norm clipping [48], geometric median [38], and FLTrust [10] have been proposed. However, these defenses are mainly evaluated against manually crafted myopic attack policies [44]. Their robustness in the face of advanced attacks remains unknown.

Due to the distributed nature of FL systems, a malicious device typically has limited knowledge about benign devices and system dynamics. To fully reveal the vulnerabilities of FL systems, it is therefore crucial to develop strong attacks that can best utilize the limited global knowledge. In this work, we take a first step in this direction by considering the white-box attack setting where the attacker has some global knowledge about the FL system and the server's algorithm, but has no access to the private data of benign devices, a reasonable assumption for real-world FL systems. To derive strong adaptive attacks, we propose to leverage the power of model-based reinforcement learning (RL) by integrating *distribution learning* and *policy learning*. A key observation of our approach is that although accurate information about individual devices can be hard to obtain in FL, it is often possible to infer their aggregated data distribution from publicly available model updates, which

---

\*Equal contribution.

is sufficient to derive strong attacks. In particular, the set of malicious devices first cooperatively estimate the aggregated data distribution through gradient inversion [22, 60]. The learned distribution is then used to build a simulator of the FL environment, which is utilized to derive an adaptive attack policy through reinforcement learning. We focus on untargeted model poisoning in this work, where the malicious devices aim to reduce the accuracy of the global model as much as possible by sending crafted gradient information to the server. However, our proposed framework can potentially be applied to other types of attacks in both the white-box and the more challenging black-box settings.

Our model-based approach distinguishes from existing work on reinforcement learning based adversarial attacks against machine learning systems [47, 58, 59]. In particular, we consider a more realistic threat model where the attacker might not always be selected due to subsampling nor does it have prior information about the distribution of the aggregated data. The attackers need to efficiently learn the distribution along the federated learning process in real time. In contrast, previous works typically assume more powerful attackers that can attack at any time and have full knowledge about the environment. Thus, they typically adopt a purely model-free approach, which is infeasible in attacking FL systems due to the large number of samples needed to be effective.

**Our contributions.** We advance the state-of-the-art in the following aspects. First, we propose a novel reinforcement learning attack framework against federated learning systems by integrating distribution learning and policy learning. Second, we theoretically quantify the effect of inaccurate distribution learning and heterogeneous local data distributions on the optimal attack performance. Third, our experiments on real-world datasets demonstrate that our RL-based attack method consistently outperforms existing model poisoning attacks [7, 15, 52] and significantly reduces the global model accuracy even when robust aggregation rules are applied. These findings indicate the importance of developing adaptive defenses for FL systems.

## 2 Related Work

**Poisoning attacks and defenses.** To compromise the integrity of the target model in federated learning, both targeted poisoning attacks [7, 4, 5] that aim to misclassify a specific set of inputs, and untargeted attacks [15, 52, 43] aiming to reduce the global model accuracy have been proposed. Existing approaches typically adopt a heuristics-based method [52]) or optimize a myopic goal [15, 43]), and are usually sub-optimal, especially when a robust aggregation rule is adopted. Further, they often require access to benign agents' local updates or the accurate global model parameters of the next round [52, 15]) to make significant attack impact. In contrast, our reinforcement learning based attack requires less global knowledge and targets a long-term attack goal.

Various defenses have been proposed for model poisoning attacks including robust-aggregation-based approaches and detection-based approaches. The former includes dimension-wise filtering that considers each dimension of local updates separately [6, 54], client-wise filtering that aims to restrict or even remove the impact of potentially malicious clients [8, 8, 38, 48], and approaches that require the server to have access to a small amount of root data [10]. Our RL-based attack is effective against all these defenses. In addition, time-coupled robust aggregation methods [2, 3, 25] that target adaptive attacks and anomaly detection-based defenses [28] have been proposed recently. Our approach can potentially be extended to compromise them by explicitly encoding the history information into states or utilizing a recurrent structure.

**RL-based attacks.** Reinforcement learning has recently been utilized for developing strong attacks in various settings, including corrupting training data of online supervised and unsupervised learning [59], manipulating the combinatorial structure of graph data [13], and injecting malicious nodes into a graph [47]. RL-based attacks have also been developed to damage the performance of reinforcement learning itself, by perturbing the reward signals during the training stage [58] and corrupting the state signals received by an agent during the testing stage [57, 45]. However, these methods typically assume that the attacker has access to an accurate MDP model or simulator and has unlimited time for training the attack policy. In contrast, our method first builds a world model by learning a data distribution from the FL model updates and then constructs an approximate simulator for training our attacks. Both distribution learning and policy learning happen when FL training is ongoing. Further, previous work has mainly focused on attacking a single RL agent by an external agent rather than an insider attack in a distributed learning environment as we consider in this work.

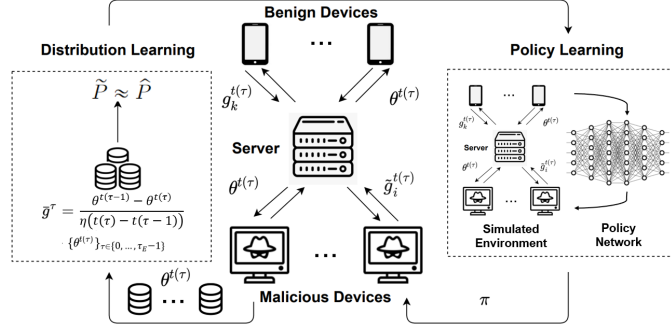


Figure 1: An overview of the RL-based attack framework against federated learning.

### 3 Approach Overview

In this section, we describe the federated learning setting considered in this work, the threat model, and the proposed RL-based attack framework.

**Federated learning.** We consider an FL setting that is similar to *federated averaging (FedAvg)* [33]. The FL system consists of a server and  $K$  workers (also known as devices or clients) in which each worker has some private data. Let  $\mathcal{K} = \{1, 2, \dots, K\}$  denote the set of workers. Coordinated by the server, the set of workers cooperate to train a machine learning model within  $\mathcal{T}$  epochs by solving the following problem:  $\min_{\mathcal{P}} \mathbb{E} F_{\mathcal{P}}$  where  $F_{\mathcal{P}} = \sum_{k=1}^K \rho_k F_k(\mathcal{P})$ . Here  $F_k(\mathcal{P})$  is the local objective of worker  $k$  and  $\rho_k$  is the weight assigned to worker  $k$  and satisfies  $\rho_k \geq 0$  and  $\sum_{k=1}^K \rho_k = 1$ . The local objective  $F_k(\mathcal{P})$  is usually defined as the empirical risk over worker  $k$ 's local data with model parameter  $\mathcal{P}$ . That is,  $F_k(\mathcal{P}) = \frac{1}{N_k} \sum_{j=1}^{N_k} \ell(\mathcal{P}; Z_{jk})$ , where  $N_k$  is the number of data samples available locally on worker  $k$ ,  $\ell(\cdot; \cdot)$  is the loss function, and  $Z_{jk} = \{x_{jk}, y_{jk}\}$  is the  $j$ th data sample that is drawn *i.i.d.* from some distribution  $P_k$ . It is typical to set  $\rho_k = \frac{N_k}{N}$ , where  $N = \sum_{k=1}^K N_k$  is the total number of data samples across workers.

If all the workers' local data distributions are the same (i.e.,  $P_k = P_{k'}$  for all  $k, k' \in \mathcal{K}$ ), we call the workers' data are *i.i.d.*; otherwise, the data are *non-i.i.d.* We write  $\mathcal{P}_k$  as the empirical distribution of the  $N_k$  data samples drawn from  $P_k$ , and let  $\mathcal{P} = \sum_{k=1}^K \frac{N_k}{N} \mathcal{P}_k$  denote the mixture empirical distribution across workers.

The FL algorithm (see Algorithm 1 in Appendix B) works as follows: at each time step  $t \geq 0$ , a random subset  $\mathcal{S}^t$  of size  $w$  is uniformly sampled without replacement from the workers set  $\mathcal{K}$  by the server for synchronous aggregation [30]. The process of selecting workers for aggregation is called *subsampling*. Let  $w/N$  denote the *subsampling rate*, i.e., the ratio of the selected workers number  $w$  to the total number of workers  $K$ . Each selected worker  $k \in \mathcal{S}^t$  then samples a minibatch  $b_k$  of size  $B$  from its local data distribution  $\mathcal{P}_k$ . The worker then calculates the average local gradient  $g_k^{t+1} = \frac{1}{B} \sum_{z \in b_k} \nabla \ell(\mathcal{P}^t; z)$  and sends the gradient to the server. The server then uses an aggregation rule to compute the aggregated gradient  $g^{t+1} = \text{Aggr}(g_{k_1}^{t+1}, \dots, g_{k_w}^{t+1})$  where  $k_i \in \mathcal{S}^t$ , and updates the global model parameters  $\mathcal{P}^{t+1} = \mathcal{P}^t + \eta g^{t+1}$  where  $\eta$  is the learning rate. The newly updated model parameters  $\mathcal{P}^{t+1}$  is then sent to the selected workers to perform another FL iteration.

**Threat Model.** We assume that among the  $K$  workers,  $M \leq K$  of them are malicious. Let  $\mathcal{A}$  denote the set of malicious workers. They are coordinated either by one leading attacker or an external agent. We refer such agent as a *leader agent*. These attackers are assumed to be *fully cooperative* and share the same goal of compromising the FL system. We consider untargeted model poisoning attacks in this work where the  $M$  cooperative attackers send crafted local updates  $g_k^t$  ( $k \in \mathcal{A}$ ) to the server in order to maximize the empirical loss  $F_{\mathcal{P}}$  (the batch size  $B$ ), and their local data distributions  $\mathcal{P}_k$  ( $k \in \mathcal{A}$ ) (but not the benign workers' local data distributions). We further assume that the attackers obtain information about the server's training algorithm (i.e., the white-box attack setting). This information includes the server's learning rate  $\eta$ , the subsampling rate  $w/N$ , the total number of workers  $K$ , the aggregation rule  $\text{Aggr}$ , and the total number of training epochs  $\mathcal{T}$ .

**RL-based online attack framework.** Our attack framework consists of the following three phases.

- **Distribution learning:** The malicious workers first jointly learn an approximation of  $\mathcal{P}$  from the model updates  $\mathbf{u}^t$  received from the server, using a gradient inversion based inference attack [19].
- **Policy learning:** The leader agent then builds a simulator of the FL environment using the attackers’ local data and the learned distribution. An optimal attack policy is then derived through reinforcement learning using data sampled from the simulator. Note that policy learning can start together with distribution learning and continue after a reasonable distribution is learned.
- **Attack execution:** The learned policy is distributed to all the malicious workers to generate attack actions. Note that attack execution can start once an initial policy is learned, which can be updated during attack execution.

We note that all the three phases happen while the federated learning process is ongoing, thus the lengths of these phases are important hyperparameters to be determined. For example, with more observations, an attacker can learn a more accurate distribution, which will help obtain a better attack policy. However, when the total time window available to attacks is limited, a longer distribution learning phase reduces the attack opportunities in Phase 3. Compared with a purely model-free approach, our model-based approach is more sample efficient, which is especially important for federated learning as a malicious worker can only attack when it is sampled by the server.

## 4 Model-based Reinforcement Learning Attack Framework

In this section, we first formulate the model poisoning attack problem as a Markov decision process (MDP). We then discuss our model-based reinforcement learning attack framework in more details.

### 4.1 Attackers’ problem as a Markov decision process

The attackers’ problem is formulated as an undiscounted MDP, denoted by  $\mathcal{M} = (\mathcal{S}; \mathbf{A}; T; r; H)$ , where

- $\mathcal{S}$  is the state space. Let  $t \in \{0; 1; \dots; H - 1\}$  denote the index of the attack step and  $t \in \{0; 1; \dots; T - 1\}$  the corresponding FL epoch when at least one attacker is selected by the server. The state at step  $t$  is defined as  $s = (\mathbf{p}^{(t)}; \mathcal{A}^{(t)})$  where  $\mathcal{A}^{(t)}$  is the set of attackers selected at time  $t$ , which is shared between all malicious workers.
- $\mathbf{A} = \mathcal{A}^M$  is the space of the attackers’ joint actions where each attacker shares the same action space  $\mathcal{A}$ . If attacker  $i$  is selected at  $t$ , its action  $a_i = \mathbf{g}_i^{t+1} \in \mathbb{R}^d$  is the local update that attacker  $i$  sends to the server at time step  $t$ , where  $d$  is the dimension of the model parameters. The only action available to an attacker not selected at  $t$  is  $K$ , indicating that the attacker does not send any information in that step.
- $T : \mathcal{S} \times \mathbf{A} \rightarrow \mathcal{P}(\mathcal{S})$  is the state transition function that represents the probability of reaching a state  $s' \in \mathcal{S}$  from the current state  $s \in \mathcal{S}$  when attackers choose actions  $a_1; \dots; a_M$ , respectively.
- $r : \mathcal{S} \times \mathbf{A} \times \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}$  is the reward function. We define the reward at step  $t$  as  $r = r(\mathbf{p}^{(t+1)}; \mathbf{p}^{(t)})$ , which is determined by the current state, the next state, and the joint attack actions and is shared by all the attackers.
- $H$  is the number of attack steps in each episode and we have  $t \in \{0; \dots; H - 1\}$ .

The attackers’ goal is to find a joint attack policy  $\pi = \pi_1; \dots; \pi_M$  that maximizes the expected total rewards over  $H$  attack steps, i.e.,  $\mathbb{E} \sum_{t=0}^{H-1} r_t$ , where  $\pi_i : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$  denotes a stationary policy of attacker  $i$  that maps the state to a probability measure over  $\mathcal{A}$ . Using the definition of  $r$ , this objective is equivalent to finding a joint policy  $\pi$  that maximizes  $\mathbb{E} \sum_{t=0}^{H-1} r_t(\pi)$ .

A key obstacle to solving the MDP is that both the transition probabilities  $T$  and the reward function  $r$  depend on the joint empirical distribution across workers  $\mathbf{p}_{k \in [K]}^{\mathbf{u}}$ , which is fixed but unknown to the attackers. Although model-free reinforcement learning can bypass this difficulty, it requires a large number of samples to be effective, which is infeasible in the online attacking scenario we consider. We therefore adopt model-based reinforcement learning as a principled approach for designing adaptive attacks in the online setting. An important observation is that although the joint empirical distribution  $\mathbf{p}_{k \in [K]}^{\mathbf{u}}$  is unknown, the attackers can learn an approximation of the mixture

distribution  $\hat{\mathcal{P}} = \sum_{k=1}^K \frac{N_k}{N} \mathcal{P}_k$ , denoted by  $\hat{\mathcal{P}}$ , from model updates shared by the server, which is often sufficient to simulate the behavior of benign agents and the server by assuming that each benign agent samples data from  $\hat{\mathcal{P}}$ . This gives rise to a new MDP  $\mathcal{M} = (\mathcal{S}; \mathbf{A}; T'; r'; H)$  where  $T'$  and  $r'$  are derived from  $\hat{\mathcal{P}}$ . Thus, our proposed model-based reinforcement learning attack framework naturally consists of the distribution learning, policy learning, and attack execution phases.

## 4.2 Distribution learning

Initially, the attackers do not perform model-poisoning attacks. Instead, they jointly learn a mixture distribution  $\hat{\mathcal{P}}$  from the model updates  $\theta^t$  using a gradient inversion based inference attack [19, 60]. Various gradient inversion attacks have been proposed in the literature. In this work, we adapt the *inverting gradients* (IG) method [19] to distribution learning. The IG method reconstructs data samples by optimizing a loss function based on the angle (i.e., cosine similarity) between the gradient generated from true data and that from the reconstructed data. The primary goal of IG is to reconstruct the original data samples, which is more ambitious than what the attackers need in our setting. On the other hand, recent works on gradient inversion including IG have focused on the server side, where the true gradients of each individual worker can be easily obtained from model updates. In contrast, the attackers only obtain approximated and aggregated gradient information from consecutive model updates received from the server, due to model aggregation and subsampling. Despite these differences, our experiment results show that the  $\hat{\mathcal{P}}$  learned using IG can help derive an effective attacker policy (see Figure 4(c)).

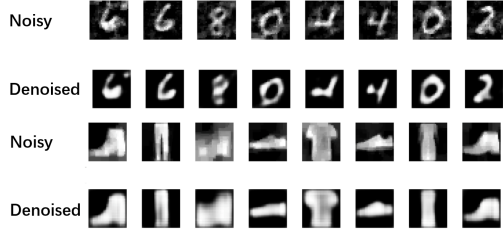


Figure 2: Examples of reconstructed images (before and after denoising) for MNIST (upper) and Fashion-MNIST (lower) datasets.

As shown in Algorithm 2 in Appendix B, for each epoch  $t \in [1, q]$  that at least one attacker is selected, the leader agent obtains the model update from one of the attackers and calculates the batch-level gradient as  $g^t = \frac{1}{B} \sum_{(x,y) \in D_{dummy}^t} \nabla_{\theta} \ell(\theta; x, y)$ . The leader agent then starts with a batch of (randomly generated) dummy data and dummy labels  $D_{dummy}$ , which is updated iteratively by solving the following optimization problem:  $\arg \min_{D_{dummy}} \frac{1}{B} \sum_{(x,y) \in D_{dummy}} \text{TV}(\rho_x, \rho_y)$ , where  $\text{TV}(\rho_x, \rho_y) = \frac{1}{2} \sum_{(x,y) \in D_{dummy}} |\rho_x(x) - \rho_y(y)|$  is the total variation [41] of  $x$ , and  $\beta$  is a fixed parameter. The process terminates after  $max\_iter$  iterations, then outputs the updated data as the reconstructed data samples. We observed that although the data samples generated by IG resemble true samples, they contain a certain amount of noise as shown in Figure 2, making the learned distribution less representative of the true distribution. To reduce noise, we adapt the method of denoising auto-encoder [50]. We utilize the clean data owned by attackers and add Gaussian noise to them to simulate the noise in the reconstructed data samples. The clean data and the synthetic noisy data are then paired to train an autoencoder for denoising, which is then used to remove the noise in reconstructed data samples as shown in the figure. The approximated mixture distribution  $\hat{\mathcal{P}}^t$  consists of the reconstructed data up to  $t \in [1, q]$  and the  $M$  attackers' local data. The distribution learning phase starts at the first FL epoch and continues for  $E$  steps. The learned distribution is shared with all the attackers.

Although we adopt the IG method in this work due to its simplicity, other more recent approaches such as the GradInversion method [55] and gradient inversion with a trained generative model [23] can be easily incorporated into our framework, which can potentially learn  $\hat{\mathcal{P}}$  in more challenging settings for complex datasets like ImageNet [14], deep networks, and large batch sizes. On the other hand, we show in the experiments that our RL-based attack trained using attackers' local data only is still effective and surpasses all the baselines, while distribution learning further boosts the attack performance. A detailed discussion on gradient inversion attacks and defenses is provided in Appendix C



new estimated distribution  $\hat{\mathcal{P}}$  is learned. Although the attackers may choose to start attacking during distribution learning, we observe that this can blur the gradient information and make distribution learning less accurate. Thus, we assume that each attacker starts attacking once the distribution learning phase is finished, and applies the latest learned attack policy during the remaining epochs of the federated learning process. During attack execution, each selected attacker first notifies other attackers so that every one knows the number of attackers that are sampled in that epoch. Each selected attacker then generates a crafted gradient according to the process described above with the parameters obtained from the latest learned policy.

The attackers' total training time (including distribution learning and policy learning) should be significantly less than the total FL training time so that the attackers have time to execute the attacks. In real-world FL training, the server usually must wait for some time (typically ranging from 1 minute to 10 minutes) before it receives responses from the clients [53, 9, 24]. In contrast, the leader agent does not incur such time cost in training attackers' policies using a simulated FL environment. Therefore, an epoch in policy learning is typically much shorter than an FL epoch, making it possible to train the attack policy with a large number of episodes. In addition, the leader agent is usually equipped with GPUs, or other parallel computing facilities and can run multiple training episodes in parallel [11]. We compare the actual running time of our RL-based attack against different defenses in our experiment setting in Appendix E.2.

## 5 Impact of Inaccurate Distribution Learning and Data Heterogeneity

Our model-based RL attack employs the estimated data distribution  $\hat{\mathcal{P}}$  to simulate the behavior of benign workers, which can suffer from two types of errors. First,  $\hat{\mathcal{P}}$  can be far away from the true mixture distribution  $\mathcal{P}$  due to inaccurate distribution learning. Second, benign workers may vary in their local data distributions  $\mathcal{P}_k$ , which cannot be fully captured by a single mixture distribution. In this section, we study how the attack performance is affected by these two factors, which provides insights into properly distributing resources between the three phases of our attack method.

Our analysis is adapted from recent works that study the impact of model inaccuracy on the performance of model-based reinforcement learning [56, 31, 58] by addressing two new challenges. First, we need to establish the connection between the inaccuracy in data distribution  $\hat{\mathcal{P}}$  and the inaccuracy in the corresponding MDP as both the reward function and the transition dynamics depend on  $\hat{\mathcal{P}}$ . Second, although there are different ways to measure the distance between two models [56], it makes more sense to use the 1-Wasserstein distance [49] to measure the distance between two data distributions. This, however, requires bounding the Lipschitz constant of the optimal value function [56]. Although this is a challenging task for general RL tasks, we are able to show that this is indeed the case in our setting under the following assumptions. The first assumption models the inaccuracy of distribution learning as well as the heterogeneity of benign workers' local data.

**Assumption 1.**  $W_1(\hat{\mathcal{P}}; \mathcal{P}_k) \leq \epsilon$  for any benign worker  $k$ , where  $W_1(p; q)$  is the 1-Wasserstein distance [49].

We further need the following standard assumptions on the loss function.

**Assumption 2.** Let  $Z$  denote the domain of data samples across all the workers. For any  $s_1; s_2 \in S$  and  $z_1; z_2 \in Z$ , the loss function  $\ell: S \times Z \rightarrow \mathbb{R}$  satisfies:

1.  $|\ell(s_1; z_1) - \ell(s_2; z_2)| \leq L \sqrt{\|s_1 - s_2\|_2^2 + \|z_1 - z_2\|_2^2}$  (Lipschitz continuity w.r.t.  $s$  and  $z$ );
2.  $\|\nabla_s \ell(s_1; z_1) - \nabla_s \ell(s_1; z_2)\|_2 \leq L_z \|z_1 - z_2\|_2$  (Lipschitz smoothness w.r.t.  $z$ );
3.  $\ell(s_2; z_1) \leq \ell(s_1; z_1) - \alpha \|\nabla_s \ell(s_1; z_1; s_2, s_1)\|_2 \sqrt{\|s_2 - s_1\|_2^2}$  (strong convexity w.r.t.  $s$ );
4.  $\|\nabla_s \ell(s_2; z_1) - \nabla_s \ell(s_1; z_1) - \alpha \nabla_s \ell(s_1; z_1; s_2, s_1)\|_2 \sqrt{\|s_2 - s_1\|_2^2}$  (strong smoothness w.r.t.  $s$ );
5.  $\ell; q$  is twice continuously differentiable with respect to  $s$ .

where  $\|\cdot\|_2$  is the L2 norm. For simplicity, we further make the following assumption on the FL environment, although our analysis can be readily applied to more general settings.

**Assumption 3.** The server adopts FedAvg without subsampling ( $w = K$ ). All workers have same amount of data ( $\rho_k = \frac{1}{K}$ ) and the local minibatch size  $B = 1$ . In each epoch of federated learning,

each normal worker’s local minibatch is sampled independently from the local empirical data distribution  $\mathcal{P}_k$ .

Since no subsampling is considered in this section, with a slight abuse of notation, we let index  $t$  denote both an attack step and the corresponding FL epoch. Let  $\mathcal{M} = \rho S; \mathbf{A}; T; r; Hq$  denote the true MDP for the attackers, and  $\hat{\mathcal{M}} = \rho S; \mathbf{A}; T'; r'; Hq$  the simulated MDP when the local distribution of any benign worker is estimated as  $\hat{\mathcal{P}}$ . The following theorem captures the attack performance loss due to inaccurate distribution learning (see Appendix D for the proof).

**Theorem 1.** Let  $\mathcal{J}_{\mathcal{M}}(\rho^* q) = \mathbb{E}_{\rho; T; \rho} \sum_{t=0}^{H-1} \rho s^t; a^t; s^{t+1} q$ s denote the expected return over  $H$  attack steps under  $\mathcal{M}$ , policy  $\rho^*$  and initial state distribution  $\rho_0$ . Let  $\rho^*$  and  $r^*$  be the optimal policies for  $\mathcal{M}$  and  $\hat{\mathcal{M}}$  respectively, with the same initial state distribution  $\rho_0$ . Then,

$$|\mathcal{J}_{\mathcal{M}}(\rho^* q) - \mathcal{J}_{\hat{\mathcal{M}}}(\rho^* q)| \leq 2H \rho L L_V L_Z \leq 2L_S;$$

where  $L_V = \frac{K-M}{K} \max_{|s| \leq 1} |u|$  is the fraction of benign nodes,  $L_V \leq \frac{H-1}{t=0} \rho K_F q^t \rho L = L K_F q$ , and  $K_F \leq \max_{|s| \leq 1} |u|$ .

In practice, the learning rate  $\rho$  is typically small enough so that  $\max_{|s| \leq 1} |u| \leq 1$ . In this case,  $L_V$  is bounded by  $\frac{L(1+K_F)}{1-K_F} \leq \frac{L(1+)}{1-}$ . Therefore, we have  $|\mathcal{J}_{\mathcal{M}}(\rho^* q) - \mathcal{J}_{\hat{\mathcal{M}}}(\rho^* q)| \leq \mathcal{O}(H \frac{1}{1-} q)$ . To ensure convergence, we typically have  $\rho \leq \frac{1}{\sqrt{H}} q$  [39], thus  $|\mathcal{J}_{\mathcal{M}}(\rho^* q) - \mathcal{J}_{\hat{\mathcal{M}}}(\rho^* q)| \leq \mathcal{O}(\frac{1}{\sqrt{H}} \frac{1}{1-} H q)$ . This result clearly shows how the loss of attack performance depends on the fraction of benign nodes, the inaccuracy of distribution learning, and the time horizon.

## 6 Experiments

In this section, we compare our RL-based attack with state-of-the-art model poisoning attacks on real-world datasets. Our code is available at <https://github.com/SliencerX/Learning-to-Attack-Federated-Learning>.

### 6.1 Experiment setup

**Datasets.** We conduct extensive experiments on four real-world datasets: MNIST [27], Fashion-MNIST [51], EMNIST [12], and CIFAR-10 [26]. Due to space limitation, experiment results for Fashion-MNIST, EMNIST, and CIFAR-10 are provided in Appendix E. For the *i.i.d.* setting, we randomly split the dataset into  $K$  groups, each of which consists of the same number of training samples. For the *non-i.i.d.* setting, we follow the method of [15] to quantify the heterogeneity of local data distribution across clients. Suppose there are  $C$  classes in the dataset, e.g.,  $C = 10$  for the MNIST and Fashion-MNIST datasets. We evenly split the worker devices into  $C$  groups (with the  $M$  attackers evenly distributed across the  $C$  groups), where each group is assigned  $\frac{1}{C}$  of training samples as follows. A training instance with label  $c$  is assigned to the  $c$ -th group with probability  $q \mathbb{1}_{\{c=C\}}$  and to every other group with probability  $\rho \mathbb{1}_{\{c \neq C\}}$ . Within each group, instances are evenly distributed. A higher  $q$  indicates a higher *non-i.i.d.* degree. We set  $q = 0.5$  as the default *non-i.i.d.* degree. To demonstrate the power of distribution learning, we assume that the set of attackers share  $m$  true data points sampled from the training instances assigned to them. We set  $m = 200$  as the default value for MNIST.

**Baselines.** We compare our RL-based attack (RL) with no attack (NA), and the state-of-the-art model poisoning FL attack methods: explicit boosting (EB) [7], inner product manipulation (IPM) [52], and local model poisoning attack (LMP) [15]. IPM manipulates the attackers’ gradients so that the inner product between the aggregation result and the true gradient is negative. This requires access to the average of normal workers’ gradients in each FL epoch, which is usually unavailable in practice. LMP generates myopic attacks by solving an optimization problem in each FL epoch. In addition to the server’s aggregation rule, it also requires access to normal workers’ local models. Although LMP with partial knowledge is also presented in [15], it performs substantially worse than the full knowledge case when the server uses the coordinate-wise median defense. We compare the RL-based attack with the more powerful full knowledge LMP below.



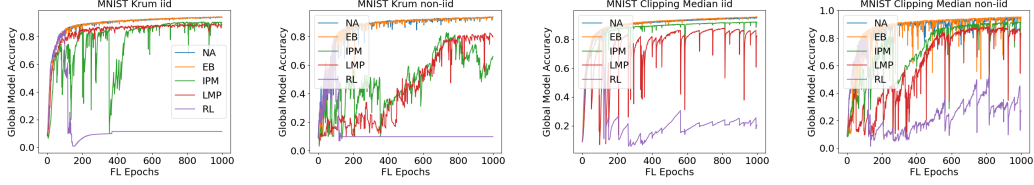


Figure 3: A comparison of global model accuracy under Krum and clipping median for both *i.i.d.* data and *non-i.i.d.* data. All parameters are set as default.

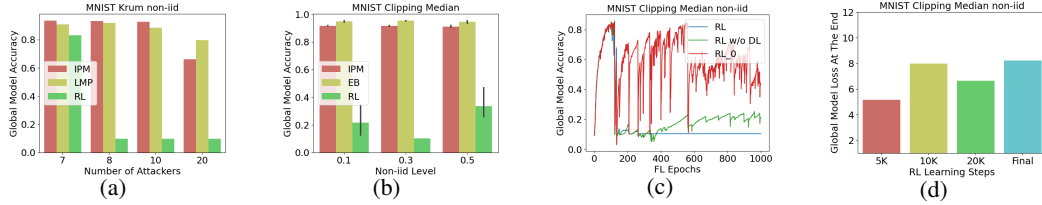


Figure 4: Attack performance on MNIST under (a) different number of attackers; (b) different non-iid degrees; (c) RL with and without distribution learning and  $RL_0$  (zero initial data for policy learning); and (d) different policy learning lengths. Non-iid degree  $q = 0.5$  in (a) and  $q = 0.3$  in (c) and (d). Other parameters are set as default.

We consider four representative robust aggregation rules of different types [44]: Krum [8], geometric median [38], both of which apply client-wise filtering to model updates, coordinate-wise median [54], which adopts a dimension-wise filtering, and FLTrust [10], which requires the server to have access to a small amount of root data. In the experiments, we actually consider an extension of the vanilla coordinate-wise median where a norm clipping [48] step is first applied. This gives a more powerful defense as we observed in experiments. We set the default norm threshold to 2.

**Default FL and RL settings.** We adopt the following default parameters for the FL models: number of total workers 100, number of attackers 20, learning rate 0.01, subsampling rate 10%, the number of total FL epochs 1,000. For our RL-based attack, both the distribution learning and policy learning phase start at the first FL epoch. The former ends at the 100th FL epoch when RL-based attack starts (all other attacks start at epoch 0). The policy learning phase ends at the 400th epoch. Since both the action space and state space are continuous in our setting, we choose the state-of-the-art Twin Delayed DDPG (TD3) [17] and Proximal Policy Optimization (PPO) [42] algorithms for training the attack policy in our experiments and find that TD3 gives better results in most cases. Below we report the results for TD3. We fix the initial model and the random seeds for subsampling and local data sampling for fair comparisons. See Appendix E for details of the datasets, experimental setups, and additional results.

## 6.2 Attack performance

Figure 3 shows how the test accuracy of the global model varies over FL epochs under different attacks when the server uses Krum and clipping median as the aggregation rule, respectively. Results for geometric median and FLTrust are provided in Appendix E.2. We observe that our RL-based attack performs significantly better in all the settings, despite the fact that IPM and LMP use the model updates of normal clients while RL does not. Note that for Krum, RL-based attack quickly drives the global model to a poor state ( $\sim 10\%$  accuracy) once the attack starts at epoch 100 under both *i.i.d.* and *non-i.i.d.* local data distributions. Attacks become harder under clipping median due to the norm clipping but our RL-based attack still reduces global model accuracy to around 50% on average. This is mainly because it targets long-term return while all other baselines are myopic. For example, in Figure 3(c), the global model accuracy drops significantly under all the attacks when five malicious devices are sampled around epoch 200. After that, the RL method keeps the accuracy at a low level, while other baselines’ accuracy rebounds rapidly.

### 6.3 Ablation studies

**Impact of the number of attackers.** Previous studies on untargeted model poisoning in federated learning typically assume a relatively large fraction of attackers. For example, the default setting is 20% in [15] and 40% in [52]. Figure 4(a) shows that our RL-based attack obtains superb performance even when the number of attackers (among 100 total clients) is as low as 8. In contrast, neither IPM nor LMP obtains meaningful attack performance even with 10 attackers. For 7 attackers, none of the baselines including the RL-based attack can cause significant damage to the FL system.

**Impact of non-i.i.d. degree.** Figure 4(b) shows the impact of data heterogeneity on attack performance. We use 5 different random seeds for all attacks and show the error bars. We observe that all the baselines obtain similar performance under different non-i.i.d. degrees and the impact of randomness in the testing environment on their performance is limited. On the other hand, we observe that the RL policies for  $q = 0.1$  and  $q = 0.5$  exhibit large variances, but even the worst-case performance of our attack outperforms the best cases of all the baselines. For  $q = 0.3$ , the RL-based attack can always lead to a model with a very high loss so that the model accuracy stays at a low level and is close to a constant, which explains the observed low variance in model accuracy. We expect that the variation across different RL policies is in part because the attackers always use the latest trained policy for attack execution, which does not necessarily give the best performance among all the intermediate policies trained (see also Figure 4(d)).

**Importance of distribution learning.** Figure 4(c) compares the global model accuracy of RL-based attack with distribution learning (RL) and that without distribution learning (RL w/o DL). We observe that in both cases, the model accuracy decreases dramatically after the attack starts at FL epoch 100. Further, the accuracy of RL w/o DL slightly increases up to 20%, while the accuracy of RL stays below 10%, which is consistent with our expectation that distribution learning allows the attackers to learn a better attack policy. Figure 4(c) also shows the attack performance of  $RL_0$ , a variant of the RL-based attack where the attackers only have 200 *unlabeled* true images used to train the denoising autoencoders, thus completely relying on distribution learning to generate labeled samples needed for policy learning. Compared with the baseline results in Figure 4(b) ( $q = 0.3$ ), we observe that  $RL_0$  still outperforms other baseline methods, further indicating the power of distribution learning. On the other hand, the fact that RL w/o DL surpasses all the baselines indicates that our approach is still applicable even when distribution learning becomes less effective in the presence of a strong defense against gradient inversion.

**Impact of training length on policy learning.** Figure 4(d) shows how the global model loss at the end of an FL training episode (in the simulated environment) varies over the RL policy training steps. We observe that longer training usually provides a better attack policy, although the training process is not stable. To fix this, one approach is to set up a separate testing environment to identify best trained policies. As mentioned above, our RL-based attack achieves promising performance even when the attackers always use the latest policy obtained during policy learning.

## 7 Conclusion

We propose a new approach for developing non-myopic attacks that can effectively compromise FL systems even with advanced defense mechanisms applied, by utilizing model-based reinforcement learning as a principled approach. While we focus on untargeted model poisoning against FL systems in this paper, our attack framework can be extended to targeted attacks (e.g., backdoor attacks) and to objectives beyond global model accuracy (e.g., fairness across clients [37, 29]). Further, our attack framework can be integrated with meta-learning [16, 21] to generalize the learned policy to different training tasks and develop black-box attacks. Another direction is to investigate novel methods to defend our adaptive attack methods. One possible solution would be to dynamically adjust FL parameters such as the subsampling rate or the aggregation rule.

## 8 Acknowledgment

This work has been funded in part by NSF grants CNS-1816495 and CNS-2146548 and Tulane University Jurist Center for Artificial Intelligence. We thank the anonymous reviewers for their valuable and insightful feedback.

## References

- [1] Utilization of FATE in Risk Management of Credit in Small and Micro Enterprises. <https://www.fedai.org/cases/utilization-of-fate-in-risk-management-of-credit-in-small-and-micro-enterprises/>.
- [2] Dan Alistarh, Zeyuan Allen-Zhu, and Jerry Li. Byzantine stochastic gradient descent. *Advances in Neural Information Processing Systems(NeurIPS)*, 31, 2018.
- [3] Zeyuan Allen-Zhu, Faeze Ebrahimiaghazani, Jerry Li, and Dan Alistarh. Byzantine-resilient non-convex stochastic gradient descent. In *International Conference on Learning Representations(ICLR)*, 2020.
- [4] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to backdoor federated learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020.
- [5] Gilad Baruch, Moran Baruch, and Yoav Goldberg. A little is enough: Circumventing defenses for distributed learning. In *Advances in Neural Information Processing Systems(NeurIPS)*, 2019.
- [6] Jeremy Bernstein, Jiawei Zhao, Kamyar Azizzadenesheli, and Anima Anandkumar. signsgd with majority vote is communication efficient and fault tolerant. In *International Conference on Learning Representations(ICLR)*, 2018.
- [7] Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. Analyzing federated learning through an adversarial lens. In *International Conference on Machine Learning(ICML)*, 2019.
- [8] Peva Blanchard, Rachid Guerraoui, Julien Stainer, et al. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Advances in Neural Information Processing Systems(NeurIPS)*, 2017.
- [9] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloé Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, Timon Van Overveldt, David Petrou, Daniel Ramage, and Jason Roselander. Towards federated learning at scale: System design. In *Proceedings of Machine Learning and Systems*, 2019.
- [10] Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. Fltrust: Byzantine-robust federated learning via trust bootstrapping. In *NDSS*, 2021.
- [11] Alfredo V Clemente, Humberto N Castejón, and Arjun Chandra. Efficient parallel methods for deep reinforcement learning. *arXiv preprint arXiv:1705.04862*, 2017.
- [12] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In *International Joint Conference on Neural Networks (IJCNN)*, 2017.
- [13] Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. Adversarial attack on graph structured data. In *International Conference on Machine Learning(ICML)*, 2018.
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition(CVPR)*, 2009.
- [15] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. Local model poisoning attacks to byzantine-robust federated learning. In *29th USENIX Security Symposium*, 2020.
- [16] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning(ICML)*, 2017.
- [17] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning(ICML)*, 2018.

- [18] Clement Fung, Chris JM Yoon, and Ivan Beschastnikh. Mitigating sybils in federated learning poisoning. *arXiv preprint arXiv:1808.04866*, 2018.
- [19] Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients - how easy is it to break privacy in federated learning? In *Advances in Neural Information Processing Systems(NeurIPS)*, 2020.
- [20] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *CoRR*, abs/1708.06733, 2017.
- [21] Abhishek Gupta, Benjamin Eysenbach, Chelsea Finn, and Sergey Levine. Unsupervised meta-learning for reinforcement learning. *arXiv preprint arXiv:1806.04640*, 2018.
- [22] Briland Hitaj, Giuseppe Ateniese, and Fernando Perez-Cruz. Deep models under the gan: information leakage from collaborative deep learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 2017.
- [23] Jinwoo Jeon, Kangwook Lee, Sewoong Oh, Jungseul Ok, et al. Gradient inversion with generative image prior. *Advances in Neural Information Processing Systems(NeurIPS)*, 2021.
- [24] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- [25] Sai Praneeth Karimireddy, Lie He, and Martin Jaggi. Learning from history for byzantine robust optimization. In *International Conference on Machine Learning(ICML)*, 2021.
- [26] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [27] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [28] Suyi Li, Yong Cheng, Wei Wang, Yang Liu, and Tianjian Chen. Learning to detect malicious clients for robust federated learning. *arXiv preprint arXiv:2002.00211*, 2020.
- [29] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In *International Conference on Machine Learning(ICML)*, 2021.
- [30] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. In *International Conference on Learning Representations(ICLR)*, 2020.
- [31] Yuping Luo, Huazhe Xu, Yuezhi Li, Yuandong Tian, Trevor Darrell, , and Tengyu Ma. Algorithmic framework for model-based deep reinforcement learning with theoretical guarantees. In *International Conference on Learning Representations(ICLR)*, 2019.
- [32] Lingjuan Lyu, Han Yu, and Qiang Yang. Threats to federated learning: A survey. *arXiv preprint arXiv:2003.02133*, 2020.
- [33] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- [34] Brendan McMahan and Daniel Ramage. Federated Learning: Collaborative Machine Learning without Centralized Training Data. <https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>.
- [35] Brendan McMahan and Daniel Ramage. Machine Learning Ledger Orchestration For Drug Discovery (MELLODDY). <https://www.melloddy.eu/>.

- [36] Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. Exploiting unintended feature leakage in collaborative learning. In *2019 IEEE Symposium on Security and Privacy*, 2019.
- [37] Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In *International Conference on Machine Learning(ICML)*, 2019.
- [38] Krishna Pillutla, Sham M Kakade, and Zaid Harchaoui. Robust aggregation for federated learning. *IEEE Transactions on Signal Processing*, 70:1142–1154, 2022.
- [39] Boris T. Polyak. *Introduction to optimization*. Optimization Software, 1987.
- [40] Nicola Rieke, Jonny Hancox, Wenqi Li, Fausto Milletari, et al. The future of digital health with federated learning. *NPJ digital medicine*, 3(1):1–7, 2020.
- [41] Leonid I Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992.
- [42] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [43] Virat Shejwalkar and Amir Houmansadr. Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning. In *NDSS*, 2021.
- [44] Virat Shejwalkar, Amir Houmansadr, Peter Kairouz, and Daniel Ramage. Back to the drawing board: A critical evaluation of poisoning attacks on production federated learning. In *IEEE Symposium on Security and Privacy*, 2022.
- [45] Yanchao Sun, Ruijie Zheng, Yongyuan Liang, and Furong Huang. Who is the strongest enemy? towards optimal and efficient evasion attacks in deep RL. In *International Conference on Learning Representations(ICLR)*, 2022.
- [46] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation from predicting 10,000 classes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition(CVPR)*, 2014.
- [47] Yiwei Sun, Suhang Wang, Xianfeng Tang, Tsung-Yu Hsieh, and Vasant Honavar. Adversarial attacks on graph neural networks via node injections: A hierarchical reinforcement learning approach. In *Proceedings of The Web Conference 2020*, 2020.
- [48] Ziteng Sun, Peter Kairouz, Ananda Theertha Suresh, and H Brendan McMahan. Can you really backdoor federated learning? *arXiv preprint arXiv:1911.07963*, 2019.
- [49] Leonid Nisonovich Vaserstein. Markov processes over denumerable products of spaces, describing large systems of automata. *Problemy Peredachi Informatsii*, 5(3):64–72, 1969.
- [50] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *International Conference on Machine Learning(ICML)*, 2008.
- [51] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [52] Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. Fall of empires: Breaking byzantine-tolerant sgd by inner product manipulation. In *Uncertainty in Artificial Intelligence (UAI)*, pages 261–270. PMLR, 2020.
- [53] Timothy Yang, Galen Andrew, Hubert Eichner, Haicheng Sun, Wei Li, Nicholas Kong, Daniel Ramage, and Françoise Beaufays. Applied federated learning: Improving google keyboard query suggestions. *arXiv preprint arXiv:1812.02903*, 2018.
- [54] Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In *Proceedings of the 35th International Conference on Machine Learning(ICML)*, 2018.

- [55] Hongxu Yin, Arun Mallya, Arash Vahdat, José Manuel Álvarez, Jan Kautz, and Pavlo Molchanov. See through gradients: Image batch recovery via gradinversion. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [56] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [57] Huan Zhang, Hongge Chen, Duane S Boning, and Cho-Jui Hsieh. Robust reinforcement learning on state observations with learned optimal adversary. In *International Conference on Learning Representations(ICLR)*, 2021.
- [58] Xuezhou Zhang, Yuzhe Ma, Adish Singla, and Xiaojin Zhu. Adaptive reward-poisoning attacks against reinforcement learning. In *International Conference on Machine Learning(ICML)*, 2020.
- [59] Xuezhou Zhang, Xiaojin Zhu, and Laurent Lessard. Online data poisoning attacks. In *Learning for Dynamics and Control*, pages 201–210. PMLR, 2020.
- [60] Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In *Advances in Neural Information Processing Systems(NeurIPS)*, 2019.