# ROBUST FEDERATED LEARNING FRAMEWORKS GUARDING AGAINST DATA FLIPPING THREATS FOR AUTONOMOUS VEHICLES

Anonymous authors

Paper under double-blind review

### ABSTRACT

Federated Learning (FL) has become an established technique to facilitate privacypreserving collaborative training across a multitude of clients. The ability to achieve collaborative learning from multiple parties containing an extensive volume of data while providing the essence of data privacy made it an attractive solution to address numerous challenges in sensitive data-driven fields such as autonomous vehicles (AVs). However, its decentralized nature exposes it to security threats, such as evasion and data poisoning attacks, where malicious participants can compromise training data. This paper addresses the challenge of defending federated learning systems against data poisoning attacks specifically focusing on data-flipping techniques in AVs by proposing a novel defense mechanism that combines anomaly detection with robust aggregation techniques. Our approach employs statistical outlier detection and model-based consistency checks to filter out compromised updates before they affect the global model. Experiments on benchmark datasets show that our method significantly enhances robustness by preventing nearly 15% of accuracy drop for our global model when confronted with a malicious participant and reduction the the attack success rate even when dealing with 20% of poisoning level. These findings provide a comprehensive solution to strengthen FL systems against adversarial threats.

032

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

## 1 INTRODUCTION

033 Autonomous driving Rajasekhar & Jaswal (2015)Martínez-Díaz & Soriguera (2018) is a rapidly de-034 veloping field that has the potential to revolutionize human transportation. the usage of machine learning in this field proved to by vary promising for it's application in various application, such 035 as autonomous driving Zhang et al. (2016), complex environment navigation Nguyen et al. (2020), 036 lane following and switching Gurghian et al. (2016) and traffic calculation Yu et al. (2021). Recent 037 advances in this area rely heavily on machine learning, that requires extensive training data. the centralized training provide more accuracy for autonomous driving solutions by using already known and controlled data. This approach neglect data privacy and third party involvement protocols. to 040 confront those issues Federated learning provide multiple advantages. 041

Federated Learning (FL) Kairouz et al. (2021) McMahan et al. (2017) Chaabene et al. (2022) has 042 become increasingly popular in machine learning, enabling models to learn from distributed devices 043 in diverse contexts. In the FL framework, participating peers undertake the task of training a global 044 model that they receive from the central server using their own local datasets. After processing their 045 local data, these peers generate model updates, which they then send back to the server. The server's 046 role is to gather and aggregate the various updates it receives from all the peers, ultimately leading 047 to the creation of an enhanced global model. Following this aggregation process, the updated global 048 model is redistributed to the peers, setting the stage for the next iteration of training. One of the 049 significant advantages of federated learning lies in its ability to enhance privacy. By keeping the local data on the devices of the peers and not transferring it to a central server, FL significantly 051 mitigates the risks associated with data breaches and privacy violations. This is particularly crucial in an era where data privacy concerns are paramount. Additionally, FL contributes to scalability 052 by distributing the computational workload across the devices of the peers, such as smartphones and other mobile devices, rather than relying solely on a central server's computing resources. This

decentralized approach not only improves efficiency but also allows for a broader range of devices to participate in the training process, making FL a highly attractive option for developing machine learning models in a secure and scalable manner Bonawitz et al. (2019).

057 But such as any technique FL faces many challenges and limitations; it's distributed nature exposes the global models to potential attacks from malicious participants. The lack of control of the service side Li et al. (2020), encourage malicious behaviour of the peers to tamper with the training 060 guidelines and conduct adversarial attack. Data poisoning represent one of the biggest challenges in 061 the usage of FL, the attackers take advantages of the model distribution to inject misclassified data 062 aiming to make the model fail or not converge, they want to make the model incorrectly classify 063 test cases with particular traits into some desired labels. One sort of targeted poisoning attack is the 064 label-flipping (LF) attack Biggio et al. (2012), in which the attackers tamper with training data in the local model by flipping the labels of select accurate instances from a transfer information from 065 a source class to a target class. Attackers train their local models after contaminating it using the 066 same hyper-parameters, loss function that that the server has supplied for the model architecture. 067 Thus, altering the training data is all that is needed to carry out the assault. That poisoned model 068 is later on send to the server to aggregate with the other intact model, causing a drop of the overall 069 performance and accuracy.

071 Multiple studies has been conducted to address those issues. Li et al. (2021) investigate peers that have the same goals as attackers, which causes a large percentage of false positives when 072 sincere peers have comparable local data. Zhou et al. (2023) present RoHFL, a hierarchical fed-073 erated learning framework for the Internet of Vehicles that uses similarity-based reputation scoring 074 and logarithm-based normalization to thwart poisoning assaults. Nevertheless, combining these 075 techniques entangle the aggregate procedure. The OQFL framework Yamany et al. (2023) uses 076 quantum-behaved particle swarm optimization (QPSO) to modify hyperparameters in order to iden-077 tify hostile cars. However, because the model must be reinitialized and retrained from start, there is a large computational overhead every time a search is conducted.

079 Our experiments, conducted on image classification datasets such as A2D2 Dataset, CIFAR-10 and a custom image collection used while training our AVs. We use Principal Component Analysis 081 (PCA) to reduce the dimensionality of update vectors and effectively differentiate between malicious and legitimate updates. And we combined it with Multiclass Classification Using Support 083 Vector Machines (SVM). Using PCA, the model can identify anomalies in the distribution of prin-084 cipal components that may indicate inconsistencies in labeling. Once the data is transformed, SVM 085 can be used to classify samples based on the extracted features. SVM decision thresholds can reveal patterns that indicate that some labels are inconsistent with the feature distribution. This may iden-087 tify flipped labels.,Furthermore, the ensemble method increases the robustness, of the search process. 088 This is because merging multiple SVM classifiers can reveal inconsistencies arising from additional labeling. It helps to be confident, that the model is resilient to adversarial attacks while maintaining 089 high classification performance... Our evaluations on the auxiliary datasets, demonstrate that our 090 defense strategy can effectively identify and block malicious participants. 091

092

# 2 RELATED WORK

094 095

Autonomous driving system based on federated learning. As participatory driving models con-096 tinue to improve in statistical accuracy, increased attention has been paid to improving the safety 097 and effectiveness of their training programs. Traditionally, driving intervention models have relied 098 on centralized training methods. However, centralized training presents challenges such as server computing capacity limitations, data security concerns, and network transmission overhead Yaacoub 100 et al. (2023). In response, the FD Framework has developed the capacity to nurture these models. 101 FL in Autonomous driving systems has been the subject of a verity of research investigations for 102 various purposes Chellapandi et al. (2023). FL is used, for example, in object detection, it makes 103 it possible for the AV framework to learn quickly and with little communication overhead, which 104 is especially useful when the amount of data is significantly more than the size of the ML model 105 while also protecting the data's privacy. In Barbieri et al. (2022), LiDAR on CAVs is utilized for object classification through a decentralized FL approach. Through V2V networks, the ML model's 106 parameters are exchanged. Comparing FL to selflearning techniques, it has been experimentally 107 demonstrated that FL is highly effective. The identification and recognition of license plates is a 108 significant additional use of FL. Applications include traffic safety and infractions, traffic monitor-109 ing, detecting unlawful or over-time parking, and parking access authentication are only a few of 110 the uses for it in ITS. It has been demonstrated that ML approaches are quite effective in identifying 111 license plates and detecting objects Kong et al. (2021)Xie et al. (2023). The Transformer model 112 has demonstrated the efficacy of the FL framework in learning spatio-temporal characteristics Zhou et al. (2022), all the while maintaining user privacy. The detection of abnormal vehicle trajectories at 113 traffic crossings has been accomplished through the use of FL in conjunction with OneClass Support 114 Vector Machine (OC-SVM) Koetsier et al. (2022). According to the published results, the federated 115 strategy enhances anomaly detection's overall accuracy while also benefiting specific data owners. 116 In other where FL provided promising results is predicting steering wheel angles and traffic control. 117 The performance of centralized learning and FL in steering angle prediction was evaluated in P et al. 118 (2021) under various noise levels, and the outcomes were equivalent. This research also took into 119 account the effects of communication load and interruptions, offering a thorough assessment of the 120 systems. Because of this, FL is appropriate for applications that include a growing number of CAVs, 121 particularly for jobs like steering wheel angle prediction. The research provided in Zhang et al. 122 (2021) showed that using FL in CAV significantly improved the quality of the edge models. In par-123 ticular, the research used optical flow and pictures as two data modalities to estimate steering wheel orientations. By employing FL to update the controller parameters dynamically, the target speed 124 has been better achieved while improving driver comfort and safetyZeng et al. (2022). Moreover, 125 FL is applied in cooperative parameter optimization between several vehicles at traffic junctions, 126 preventing crashes and enhancing driving comfort Wu et al. (2021). By precisely calculating the 127 road friction coefficients, FL is used in Liu et al. (2021) to improve brake performance in a variety 128 of driving scenarios and settings. This method maximizes the braking action while protecting the 129 driver's privacy. To optimize the controller design for AVs with variable vehicle participation in the 130 FL training process, a FL framework is proposed in Zeng et al. (2022). 131

Lable Flipping in Federated Learning. The rising popularity of FL has led to the exploration 132 of various attacks in this context, such as backdoor attacks Zhuang et al. (2024), gradient leakage 133 attacks Yang et al. (2024), and membership inference attacks Zhu et al. (2024). In work we focus 134 on data poisoning attacks[Data Poisoning Attacks Against Federated Learning Systems], such as 135 label-flipping (LF) Li et al. (2023a), Biggio et al. (2012) and feature perturbation (FP), are critical 136 areas of research. LF attacks have been widely applied in image processingPaudice et al. (2019). 137 For instance, Nowroozi et al. (2023) evaluated LF attacks and proposed a defense mechanism using 138 real datasets from the UCI repository, including MNIST Deng (2012) and Spambase Hopkins & 139 Suermondt (1999). Further experiments Rosenfeld et al. (2020) manipulated the MNIST dataset 140 using LF attacks and found a slight increase in classification error after injecting ten poisoning points. The study was repeated with Multiclass Logistic Regression, revealing an error increase 141 from 2% to 2.1% due to a random LF attack. The approach in Tolpegin et al. (2020) extended to 142 label-specific scenarios, where adversaries could adjust predictions based on predetermined rules. 143 Experiments on CIFAR-10 and a reduced version of ImageNet confirmed the effectiveness of the 144 proposed method. In order to protect against poisoning assaults, a number of studies concentrate 145 on evaluating certain updates. To differentiate between faulty and accurate updates, Jebreel et al. 146 (2020) suggest examining the biases in the output layer. But only in the IID setting does it take 147 model poisoning attacks into account. In order to prevent data poisoning attacks, FGold Fung et al. 148 (2020) and Awan et al. (2021) evaluate the weights of the output layer; But these techniques also 149 frequently penalize identical but good updates mistakenly, which causes the model's performance to 150 significantly decline. Using a kernel density estimator, Li et al. (2023b) calculates how harmful each 151 local update is in relation to its k-nearest neighbors. After that, it uses an asymptotic threshold to determine if updates are benign or poisoned. Not only is it difficult to choose a threshold of this kind, 152 but this approach has not been validated with big DL models or non-IID data. In order to identify 153 the LF attack, Qayyum et al. (2022) suggests a method for discovering the correlation between the 154 latent features of training data and updates. However, the strategy imposes an additional cost on all 155 parties to train another model that learns such relationship. Furthermore, it is unrealistic to believe 156 that throughout the early training rounds, all peers will behave appropriately. 157

158

# 3 STUDY DESIGN

159 160

For our experiment we used three SunFounder Picar Sunfounder that we trained using a CNN (Convolutional Neural Network) model. The collected data from the integrated camera where used to



Figure 1: Proposed Federated Learning model architecture for Label flipping Attack on for AV

create our own dataset that contained 5K images of traffic light, lines and different types of obstacle.
Moreover to address the data shortage we used the CIFAR10 dataset CIFAR-10 with it 60 K colored images of 10 different classes and we divided it into 50K data for training and 10K for testing and The A2D2 Dataset Audi that features over 40k labeled with 38 features.

We implement our FL framework for malicious vehicle detection using N = 3 participants, one central aggregator, and k = 5 each . We use an independent and identically distributed (iid ) data distribution, we assume the total training dataset is uniformly randomly distributed among all participants with each participant receiving a unique subset of the training data. The testing data is used for model evaluation only and is therefore not included in any participant Pi's train dataset. Observing that both CNN models converge after fewer than 200 training rounds, we set our FL experiments to run for R = 200 rounds total.

197 We trained our federated learning models on each dataset without adversarial settings. Next, the appropriate global models on test samples for each dataset. We first determined samples with high prediction confidence by computing softmax probabilities and choosing the examples that correctly 199 forecasted in order to generate the Complementary dataset. If the projected class matched the ac-200 tual label and its corresponding probability exceeded the threshold, we added the sample to the our 201 dataset. Participants within the federated learning framework must maintain continuous commu-202 nication and collaboration with the aggregation server. The model M is finished with parameters 203  $\theta_R$  at the end of R rounds of FL. The test dataset used to evaluate M is denoted by  $D_{\text{test}}$ , where 204  $D_{\text{test}} \cap D_i = \emptyset$  for each participant dataset  $D_i$ . We present an in-depth analysis of label flipping 205 attacks in FL in the following sections.

206 207 208

185

#### 3.1 LABEL FLIPPING ATTACK

We use a label flipping attack to implement targeted data poisoning in FL. Given a source class  $C_{\text{Source}}$  and a target class  $C_{\text{target}}$  from C, each malicious participant Pi modifies their dataset Di as follows: For all instances in Di whose class is  $C_{\text{Source}}$ , change their class to  $C_{\text{target}}$ . We denote this attack by  $C_{\text{Source}} \rightarrow C_{\text{target}}$ . For example, images with initial red light class labels may be altered to have a green light class by malicious participants, according to the CIFAR-10 image classification sign red light  $\rightarrow$  green light. The attack tries to increase the possibility that, during testing, the final global model would mistakenly classify traffic signals. The threat of label flipping is well-known in centralized machine learning. It's also acceptable in the FL scenario given the hostile objective

Ree	<b>[uire:</b> $n > 0$ the amount of clients
(	$\overline{C}NN \leftarrow create\_model()$
f	or i in [0, n] do
	$CNN[i].initiate(initial_packets)$
e	nd for
f	or 1 in $[0, n]$ in parallel <b>do</b>
	open_port()
	$awaii\_crient\_connection()$ CNN[i] = massive elient CNN()
	$CNN[i] \leftarrow receive\_crieni_CNN()$
Re	mire suspicious packet
r	eturn CNN.predict(suspicious packet)
Alg	orithm 2 Client-side code
	$CNN \leftarrow create \ decision \ CNN()$
(	CNN initiate (initial packets)
$\tilde{c}$	CNN.train(local_Data)
0	onnect_to_server()
s	end_CNN_to_server()
	v'
and	capabilities indicated above. Label flipping is different from other poisoning methods in that
the	adversary does not need to know the CNN architecture. loss function L, global distribution of D.
etc.	Its time and energy efficiency make it a desirable feature, especially since FL is frequently used
wit	a edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve
wit cha	n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a
wit cha fed	n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is P malicious, we proceed as follows. At
wit cha fed the	n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total P as
with cha fed the mat	n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then
wit cha fed the mal inje	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P = 3$ , where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results.
wit cha fed the mal inje We	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P = 3$ , where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions:
wit cha fed the mal inje We	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions:
with cha fed the mal inje We	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas-
with cha fed the mal inje We	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training.
with cha fed the mal inje We	<ul> <li>n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with P<sup>"</sup>= 3, where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate N × m% of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions:</li> <li>A source class → C<sub>target</sub> class pairing where the source class was very frequently misclassified as the C<sub>target</sub> class in federated, non-poisoned training.</li> <li>A pairing where the source class was very infrequently misclassified as the C<sub>target</sub> class.</li> </ul>
wit cha fed the mal inje We	<ul> <li>n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with P<sup>"</sup>= 3, where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate N × m% of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions:</li> <li>A source class → C<sub>target</sub> class pairing where the source class was very frequently misclassified as the C<sub>target</sub> class in federated, non-poisoned training.</li> <li>A pairing where the source class was very infrequently misclassified as the C<sub>target</sub> class.</li> </ul>
wit cha fed the mail inje We	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P = 3$ , where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes.
witi chaa fed the mal inje We	<ul> <li>n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with P<sup>*</sup> = 3, where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate N × m% of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions:</li> <li>A source class → C<sub>target</sub> class pairing where the source class was very frequently misclassified as the C<sub>target</sub> class in federated, non-poisoned training.</li> <li>A pairing where the source class was very infrequently misclassified as the C<sub>target</sub> class.</li> <li>A pairing between these two extremes.</li> </ul>
witt cha fed the mainje We The	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes.
witt cha fed the mal inje We The flip	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes.
witt chaa fed the mai inje We The flip	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes.
witi cha fed the malinje We The flip 3.2	<ul> <li>n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a crated learning (FL) system with P<sup>**</sup> = 3, where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate N × m% of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions:</li> <li>A source class → C<sub>target</sub> class pairing where the source class was very frequently misclassified as the C<sub>target</sub> class in federated, non-poisoned training.</li> <li>A pairing where the source class was very infrequently misclassified as the C<sub>target</sub> class.</li> <li>A pairing between these two extremes.</li> </ul>
witi cha fed the malinje We The flip 3.2	<ul> <li>n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a trated learning (FL) system with P"= 3, where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate N × m% of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions:</li> <li>A source class → C<sub>target</sub> class pairing where the source class was very frequently misclassified as the C<sub>target</sub> class in federated, non-poisoned training.</li> <li>A pairing where the source class was very infrequently misclassified as the C<sub>target</sub> class.</li> <li>A pairing between these two extremes.</li> </ul>
witt cha fed the malinje We The flip 3.2 We	<ul> <li>n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a trated learning (FL) system with P"= 3, where one is P malicious, we proceed as follows. At start of each experiment, we randomly designate N × m% of the participants from the total P as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions:</li> <li>A source class → C<sub>target</sub> class pairing where the source class was very frequently misclassified as the C<sub>target</sub> class in federated, non-poisoned training.</li> <li>A pairing where the source class was very infrequently misclassified as the C<sub>target</sub> class.</li> <li>A pairing between these two extremes.</li> <li>se conditions provide a diverse set of scenarios to evaluate the effectiveness and impact of label ping attacks in the federated learning environment.</li> </ul>
witt chaa fed the malinje We The flip 3.2 We Glo	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes. se conditions provide a diverse set of scenarios to evaluate the effectiveness and impact of label ping attacks in the federated learning environment. ATTACK EVALUATION METRICS employ several evaluation indicators to do this. bal Model Accuracy ( $M_{acc}$ ): The global model accuracy is the percentage of instances $x \in D_{trest}$
witi chaa fed the mai inje We The flip 3.2 We Glo who	n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P^{"}=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes. see conditions provide a diverse set of scenarios to evaluate the effectiveness and impact of label ping attacks in the federated learning environment. ATTACK EVALUATION METRICS employ several evaluation indicators to do this. bal Model Accuracy ( $M_{acc}$ ): The global model accuracy is the percentage of instances $x \in D_{test}$ are the global model $M$ with final parameters $\theta_R$ predicts $M_{\theta_P}(x) = c_i$ and $c_i$ is indeed the true
witi cha fed the mal inje We The flip 3.2 We <b>Gle</b> who class	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P^{"}=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes. se conditions provide a diverse set of scenarios to evaluate the effectiveness and impact of label ping attacks in the federated learning environment. ATTACK EVALUATION METRICS employ several evaluation indicators to do this. <b>bal Model Accuracy</b> ( $M_{acc}$ ): The global model accuracy is the percentage of instances $x \in D_{test}$ ere the global model $M$ with final parameters $\theta_R$ predicts $M_{\theta_R}(x) = c_i$ and $c_i$ is indeed the true s label of $x$ .
witi chaa fed the mal inje We The flip 3.2 We <b>Glo</b> who class	n edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P''=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes. see conditions provide a diverse set of scenarios to evaluate the effectiveness and impact of label ping attacks in the federated learning environment. ATTACK EVALUATION METRICS employ several evaluation indicators to do this. <b>bal Model Accuracy</b> ( $M_{acc}$ ): The global model accuracy is the percentage of instances $x \in D_{test}$ or the global model $M$ with final parameters $\theta_R$ predicts $M_{\theta_R}(x) = c_i$ and $c_i$ is indeed the true s label of $x$ .
witt chaa fed the mainje We The flip 3.2 We Gld who clas Cla	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a erated learning (FL) system with $P^{-}=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes. see conditions provide a diverse set of scenarios to evaluate the effectiveness and impact of label ping attacks in the federated learning environment. ATTACK EVALUATION METRICS employ several evaluation indicators to do this. bal Model Accuracy ( $M_{acc}$ ): The global model accuracy is the percentage of instances $x \in D_{test}$ ere the global model $M$ with final parameters $\theta_R$ predicts $M_{\theta_R}(x) = c_i$ and $c_i$ is indeed the true s label of $x$ . ss Recall ( $c_{recall_t}$ ): Where the percentage
witt chaa fed the mainje We The flip 3.2 We Glo who clas Clas	he edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a parated learning (FL) system with $P''=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes. se conditions provide a diverse set of scenarios to evaluate the effectiveness and impact of label ping attacks in the federated learning environment. ATTACK EVALUATION METRICS employ several evaluation indicators to do this. <b>bal Model Accuracy</b> ( $M_{acc}$ ): The global model accuracy is the percentage of instances $x \in D_{test}$ re the global model $M$ with final parameters $\theta_R$ predicts $M_{\theta_R}(x) = c_i$ and $c_i$ is indeed the true s label of $x$ . <b>ss Recall</b> ( $c_{recall_i}$ ): Where the percentage
witt chaa fed the mainje We The flip 3.2 We <b>Gla</b>	he edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a grated learning (FL) system with $P''=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then icious, while the remaining participants are considered honest. The malicious participants is then icious, while the remaining participants are considered honest. The malicious participants is then icious, while the remaining participants are considered honest. The malicious participants is then icious with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclas- sified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes. se conditions provide a diverse set of scenarios to evaluate the effectiveness and impact of label ping attacks in the federated learning environment. ATTACK EVALUATION METRICS employ several evaluation indicators to do this. <b>bal Model Accuracy</b> ( $M_{acc}$ ): The global model accuracy is the percentage of instances $x \in D_{test}$ re the global model $M$ with final parameters $\theta_R$ predicts $M_{\theta_R}(x) = c_i$ and $c_i$ is indeed the true s label of $x$ . <b>ss Recall</b> ( $c_{recall_i}$ ): Where the percentage $\frac{T_{P_i}}{T_{P_i} + F_{N_i}} \cdot 100\%$
witi chaa fed the mainje We The flip 3.2 We Glo who clas Clas	h edge devices. In addition, it is simple enough for non-experts to perform and doesn't involve nging or meddling with participant-side FL software. To simulate the label flipping attack in a strated learning (FL) system with $P^*=3$ , where one is $P$ malicious, we proceed as follows. At start of each experiment, we randomly designate $N \times m\%$ of the participants from the total $P$ as icious, while the remaining participants are considered honest. The malicious participants is then cted with flipped labels, each experiment is repeated 10 times, and we report the average results. examine three label flipping attack settings that represent a range of adversarial conditions: • A source class $\rightarrow C_{target}$ class pairing where the source class was very frequently misclassified as the $C_{target}$ class in federated, non-poisoned training. • A pairing where the source class was very infrequently misclassified as the $C_{target}$ class. • A pairing between these two extremes. se conditions provide a diverse set of scenarios to evaluate the effectiveness and impact of label ping attacks in the federated learning environment. ATTACK EVALUATION METRICS employ several evaluation indicators to do this. <b>bal Model Accuracy</b> ( $M_{acc}$ ): The global model accuracy is the percentage of instances $x \in D_{test}$ are the global model $M$ with final parameters $\theta_R$ predicts $M_{\theta_R}(x) = c_i$ and $c_i$ is indeed the true s label of $x$ . <b>ss Recall</b> ( $c_{recall_i}$ ): Where the percentage

represents the class recall for any class  $c_i \in C$ . Whereas  $F_{N_i}$  is the number of examples  $x \in D_{\text{test}}$ where  $M_{\theta_R}(x) \neq c_i$  and the true class label of x is  $c_i$ . The number of instances  $x \in D_{\text{test}}$  is  $T_{P_i}$ , where  $M_{\theta_R}(x) = c_i$  and  $c_i$  is the true class label of x.

269



Alg	orithm 3 Federated Learning Algorithm with Adversarial Mitigation and PCA
Req	uire: P: Total number of AV nodes
Req	<b>uire:</b> K: Participating vehicles during aggregation
Req	uire: R: Federation rounds
Req	<b>uire:</b> $D_C$ : Initial dataset
Req	uire: $\mathcal{A}$ : Adversary-controlled vehicles list
Req	uire: $M_s$ : MCSVM model
Req	uire: $\mathbf{W}_l$ : Local model
Ens	<b>ure:</b> $\mathbf{W}_t$ : Global model shared with all peers in $R$ -federation round
1:	$\mathbf{W}_0 \leftarrow \text{initialize global model}$
2:	for $t = 0$ to $T - 1$ do
3:	$S \leftarrow$ random set of K participants (Client Identifier)
4:	$C$ sends $\mathbf{W}_t$ to all participants in $S$
5:	for each participant $k \in S$ in parallel <b>do</b>
6:	if $k \in \mathcal{A}$ then
7:	Poison data $D_k$ through LF
8:	end if
9:	$\mathbf{W}_{t+1}^{k} \leftarrow \text{ClientUpdate}(k, \mathbf{W}_{t})$
10:	end for
11:	for each participant $k \in S$ do
12:	<b>Apply PCA:</b> Transform $\mathbf{W}_l$ using PCA to reduce dimensionality:
	$\mathbf{W}_l = PCA_{\mathbf{W}_l}$
13:	for each sample <i>i</i> in $W_i$ do
	$M_{\rm s}(W_l)$
14:	end for
15:	Test the <b>w</b> model and compute $\hat{y}_{i,k}^L$ for each $M_s$ ( $W_l$ )
16:	Test the model $\mathbf{W}_{t+1}^k$
17:	Compute Outlier_Score <sub>t+1</sub>
18:	end for
19:	$O \leftarrow \text{Select } \tau \text{ participants with the highest Outlier}_Score_{t+1} \text{ scores}$
20:	$G \leftarrow P - O$ {Remove malicious participants from P}
21:	Perform aggregation of models $\mathbf{W}_{t+1}$
22:	end for

### 3.4 EVALUATION INDICATORS

359

360

361

364

365

366

367

368

369

370

371

372

373 374

375

Throughout this paper, we employ metrics to enhance our understanding of both the security and utility provided by the models under scrutiny, which we subject to experimental manipulation. These 362 metrics have been defined as follows: 363

- Source Class Recall: This metric calculates the number of correct positive predictions made out of all positive predictions that could have been made by the model. =In the event of label tampering by a malicious user, this metric will decrease, as fewer (or none) correct positive predictions will be made for the specific class  $C_{\text{target}}$  by the attacker.
  - Sparse Categorical Accuracy: This metric evaluates the accuracy of a model's predictions by comparing the predicted class labels with the actual ground truth labels.
- CrossEntropy Loss: This measures the disparity between the predicted probability distribution and the true probability distribution of the classes. In our models, it serves to quantify how well the predicted probabilities align with the actual class labels.

#### RESULTS 4

376 First, we study the effect of a single malicious participant on our federated learning framework, FL. 377 We find that just a single malicious participant can significantly degrade global model performance-



Figure 3: Evolution of the source class recall by Figure 4: Evaluation of the Global model accuround when  $\alpha = 0.8$  racy with a malicious participant

395 source class recall losses of over 25% are possible when this adversary is consistently well-396 represented in the participant pool. This impact on source class recall is highest for a high availability 397 level of  $\alpha = 0.9$ . In this respect, the effect becomes smaller if availability goes down, meaning here that lower losses can be obtained for lower values:  $\alpha = 0.7$ ,  $\alpha = 0.6$  or  $\alpha = 0.5$ . Thus, to max-398 imize effectiveness of the attack it is beneficial for the malicious participant to remain as available 399 as possible-especially in the later training rounds. To further demonstrate this effect of availability, 400 we report source class recall by round for  $\alpha = 0.7$  and  $\alpha = 0.9$ . A higher availability of the mali-401 cious participant results in a noticeable degradation in source class recall, along with lowering the 402 value of recall for  $\alpha = 0.9$  compared to  $\alpha = 0.7$ . The probabilistic selection of the participants 403 can be considered one of the prime reasons for the variability of recall across rounds. A round with 404 fewer malicious participants tends to increase source recall, and a higher number falls back. Each 405 experimental condition is run three times, and the outcomes are averaged to remove round-to-round 406 variability. As our results indicate, even a single malicious participant can significantly reduce global 407 model performance, with source class recall losses of over 25% possible under high availability. Indeed, with high availability, there is a negative effect, while a decrease in availability tends to grant 408 considerably better results. Importantly, for values of k significantly larger than  $N \times m\%$ , increas-409 ing availability ( $\alpha$ ) becomes less effective for meaningful impacts in individual training rounds. As 410 for the accuracy we notice that each round of training using the malicious participant effected the 411 model accuracy by a drop of 0.1% with each training round leading to an overall loss of 20% when 412 to model finish training. 413

414 Using our proposed defense mechanism enables the detection such a malicious participant and never allows any updates from that participant or blacklists the participant for further usage in rounds. 415 Leading to no accuracy lost and the integrity of the global model training. Using Principal Com-416 ponent Analysis (PCA) and Multi-class Support Vector Machine (SVM) classifiers together offers a 417 strong way to defend against label flipping attacks in federated learning. PCA cuts down the number 418 of features and gets rid of noise. This proves important in federated setups where data quality often 419 changes from one user to another. Focusing on the most useful features helps the model work faster 420 and spot unusual data more easily. After the data gets changed by PCA, the usage of Multiclass 421 Classification SVM classifiers to draw complex lines is helpful for our agent in telling apart good 422 labels from malicious ones. This method really works in finding strange data and keeping correct 423 results even with malicious participation.

424 425

392

393 394

5 CONCLUSION

426

In this paper, we investigated data poisoning attacks targeting Federated Learning (FL) systems in autonomous vehicles. Our study reveals the susceptibility of FL systems to label flipping poisoning attacks, highlighting their significant adverse effects on the global model. We established a defence mechanism biased on PCA and MCSVM do to their ability to separate outliers. We further investigated the results of our mechanism which lead to the avoidance of the attack and the preservation of the global model integrity.

432 Future Work: Since this approach was based on a real life simulation with real cars the number 433 of participant was limited to three do the high necessity of computation cost. we aim to extend our 434 approach on multiple participant. Also the recreation of our strategy with a computer simulation 435 using different nodes and compare it to stats of the art solution. Another critical area is to test our 436 defence strategy against other types of adversarial attack such as backdoor and noise injection. One other critical area to look into is domain adaptation since a participant can contain different data 437 from the original but not harmless. Creating a distinguish between outlier and other domain data is 438 crucial. 439

- 440 441 **R**EEE
- 441 REFERENCES
- Audi. Driving Dataset. URL https://www.a2d2.audi/a2d2/en.html.
- Sana Awan, Bo Luo, and Fengjun Li. CONTRA: Defending Against Poisoning Attacks in Federated
  Learning. In Elisa Bertino, Haya Shulman, and Michael Waidner (eds.), *Computer Security – ESORICS 2021*, pp. 455–475, Cham, 2021. Springer International Publishing. ISBN 978-3-03088418-5. doi: 10.1007/978-3-030-88418-5\_22.
- Luca Barbieri, Stefano Savazzi, Mattia Brambilla, and Monica Nicoli. Decentralized federated learning for extended sensing in 6G connected vehicles. *Veh. Commun.*, 33(C), January 2022.
  ISSN 2214-2096. doi: 10.1016/j.vehcom.2021.100396. URL https://doi.org/10. 1016/j.vehcom.2021.100396.
- Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning Attacks against Support Vector Machines. June 2012.
- Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, 455 Vladimir Ivanov, Chloé Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMa-456 han, Timon Van Overveldt, David Petrou, Daniel Ramage, and Jason Roselander. To-457 wards federated learning at scale: System design. In A. Talwalkar, V. Smith, and 458 M. Zaharia (eds.), Proceedings of Machine Learning and Systems, volume 1, pp. 374-459 388, 2019. URL https://proceedings.mlsys.org/paper\_files/paper/2019/ 460 file/7b770da633baf74895be22a8807f1a8f-Paper.pdf. 461
- Riadh Ben Chaabene, Darine Amayed, and Mohamed Cheriet. Leveraging Centric Data Federated
   Learning Using Blockchain For Integrity Assurance, June 2022. URL http://arxiv.org/
   abs/2206.04731. arXiv:2206.04731 [cs].
- Vishnu Pandi Chellapandi, Liangqi Yuan, Stanislaw H Żak, and Ziran Wang. A Survey of Federated Learning for Connected and Automated Vehicles. In 2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC), pp. 2485–2492, September 2023. doi: 10. 1109/ITSC57777.2023.10421974. URL https://ieeexplore.ieee.org/document/ 10421974. ISSN: 2153-0017.
- 470 471 CIFAR-10. CIFAR-10 and CIFAR-100 datasets. URL https://www.cs.toronto.edu/ ~kriz/cifar.html.
- 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/ document/6296535. Conference Name: IEEE Signal Processing Magazine.
- Clement Fung, Chris J. M. Yoon, and Ivan Beschastnikh. The limitations of federated learning in sybil settings. In 23rd International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2020), pp. 301–316, San Sebastian, October 2020. USENIX Association. ISBN 978-1-939133-18-2. URL https://www.usenix.org/conference/raid2020/presentation/fung.
- Alexandru Gurghian, Tejaswi Koduri, Smita V. Bailur, Kyle J. Carey, and Vidya N. Murali.
   DeepLanes: End-To-End Lane Position Estimation Using Deep Neural Networks. In 2016 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 38– 45, June 2016. doi: 10.1109/CVPRW.2016.12. URL https://ieeexplore.ieee.org/ document/7789502. ISSN: 2160-7516.

486 Reeber Erik Forman George Hopkins, Mark and Jaap Suermondt. Spambase. UCI Machine Learning 487 Repository, 1999. DOI: https://doi.org/10.24432/C53G6X. 488

- Najeeb Jebreel, Alberto Blanco-Justicia, David Sánchez, and Josep Domingo-Ferrer. Efficient De-489 tection of Byzantine Attacks in Federated Learning Using Last Layer Biases. In Vicenç Torra, 490 Yasuo Narukawa, Jordi Nin, and Núria Agell (eds.), Modeling Decisions for Artificial Intelli-491 gence, pp. 154–165, Cham, 2020. Springer International Publishing. ISBN 978-3-030-57524-3. 492 doi: 10.1007/978-3-030-57524-3\_13. 493
- 494 Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin 495 Bhagoji, Kallista Bonawit, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael G. L. 496 D'Oliveira, Hubert Eichner, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett, 497 Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaid Harchaoui, Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi, 498 Mikhail Khodak, Jakub Konecný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo, 499 Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus 500 Pagh, Hang Qi, Daniel Ramage, Ramesh Raskar, Mariana Raykova, Dawn Song, Weikang Song, 501 Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma, 502 Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. 2021.
- 504 Christian Koetsier, Jelena Fiosina, Jan N. Gremmel, Jörg P. Müller, David M. Woisetschläger, and 505 Monika Sester. Detection of anomalous vehicle trajectories using federated learning. ISPRS 506 Open Journal of Photogrammetry and Remote Sensing, 4:100013, April 2022. ISSN 2667-3932. doi: 10.1016/j.ophoto.2022.100013. URL https://www.sciencedirect.com/ 507 science/article/pii/S2667393222000023. 508
- 509 Xiangjie Kong, Kailai Wang, Mingliang Hou, Hao Xinyu, Guojiang Shen, Chen Xin, and Feng 510 Xia. A Federated Learning-Based License Plate Recognition Scheme for 5G-Enabled Internet of 511 Vehicles. *IEEE Transactions on Industrial Informatics*, PP:1–1, March 2021. doi: 10.1109/TII. 512 2021.3067324. 513
- 514 Qingru Li, Xinru Wang, Fangwei Wang, and Changguang Wang. A Label Flipping Attack on Machine Learning Model and Its Defense Mechanism. pp. 490-506. January 2023a. ISBN 978-3-515 031-22676-2. doi: 10.1007/978-3-031-22677-9\_26. 516
- 517 Shenghui Li, Edith Ngai, Fanghua Ye, and Thiemo Voigt. Auto-weighted Robust Federated Learning 518 with Corrupted Data Sources. January 2021. doi: 10.48550/arXiv.2101.05880. 519
- 520 Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated Learning: Challenges, Methods, and Future Directions. IEEE Signal Processing Magazine, 37(3):50-60, May 2020. 521 ISSN 1558-0792. doi: 10.1109/MSP.2020.2975749. URL https://ieeexplore.ieee. 522 org/document/9084352. Conference Name: IEEE Signal Processing Magazine. 523
- 524 Xingyu Li, Zhe Qu, Shangqing Zhao, Bo Tang, Zhuo Lu, and Yao Liu. LoMar: A Lo-525 cal Defense Against Poisoning Attack on Federated Learning. IEEE Transactions 526 on Dependable and Secure Computing, 20(1):437-450, January 2023b. ISSN 1941-527 0018. 10.1109/TDSC.2021.3135422. URL https://ieeexplore.ieee. doi: 528 org/abstract/document/9650669?casa\_token=dpPrV6pRoF8AAAAA: 529 qhN3oGXR6vbJNj2xm0Z7TuAhPF1Ss0AhTvWQnU2cdFG7ew-qT8k3VvVCJQwo8XqLk7wtwb04zA. 530 Conference Name: IEEE Transactions on Dependable and Secure Computing.
- Sha Liu, Yuchuan Fu, Pincan Zhao, Fan Li, and Changle Li. Autonomous Braking Al-532 gorithm for Rear-End Collision via Communication-Efficient Federated Learning. In 2021 IEEE Global Communications Conference (GLOBECOM), pp. 01–06, December 2021. doi: 534 10.1109/GLOBECOM46510.2021.9685298. URL https://ieeexplore.ieee.org/ document/9685298. 536

531

535

Margarita Martínez-Díaz and Francesc Soriguera. Autonomous vehicles: theoretical and practical challenges. Transportation Research Procedia, 33:275-282, January 2018. ISSN 2352-1465. 538 doi: 10.1016/j.trpro.2018.10.103. URL https://www.sciencedirect.com/science/ article/pii/S2352146518302606.

551

564

565

566

567

- 540 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Ar-541 cas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In 542 Aarti Singh and Jerry Zhu (eds.), Proceedings of the 20th International Conference on Artifi-543 cial Intelligence and Statistics, volume 54 of Proceedings of Machine Learning Research, pp. 544 1273-1282. PMLR, 20-22 Apr 2017. URL https://proceedings.mlr.press/v54/ mcmahan17a.html.
- 546 Anh Nguyen, Ngoc Nguyen, Kim Tran, Erman Tjiputra, and Quang D. Tran. Autonomous Naviga-547 tion in Complex Environments with Deep Multimodal Fusion Network. pp. 5824–5830, October 548 2020. doi: 10.1109/IROS45743.2020.9341494. 549
- Ehsan Nowroozi, Nada Jadalla, Samaneh Ghelichkhani, and Alireza Jolfaei. Mitigating Label 550 Flipping Attacks in Malicious URL Detectors Using Ensemble Trees. December 2023. doi: 10.13140/RG.2.2.33453.26082. 552
- 553 Aparna P, Gandhiraj Rajendran, and Manoj Panda. Steering Angle Prediction for Autonomous Driving using Federated Learning: The Impact of Vehicle-To-Everything Communication. pp. 554 1-7, July 2021. doi: 10.1109/ICCCNT51525.2021.9580097. 555
- 556 Andrea Paudice, Luis Muñoz-González, and Emil Lupu. Label Sanitization Against Label Flipping Poisoning Attacks. pp. 5–15. February 2019. ISBN 978-3-030-13452-5. doi: 10.1007/ 558 978-3-030-13453-2\_1.
- 559 Adnan Qayyum, Muhammad Umar Janjua, and Junaid Qadir. Making federated learning robust to 560 adversarial attacks by learning data and model association. Computers & Security, 121:102827, 561 October 2022. ISSN 0167-4048. doi: 10.1016/j.cose.2022.102827. URL https://www. sciencedirect.com/science/article/pii/S0167404822002218. 563
  - M V Rajasekhar and Anil Kumar Jaswal. Autonomous vehicles: The future of automobiles. In 2015 IEEE International Transportation Electrification Conference (ITEC), pp. 1–6, August 2015. doi: 10.1109/ITEC-India.2015.7386874. URL https://ieeexplore.ieee.org/ document/7386874.
- 568 Elan Rosenfeld, Ezra Winston, Pradeep Ravikumar, and J. Zico Kolter. Certified robustness to labelflipping attacks via randomized smoothing. In Proceedings of the 37th International Conference 569 on Machine Learning, volume 119 of ICML'20, pp. 8230-8241. JMLR.org, 2020. 570
- 571 SunFounder Picar-X Video Robot Car Kit for Raspberry Pi 5/4/3B+/3B, Sunfounder. 572 Python/Blockly (Scratch), Video Courses, Rechargeable Batterry (Raspberry Pi NOT Included). 573 URL https://www.sunfounder.com/products/picar-x.
- 574 Vale Tolpegin, Stacey Truex, Mehmet Emre Gursoy, and Ling Liu. Data Poisoning Attacks Against 575 Federated Learning Systems. volume 12308, pp. 480–501, Cham, 2020. Springer International 576 Publishing. ISBN 978-3-030-58950-9 978-3-030-58951-6. doi: 10.1007/978-3-030-58951-6\_ 577 24. URL http://link.springer.com/10.1007/978-3-030-58951-6\_24. Book 578 Title: Computer Security - ESORICS 2020 Series Title: Lecture Notes in Computer Science. 579
- Tianhao Wu, Mingzhi Jiang, Yinhui Han, Zheng Yuan, Xinhang Li, and Lin Zhang. A Traffic-580 Aware Federated Imitation Learning Framework for Motion Control at Unsignalized Intersections 581 with Internet of Vehicles. *Electronics*, 10(24):3050, January 2021. ISSN 2079-9292. doi: 10. 582 3390/electronics10243050. URL https://www.mdpi.com/2079-9292/10/24/3050. 583 Number: 24 Publisher: Multidisciplinary Digital Publishing Institute. 584
- Renyou Xie, Chaojie Li, Xiaojun Zhou, and Zhaoyang Dong. Asynchronous Federated Learn-585 ing for Real-Time Multiple Licence Plate Recognition Through Semantic Communication. In 586 ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1-5, June 2023. doi: 10.1109/ICASSP49357.2023.10097251. URL 588 https://ieeexplore.ieee.org/document/10097251. ISSN: 2379-190X. 589
- Jean-Paul A. Yaacoub, Hassan N. Noura, and Ola Salman. Security of federated learning with IoT systems: Issues, limitations, challenges, and solutions. Internet of Things and Cyber-Physical Systems, 3:155–179, January 2023. ISSN 2667-3452. doi: 10.1016/j.iotcps. 592 2023.04.001. URL https://www.sciencedirect.com/science/article/pii/ S2667345223000226.

- Waleed Yamany, Nour Moustafa, and Benjamin Turnbull. OQFL: An Optimized Quantum-Based Federated Learning Framework for Defending Against Adversarial Attacks in Intelligent Transportation Systems. *IEEE Transactions on Intelligent Transportation Systems*, 24(1): 893–903, January 2023. ISSN 1558-0016. doi: 10.1109/TITS.2021.3130906. URL https: //ieeexplore.ieee.org/document/9641742. Conference Name: IEEE Transactions on Intelligent Transportation Systems.
- Haomiao Yang, Mengyu Ge, Dongyun Xue, Kunlan Xiang, Hongwei Li, and Rongxing Lu. Gradient Leakage Attacks in Federated Learning: Research Frontiers, Taxonomy, and Future Directions. *IEEE Network*, 38(2):247–254, March 2024. ISSN 1558-156X. doi: 10.1109/MNET.001. 2300140. URL https://ieeexplore.ieee.org/document/10107713. Conference Name: IEEE Network.
- Haiyang Yu, Rui Jiang, Zhengbing He, Zuduo Zheng, Li Li, Runkun Liu, and Xiqun Chen. Automated vehicle-involved traffic flow studies: A survey of assumptions, models, speculations, and perspectives. *Transportation Research Part C: Emerging Technologies*, 127:103101, June 2021.
  ISSN 0968-090X. doi: 10.1016/j.trc.2021.103101. URL https://www.sciencedirect. com/science/article/pii/S0968090X21001224.
- Tengchan Zeng, Omid Semiari, Mingzhe Chen, Walid Saad, and Mehdi Bennis. Federated Learn ing on the Road Autonomous Controller Design for Connected and Autonomous Vehicles.
   *IEEE Transactions on Wireless Communications*, 21(12):10407–10423, December 2022. ISSN 1558-2248. doi: 10.1109/TWC.2022.3183996. URL https://ieeexplore.ieee.org/
   document/9806308/?arnumber=9806308. Conference Name: IEEE Transactions on Wireless Communications.
- Hongyi Zhang, Jan Bosch, and Helena Olsson. End-to-End Federated Learning for Autonomous
   Driving Vehicles. pp. 1–8, July 2021. doi: 10.1109/IJCNN52387.2021.9533808.
- Jingwei Zhang, Jost Springenberg, Joschka Boedecker, and Wolfram Burgard. Deep Reinforcement Learning with Successor Features for Navigation across Similar Environments. December 2016. doi: 10.48550/arXiv.1612.05533.
- Hongliang Zhou, Yifeng Zheng, Hejiao Huang, Jiangang Shu, and Xiaohua Jia. Toward Robust Hier archical Federated Learning in Internet of Vehicles. *IEEE Transactions on Intelligent Transporta- tion Systems*, 24(5):5600–5614, May 2023. ISSN 1558-0016. doi: 10.1109/TITS.2023.3243003.
   URL https://ieeexplore.ieee.org/abstract/document/10046398?casa\_
   token=qsufVdnlUg8AAAAA:aE1Gx\_C2WmFp4raGvJriYt605HDzbPrmIPZS\_
   K-JkW-KZmsbdi\_lLHDaeM8HqOfQId8QvLiYEQ. Conference Name: IEEE Transactions
- on Intelligent Transportation Systems.
- Kuehan Zhou, Ruimin Ke, Zhiyong Cui, Qiang Liu, and Wenxing Qian. STFL:Spatio-temporal Fed erated Learning for Vehicle Trajectory Prediction. In 2022 IEEE 2nd International Conference on Digital Twins and Parallel Intelligence (DTPI), pp. 1–6, October 2022. doi: 10.1109/DTPI55838.
   2022.9998967. URL https://ieeexplore.ieee.org/document/9998967.
- Gongxi Zhu, Donghao Li, Hanlin Gu, Yuxing Han, Yuan Yao, Lixin Fan, and Qiang Yang.
   Evaluating Membership Inference Attacks and Defenses in Federated Learning. 2024. doi: 10.48550/ARXIV.2402.06289. URL https://arxiv.org/abs/2402.06289. Publisher: arXiv Version Number: 1.
- Haomin Zhuang, Mingxian Yu, Hao Wang, Yang Hua, Jian Li, and Xu Yuan. Backdoor Federated
  Learning by Poisoning Backdoor-Critical Layers, April 2024. URL http://arxiv.org/
  abs/2308.04466. arXiv:2308.04466 [cs].
- 642

619

- 643
- 644
- 645
- 646
- 647