

# 000 001 002 003 004 005 FEDTAP: TRUST-AWARE TEMPORAL AGGREGATION 006 FOR ROBUST FEDERATED LEARNING 007

008 NODI Lab., Jilin University  
009

## 010 ABSTRACT 011

012 Federated learning (FL) enables multiple clients to collaboratively train a shared  
013 global model without exposing private data, but its distributed nature makes it vul-  
014 nerable to adversarial and unreliable updates. Existing robust aggregation methods  
015 focus on single-round defenses and often fail against adaptive or stealthy attacks  
016 that evolve gradually over time. In this work, we propose **FedTAP** (*Federated*  
017 *Trust-aware Adaptive Prediction and Aggregation*), a temporal modeling frame-  
018 work that enhances robustness by explicitly incorporating the temporal dynamics  
019 of global model updates. FedTAP treats the server as an observer that predicts the  
020 expected benign evolution of the global model and measures deviations as temporal  
021 residuals. These residuals are combined with spatial consistency cues to compute  
022 a unified anomaly score, which drives a temporal trust propagation mechanism  
023 that accumulates each client's credibility across rounds. The resulting trust values  
024 are then used to assign adaptive aggregation weights, allowing reliable clients to  
025 maintain influence while gradually suppressing persistent adversaries.  
026

## 027 1 INTRODUCTION 028

029 Federated learning (FL) has become an essential paradigm for distributed machine learning, where a  
030 central server coordinates multiple clients to train a shared global model without directly accessing  
031 their local data (McMahan et al., 2017). This design protects data privacy and enables large-scale  
032 learning across heterogeneous devices and organizations. However, the distributed nature of FL  
033 also introduces security and reliability challenges. Some clients may behave maliciously or produce  
034 unreliable updates due to adversarial intent, hardware failure, or data corruption. These abnormal  
035 updates can significantly degrade the global model, cause convergence instability, or even implant  
036 hidden backdoors (Bagdasaryan et al., 2020; Wang et al., 2020). Since the server can only observe  
037 model updates and has no visibility into the local data, detecting such behaviors remains a difficult  
038 and open problem.

039 To mitigate these vulnerabilities, many robust aggregation algorithms have been proposed. Early  
040 works such as Krum (Blanchard et al., 2017) and the coordinate-wise median or trimmed mean  
041 estimators (Yin et al., 2018) reduce the influence of corrupted updates by selecting or reweighting  
042 client gradients based on spatial consistency. Later, the geometric-median aggregation (RFA) (Pillutla  
043 et al., 2022) improves stability under heterogeneous data distributions, and divergence-based methods  
044 such as the  $\gamma$ -mean estimator (Li et al., 2022) achieve robustness against heavy-tailed or adversarial  
045 updates. More recent studies employ Huber-loss based aggregation to balance robustness and conver-  
046 gence under non-IID data (Zhao et al., 2024). Although these methods achieve strong performance  
047 against Byzantine attacks, they operate at the level of individual communication rounds and do not  
048 consider temporal information across rounds. As a result, they remain vulnerable to adaptive and  
049 temporally coordinated attacks that make small but consistent changes over time.

050 Several studies have investigated poisoning and backdoor attacks that exploit this temporal weakness.  
051 Attackers can gradually modify model parameters in multiple rounds while keeping each individual  
052 update close to the normal distribution (Wang et al., 2020). This strategy allows the malicious  
053 clients to remain undetected while progressively biasing the global model. To defend against such  
054 attacks, approaches such as FLTrust (Cao et al., 2020) use a small trusted dataset at the server for  
055 calibration, and clustering-based detection methods such as MUDGuard (Wang et al., 2024) identify  
056 anomalies through spatial analysis of client gradients. However, these defenses treat each training  
057 round independently and cannot capture long-term behavioral patterns. Recent methods attempt to

use historical information, such as FoolsGold, which detects sybil attacks by analyzing the similarity of client updates over time (Fung et al., 2018), and recursive aggregation methods that adjust client contributions based on previous model deviations (Herath et al., 2023). Yet, these approaches still rely on short-term memory and lack an explicit model for the temporal evolution of trust or the prediction of the global trajectory.

In this work, we propose **FedTAP** (*Federated Trust-aware Adaptive Prediction and Aggregation*), a temporal modeling framework for robust federated learning. We view the server as an *observer* that continuously monitors the global model dynamics. FedTAP first builds a lightweight predictor that estimates the expected benign evolution of the global model. Deviations from this forecast define temporal residuals that measure inconsistency between the observed and predicted behaviors. To complement this, a spatial subspace consistency check captures deviations in the geometry of client updates. The two signals are combined into a unified multi-scale anomaly score that drives a *temporal trust propagation* mechanism. This mechanism accumulates each client's credibility over time, allowing the system to recognize persistent anomalies while avoiding overreaction to temporary noise or statistical fluctuations. Finally, the accumulated trust scores are mapped to adaptive weights for the aggregation step, which enables smooth and interpretable adjustment of each client's influence.

FedTAP advances robust federated learning by integrating temporal prediction, trust propagation, and adaptive aggregation within a single framework. It has three main advantages. First, temporal prediction captures long-term consistency of the global model and detects slow, stealthy attacks. Second, trust propagation maintains stability and prevents drastic changes caused by transient noise. Third, the soft weighting scheme ensures compatibility with existing aggregation rules without modifying the standard communication process. Experimental results on benchmark datasets such as MNIST, CIFAR-10, and Tiny-ImageNet show that FedTAP significantly improves resilience against both fast and slow poisoning attacks while maintaining high model accuracy and stable convergence.

## 2 RELATED WORK

Federated learning allows multiple clients to train a shared model collaboratively without exchanging their private data (McMahan et al., 2017). This decentralized paradigm protects data privacy and improves scalability but is highly sensitive to malicious or unreliable clients. Robust aggregation methods aim to mitigate this problem by designing aggregation schemes that can tolerate a certain proportion of corrupted updates. Krum selects the update that is most consistent with others in Euclidean distance, offering formal Byzantine resilience guarantees (Blanchard et al., 2017). Coordinate-wise median and trimmed mean estimators further improve robustness by achieving optimal statistical error bounds under bounded adversarial perturbations (Yin et al., 2018). Pillutla et al. introduce the geometric-median based RFA method, which shows strong empirical robustness across non-IID datasets (Pillutla et al., 2022). Subsequent research explores the  $\gamma$ -mean aggregation rule, which uses minimum  $\gamma$ -divergence estimation to enhance stability under heavy-tailed and adversarial noise (Li et al., 2022). More recent work adopts Huber-loss based aggregation to balance robustness and accuracy under heterogeneous data conditions (Zhao et al., 2024). Although these algorithms perform well for single-round aggregation, they generally overlook the temporal evolution of model parameters across rounds, which is crucial for detecting slow or adaptive attacks.

Beyond aggregation design, many studies investigate poisoning and backdoor attacks in FL. Bagdasaryan et al. demonstrate that a small number of adversarial clients can perform model replacement to embed backdoors while maintaining high accuracy on clean data (Bagdasaryan et al., 2020). Wang et al. show that coordinated multi-round manipulations can produce stealthy and persistent backdoor effects (Wang et al., 2020). Surveys by Nguyen et al. (Nguyen et al., 2024) and Manzoor et al. (Manzoor et al., 2024) summarize diverse poisoning and backdoor strategies as well as their countermeasures. On the defensive side, FLTrust leverages a small trusted dataset at the server to calibrate client updates and mitigate poisoning risk (Cao et al., 2020). Clustering-based anomaly detection methods, such as MUDGuard, identify malicious clients by analyzing the spatial similarity of gradient directions (Wang et al., 2024). These approaches, however, focus on static aggregation within each round and do not explicitly model how malicious behaviors evolve temporally.

Several recent studies begin to integrate temporal or reputational information into the defense process. FoolsGold tracks the similarity of client updates across rounds to detect sybil-based poisoning attacks (Fung et al., 2018). Herath et al. propose recursive Euclidean distance based aggregation that adjusts

108 the contribution of each client using their historical deviation from previous global models (Herath  
 109 et al., 2023). Although these methods utilize history, they still rely on local pairwise consistency or  
 110 short-term correlations rather than forecasting the global model trajectory or formalizing long-term  
 111 trust. Building on these observations, we propose **FedTAP**, which explicitly treats the server as an  
 112 observer of a dynamic system. It forecasts the expected benign evolution of the global model and  
 113 accumulates temporal trust based on multi-round consistency.

### 115 3 PROPOSED DESIGN OF FEDTAP

117 Federated learning aggregates client updates to evolve a global parameter  $g_t \in \mathbb{R}^d$ :

$$119 \quad g_{t+1} = \mathcal{F}_{\text{base}}(g_t, \{\Delta_{t,i}\}_{i \in S_t}), \quad (1)$$

120 where  $S_t \subseteq [N]$  is the set of participating clients and  $\mathcal{F}_{\text{base}}$  is an arbitrary base aggregator (FedAvg,  
 121 Trimmed-Mean, Bulyan, etc.). We consider *temporally stealthy* adversaries that control a subset  $\mathcal{M}$   
 122 with  $|\mathcal{M}|/N < 0.5$ . These adversaries inject small but coordinated changes in each round, which  
 123 gradually alter the global model while keeping single-round statistics almost unchanged.  
 124

125 We view the server as an *observer* of a dynamical system and monitor its *temporal consistency*. The  
 126 proposed method has two main parts: (i) a light-weight predictor that estimates the expected benign  
 127 evolution of the global model, and (ii) a *temporal trust propagation* mechanism that aggregates  
 128 evidence across rounds. This design allows the server to separate persistent adversaries from normal  
 129 clients and prevents overreaction to random noise or single-round outliers.

#### 130 3.1 OBSERVER-BASED PREDICTION AND TEMPORAL RESIDUALS

132 In each communication round  $t$ , the server observes the current global model  $g_t$  and a batch of client  
 133 updates  $\{\Delta_{t,i}\}_{i \in S_t}$ . Under normal conditions, the global parameters change smoothly because both  
 134 optimization noise and data heterogeneity vary gradually. We use this temporal smoothness to predict  
 135 where the global model is expected to move next. A client update that moves the model away from  
 136 this predicted direction is likely to be abnormal or malicious.

138 **Modeling benign dynamics.** We model the normal evolution of the global model as

$$140 \quad g_{t+1} = F(g_t, \{\Delta_{t,i}\}_{i \in S_t}) + \xi_t, \quad (2)$$

141 where  $F(\cdot)$  represents the aggregation mapping and  $\xi_t$  accounts for randomness caused by client  
 142 sampling or mini-batch optimization. Because  $F$  is unknown and may be nonlinear, the server learns  
 143 a data-driven approximation of its dynamics. We use a **linear autoregressive (AR)** predictor:

$$145 \quad \hat{g}_{t+1} = A_0 g_t + A_1 g_{t-1} + \cdots + A_{K-1} g_{t-K+1}, \quad (3)$$

146 where each  $A_k$  is a small coefficient matrix that captures short-term relationships between previous  
 147 and future global parameters. The predictor learns how the model tends to evolve during the recent  
 148  $W$  rounds. To estimate  $\{A_k\}$ , we minimize the mean squared prediction error:  
 149

$$150 \quad \min_{\{A_k\}} \sum_{u=t-W+1}^t \|g_{u+1} - \sum_{k=0}^{K-1} A_k g_{u-k}\|_2^2 + \rho \sum_k \|A_k\|_F^2, \quad (4)$$

153 where the small ridge term  $\rho$  (typically  $10^{-4}$ ) improves stability and prevents overfitting. The  
 154 predictor can be updated each round using **recursive least squares (RLS)** with a forgetting factor of  
 155 0.95, which requires only  $O(Wd_{\text{eff}})$  operations per layer. Empirical results show that  $K \in [2, 4]$  and  
 156  $W \in [10, 30]$  are sufficient to capture the evolution of federated learning dynamics.

158 **Candidate global and residual computation.** Once the server obtains the prediction  $\hat{g}_{t+1}$ , it  
 159 evaluates each client's update separately. For client  $i$ , the server imagines what the global model  
 160 would become if only this client's update were applied:

$$161 \quad g_{t+1}^{(i)} = g_t + \eta \Delta_{t,i}, \quad (5)$$

162 where  $\eta$  is the server step size (or 1 if  $\Delta_{t,i}$  already includes the step). The virtual model  $g_{t+1}^{(i)}$   
 163 represents the trajectory that client  $i$  alone would produce. The difference between the predicted  
 164 benign model  $\hat{g}_{t+1}$  and the client-induced model  $g_{t+1}^{(i)}$  defines the *temporal residual*:  
 165

$$r_{t,i} = \|g_{t+1}^{(i)} - \hat{g}_{t+1}\|_2. \quad (6)$$

166 A large residual indicates that the client update would move the model in a direction that is inconsistent  
 167 with recent global behavior. To ensure balance across network layers,  $r_{t,i}$  is computed for each layer  
 168 and then normalized before summation.  
 169

170 **Takeaways 3.1.** *The predictor can be updated online without storing the full training history. The*  
 171 *server maintains a rolling buffer of  $W$  past global models and solves equation 4 using standard RLS*  
 172 *updates. In practice, this requires only a few matrix–vector multiplications per round. The matrices*  
 173  *$A_k$  can also be restricted to block-diagonal form so that each network layer evolves independently,*  
 174 *which reduces computation and improves stability.*

### 176 3.2 MULTI-SCALE DETECTION AND TEMPORAL TRUST

177 Single-round detection often fails when attackers behave slowly or adaptively. To address this  
 178 problem, we combine two complementary consistency checks: (i) a *spatial cue* that measures how far  
 179 a client update lies from the main subspace spanned by recent benign updates, and (ii) a *temporal cue*  
 180 that tracks how each client’s behavior evolves across rounds.  
 181

182 **Subspace consistency (spatial cue).** During normal training, most client updates lie within a  
 183 low-dimensional subspace because of the shared model architecture and similar data distributions.  
 184 Let  $\{\bar{\Delta}_u\}_{u=t-W+1}^t$  denote the aggregated updates from the most recent  $W$  rounds. We maintain  
 185 an incremental Principal Component Analysis (PCA) to estimate this dominant subspace  $\mathcal{S}$ . The  
 186 subspace dimension  $d_s$  is selected to retain 90–95% of the total variance. The deviation of client  $i$ ’s  
 187 update from this subspace is measured as

$$s_{t,i} = \|\Delta_{t,i} - \mathbf{P}_{\mathcal{S}}(\Delta_{t,i})\|_2, \quad (7)$$

188 where  $\mathbf{P}_{\mathcal{S}}(\cdot)$  is the projection onto  $\mathcal{S}$ . A large  $s_{t,i}$  indicates that the client update moves in a direction  
 189 that is not well aligned with recent benign updates, which may suggest manipulation. The PCA basis  
 190 is updated online each round and requires  $O(d_s^2 d_{\text{eff}})$  operations.  
 191

192 Because the scale of residuals can vary across rounds, we apply robust normalization using the  
 193 median and the median absolute deviation (MAD), which are resistant to outliers:  
 194

$$\tilde{r}_{t,i} = \frac{r_{t,i} - \text{Med}(r_t)}{1.4826 \text{ MAD}(r_t)}, \quad \tilde{s}_{t,i} = \frac{s_{t,i} - \text{Med}(s_t)}{1.4826 \text{ MAD}(s_t)}. \quad (8)$$

195 Here  $r_{t,i}$  is the temporal residual defined in Eq. equation 6. After normalization, typical benign  
 196 clients have values near 0–1, while clear outliers often exceed 3–4. We combine both cues into a  
 197 single *anomaly score*:

$$z_{t,i} = \alpha \tilde{r}_{t,i} + \beta \tilde{s}_{t,i}, \quad \text{where } \alpha, \beta \geq 0, \text{ and } \alpha=\beta=1 \text{ by default.} \quad (9)$$

201 This score balances spatial and temporal information for reliable detection. We convert this anomaly  
 202 score into an *instantaneous credibility*  $c_{t,i} \in (0, 1)$  through a smooth logistic mapping:  
 203

$$c_{t,i} = \frac{1}{1 + \exp((z_{t,i} - \tau_t)/\kappa_t)}, \quad (10)$$

204 where  $\tau_t$  is an adaptive threshold and  $\kappa_t$  controls the sharpness of the transition. A value of  $c_{t,i}$  close  
 205 to 1 means the update appears normal, while smaller values indicate that the update is less reliable.  
 206 We set  $\kappa_t = 1.4826 \text{ MAD}(\{z_{t,i}\})$  so that the mapping automatically adapts to the score dispersion  
 207 in each round.

208 **Temporal trust propagation (temporal cue).** Instantaneous credibility alone can be unreliable  
 209 because a benign client may appear abnormal in a single noisy round. To improve long-term  
 210 robustness, we maintain a latent *trust variable*  $\theta_{t,i} \in [0, 1]$  for each client. This variable integrates the  
 211 client’s credibility over time:

$$\theta_{t+1,i} = \beta_{\theta} \theta_{t,i} + (1 - \beta_{\theta}) c_{t,i}, \quad (11)$$

216 where  $\beta_\theta \in [0.8, 0.99]$  determines the memory length (default  $\beta_\theta=0.9$ ). If a client remains reliable  
 217 ( $c_{t,i} \approx 1$ ), its trust gradually returns to 1. If it stays suspicious for several rounds,  $\theta_{t,i}$  decays exponentially  
 218 and marks the client as low trust. This process smooths short-term noise and captures  
 219 persistent misbehavior. We also apply a mild hysteresis rule: when  $\theta_{t,i}$  is below 0.3, small one-time  
 220 improvements do not immediately restore full trust.

221 Client behaviors and gradient magnitudes may change over time. Therefore, the decision threshold  $\tau_t$   
 222 should adjust smoothly. We update it using robust exponential smoothing:

$$224 \quad \tau_{t+1} = \gamma_\tau \tau_t + (1 - \gamma_\tau) \text{Med}(\{z_{t,i}\}) + c_\tau \text{MAD}(\{z_{t,i}\}), \quad (12)$$

225 where  $\gamma_\tau \in [0.8, 0.95]$  defines the memory length and  $c_\tau \in [2, 4]$  controls sensitivity. This update  
 226 allows the threshold to follow gradual distribution shifts while ignoring short-term noise.

227 **Takeaways 3.2.** *Each round provides an instantaneous credibility  $c_{t,i}$  based on current statistics,  
 228 and a trust variable  $\theta_{t,i}$  that aggregates credibility over time. Benign clients maintain high trust,  
 229 while adversarial clients accumulate low trust and are gradually downweighted in the aggregation.  
 230 Together, this multi-scale mechanism identifies both short-term anomalies and slow, stealthy attacks.*

### 232 3.3 TRUST-AWARE ADAPTIVE AGGREGATION AND UPDATING

234 Once the server obtains the latent trust scores  $\{\theta_{t,i}\}$ , it must incorporate them into the aggregation  
 235 process in a smooth and interpretable way. Instead of using hard rejection or strict thresholds, we  
 236 translate trust into *soft aggregation weights*. This approach allows reliable clients to have stronger  
 237 influence, while clients that behave abnormally are gradually suppressed. The design also provides  
 238 temporal memory, so that the model does not overreact to noise in a single round.

239 **Trust-guided aggregation.** We convert each client's trust value  $\theta_{t,i} \in [0, 1]$  into a positive aggrega-  
 240 tion weight:

$$242 \quad w_{t,i} = \exp(-\lambda(1 - \theta_{t,i})), \quad (13)$$

243 where  $\lambda > 0$  controls how strongly low-trust clients are downweighted (typical  $\lambda \in [4, 6]$ ). When  
 244  $\theta_{t,i} = 1$ , the client is fully trusted ( $w_{t,i} = 1$ ); when  $\theta_{t,i} = 0.5$ , the weight decreases to  $e^{-\lambda/2}$ . Since  
 245  $\theta_{t,i}$  evolves smoothly according to Eq. equation 11, this mapping makes the aggregation stable across  
 246 rounds. A client with a temporary anomaly will not be severely penalized, but persistent abnormal  
 247 behavior will lead to exponential suppression.

248 The adaptive weights are then integrated into the base aggregator:

$$250 \quad g_{t+1} = \mathcal{F}_{\text{base}}\left(g_t, \{w_{t,i} \cdot \Delta_{t,i}\}_{i \in S_t}\right), \quad (14)$$

252 where  $\mathcal{F}_{\text{base}}$  can represent FEDAVG, Trimmed-Mean, Bulyan, or other robust methods. This structure  
 253 is plug-and-play: the proposed trust mechanism acts as a *reweighting layer* that can be added to any  
 254 existing aggregation framework. For long-term stability, a client is temporarily isolated if its trust  
 255 remains below a minimum level  $\theta_{\min}$  for  $L$  consecutive rounds. We set  $\theta_{\min} = 0.3$  and  $L = 3$  by  
 256 default. Such clients are re-evaluated after a cooldown period to avoid permanent exclusion.

257 **Stability and analysis.** To analyze the mechanism, assume benign residuals follow a sub-Gaussian  
 258 distribution with variance proxy  $\sigma^2$ , and that an adversary causes a mean shift  $\delta > 0$  in the anomaly  
 259 score  $z_{t,i}$  for  $L$  consecutive rounds. Then the expected credibility decreases by  $\Omega(\delta/\kappa_t)$ . By  
 260 repeatedly applying Eq. equation 11, the trust variable evolves as

$$262 \quad \theta_{t+L,i} \leq \beta_\theta^L \theta_{t,i} + (1 - \beta_\theta^L) \left(1 - \Omega\left(\frac{\delta}{\kappa_t}\right)\right), \quad (15)$$

264 which shows that persistent anomalies cause an exponential decay in trust. Under the standard  
 265 Lipschitz continuity of  $\mathcal{F}_{\text{base}}$ , the expected model deviation satisfies

$$266 \quad \mathbb{E}[\|g_{t+1} - g^*\|] \leq \rho \mathbb{E}[\|g_t - g^*\|] + \epsilon, \quad (16)$$

268 where  $\rho < 1$  because adversarial updates receive very small weights as  $\theta_{t,i}$  decreases. This guarantees  
 269 the stability and convergence of the global model even when bounded adversarial perturbations are  
 present.

270       **Takeaways 3.3.** This stage transforms trust into quantitative influence. Clients with consistent  
271       behavior keep high trust and near-unit weights, while long-term or stealthy attackers experience  
272       exponential trust decay, which limits their effect on the global model. This time-weighted aggregation  
273       extends standard robust aggregators with a temporal perspective, combining statistical robustness  
274       with adaptive control.  
275

## 276       REFERENCES 277

- 278       Eugene Bagdasaryan, Alexander Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to  
279       backdoor federated learning. In *Proceedings of the 23rd International Conference on Artificial  
280       Intelligence and Statistics (AISTATS)*, pp. 2938–2948, 2020.  
281
- 282       Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Byzantine-tolerant  
283       machine learning. *arXiv preprint arXiv:1703.02757*, 2017.  
284
- 285       Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. Fltrust: Byzantine-robust federated  
286       learning via trust bootstrapping. In *Proceedings of the Network and Distributed System Security  
Symposium (NDSS)*, 2020.  
287
- 288       Clement Fung, Chih-Jen Yoon, and Ivan Beschastnikh. Mitigating sybils in federated learning  
289       poisoning. In *arXiv preprint arXiv:1808.04866*, 2018.  
290
- 291       Charuka Herath, Yogachandran Rahulamathan, and Xiaolan Liu. Recursive euclidean distance  
292       based robust aggregation technique for federated learning. In *Proceedings of 2023 IEEE IAS  
Global Conference on Emerging Technologies (GlobConET)*, pp. 1–6, 2023.  
293
- 294       Cen-Jhih Li, Pin-Han Huang, Yi-Ting Ma, Hung Hung, and Su-Yun Huang. Robust aggregation for  
295       federated learning by minimum  $\gamma$ -divergence estimation. *Entropy*, pp. 686, 2022.  
296
- 297       Habib Ullah Manzoor, Attia Shabbir, Ao Chen, David Flynn, and Ahmed Zoha. A survey of security  
298       strategies in federated learning: Defending models, data, and privacy. *Future Internet*, 16:374,  
2024.  
299
- 300       Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.  
301       Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the  
International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 1273–1282, 2017.  
302
- 303       Thuy Dung Nguyen, Tuan Nguyen, Phi Le Nguyen, Hieu H. Pham, Khoa D. Doan, and Kok-Seng  
304       Wong. Backdoor attacks and defenses in federated learning: A survey, challenges and future  
305       directions. *Engineering Applications of Artificial Intelligence*, 127:107166, 2024.  
306
- 307       Krishna Pillutla, Sham M. Kakade, and Zaid Harchaoui. Robust aggregation for federated learning.  
308       *IEEE Transactions on Signal Processing*, 70:1142–1154, 2022.  
309
- 310       Hongyi Wang, Kartik Sreenivasan, Shashank Rajput, Harit Vishwakarma, Saurabh Agarwal, Jy yong  
311       Sohn, Kangwook Lee, and Dimitris Papailiopoulos. Attack of the tails: Yes, you really can  
312       backdoor federated learning. In *Proceedings of Advances in Neural Information Processing  
Systems (NeurIPS)*, pp. 16070–16084, 2020.  
313
- 314       Rui Wang, Xingkai Wang, Huanhuan Chen, Jérémie Decouchant, Stjepan Picek, Nikolaos Laoutaris,  
315       and Kaitai Liang. Mudguard: Taming malicious majorities in federated learning using privacy-  
316       preserving byzantine-robust clustering. *Proceedings of the ACM on Measurement and Analysis of  
Computing Systems (POMACS)*, pp. 1–41, 2024.  
317
- 318       Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-robust distributed  
319       learning: Towards optimal statistical rates. In *Proceedings of the 35th International Conference on  
Machine Learning (ICML)*, volume 80, pp. 5650–5659, 2018.  
320
- 321       Puning Zhao, Fei Yu, and Zhiguo Wan. A huber loss minimization approach to byzantine robust  
322       federated learning. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp.  
323       21806–21814, 2024.