

























































FL Directed Deviation Attack

Fang Attack [USENIX Sec '20]: Computes a direction vector along the inverse of the average benign direction.

- Krum Attack: This attack ensures all malicious models are close to each other with small mutual distance, fooling the Krum aggregator to choose the poisoned model.

- Trim Attack: This attack samples model updates per parameter in a way that skews the distribution toward the malicious direction.

Shejwalkar Attack [NDSS '21]: Computes a perturbation (inverse unit) vector and scales it up before adding to the benign updates. The scaling factor is tuned depending on the dataset and the model used.











Macro Results

Table 1: Impact of Directed Deviation Model Poisoning Attacks: This table presents the test accuracy for directed deviation model poisoning attacks (Full-Krum; Full-Trim) on various datasets with a c/m ratio of 0.2. For the Shakespeare dataset, we report the test loss instead. The results highlight the damaging impact of Full-Trim attacks on mean-like aggregations (FedSGD, Trimmed mean, Median) and Full-Krum attacks on Krum-like aggregations (Krum, Bulyan). While existing defenses such as FABA, FoolsGold, and FLTrust show mixed results, our proposed method, FLAIR, consistently outperforms in all cases.

Attack	Defense	Metrics				
		MNIST+ DNN	CIFAR-10+ ResNet-18	Shakespeare+ GRU	FEMNIST+ DNN	
None	FedSGD	92.45	71.17	1.62	83.60	
	FLAIR	92.52	66.92	1.64	83.58	
	FABA	91.77	69.94	1.76	82.69	
	FoolsGold	91.20	70.71	1.63	83.80	
	FLTrust	87.70	68.08	1.62	82.72	
Full-Krum	FedSGD	82.97	39.68	1.62	29.87	
	Krum	8.92	9.81	11.98	5.62	
	Bulyan	10.14	13.24	9.23	9.91	
	FLAIR	87.73	61.26	1.64	80.19	
	FABA	86.99	55.96	1.75	55.61	
	FoolsGold	47.12	42.28	1.63	0.07	
	FLTrust	82.50	65.25	1.67	79.53	
Full-Trim	FedSGD	65.25	47.32	1.74	32.34	
	Trim	36.36	55.25	3.28	13.03	
	Median	28.37	50.54	3.30	45.6	
	FLAIR	90.55	67.65	1.66	82.51	
	FABA	91.84	67.31	1.64	79.66	
	FoolsGold	91.61	69.24	1.66	83.09	
	FLTrust	34.20	64.23	1.68	79.28	

Macro Results

Table 3: Comparison of test accuracies for FLAIR under SHE-JWALKAR attack, with and without the knowledge of the aggregator. The performance of FLAIR is compared against the baseline FedSGD model on two datasets: MNIST and CIFAR-10. The table illustrates FLAIR's robustness in the face of attacks and its ability to maintain high accuracy rates.

AGR	AGR-knowledge	MNIST	CIFAR-10
FedSGD	No	10.10	10.00
	Yes	10.09	10.00
FLAIR	No	92.25	69.35
	Yes	92.83	69.98









































