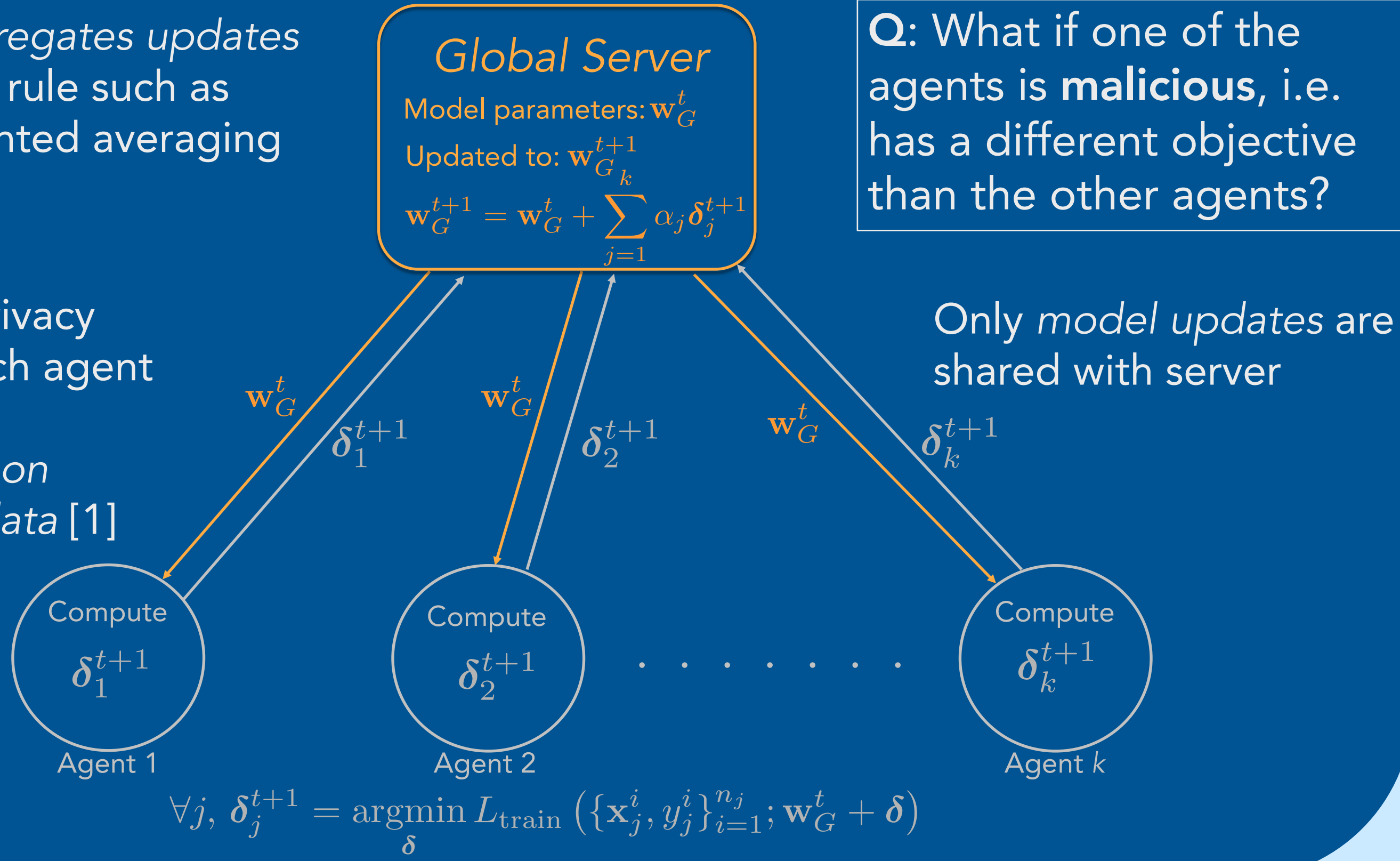


Federated Learning

Server aggregates updates based on a rule such as linear weighted averaging (shown)

Guided by privacy concerns, each agent performs computation on locally held data [1]

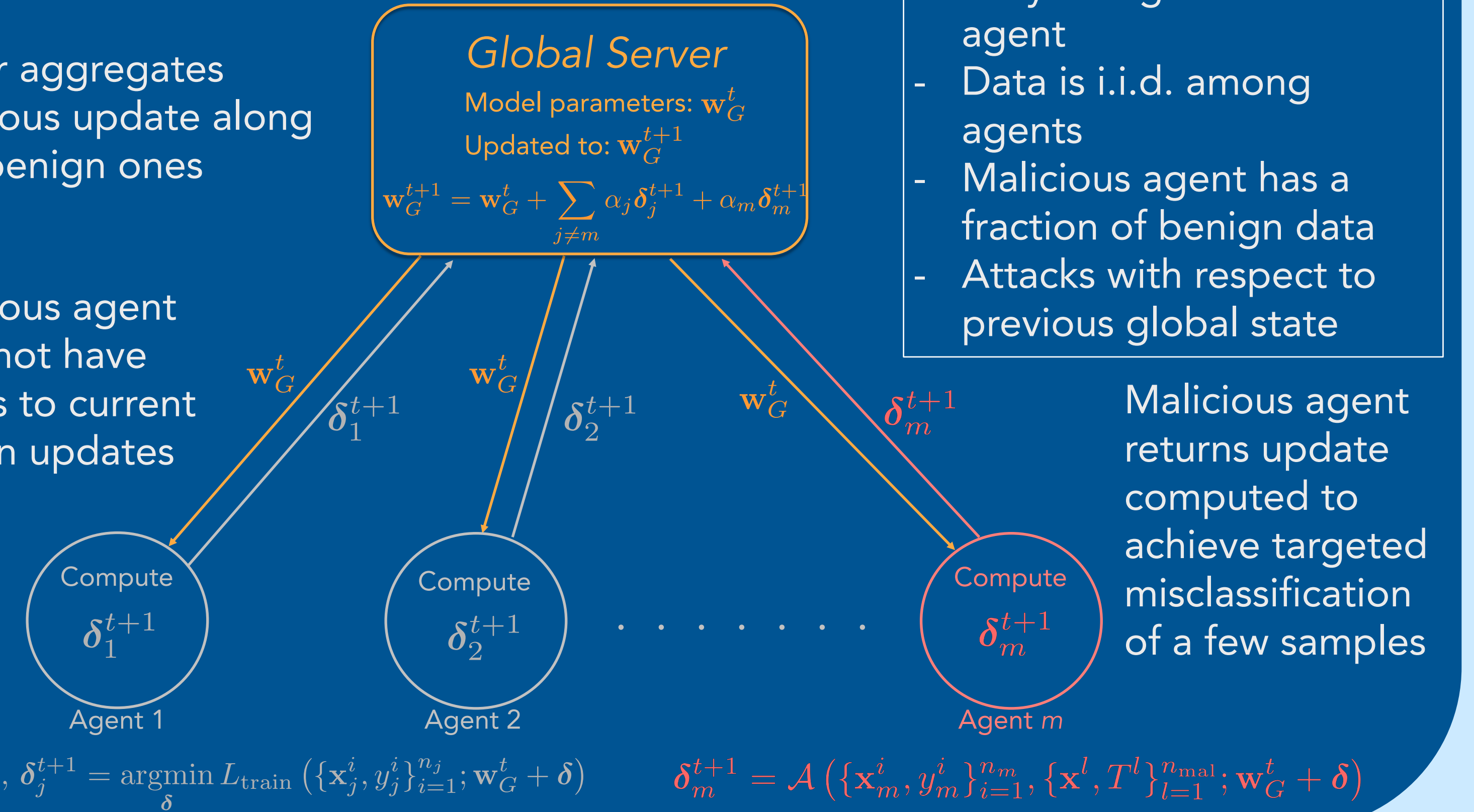


$$\forall j, \delta_j^{t+1} = \argmin_{\delta} L_{\text{train}}(\{x_j^i, y_j^i\}_{i=1}^{n_j}; w_G^t + \delta)$$

Threat Model

Server aggregates malicious update along with benign ones

Malicious agent does not have access to current benign updates

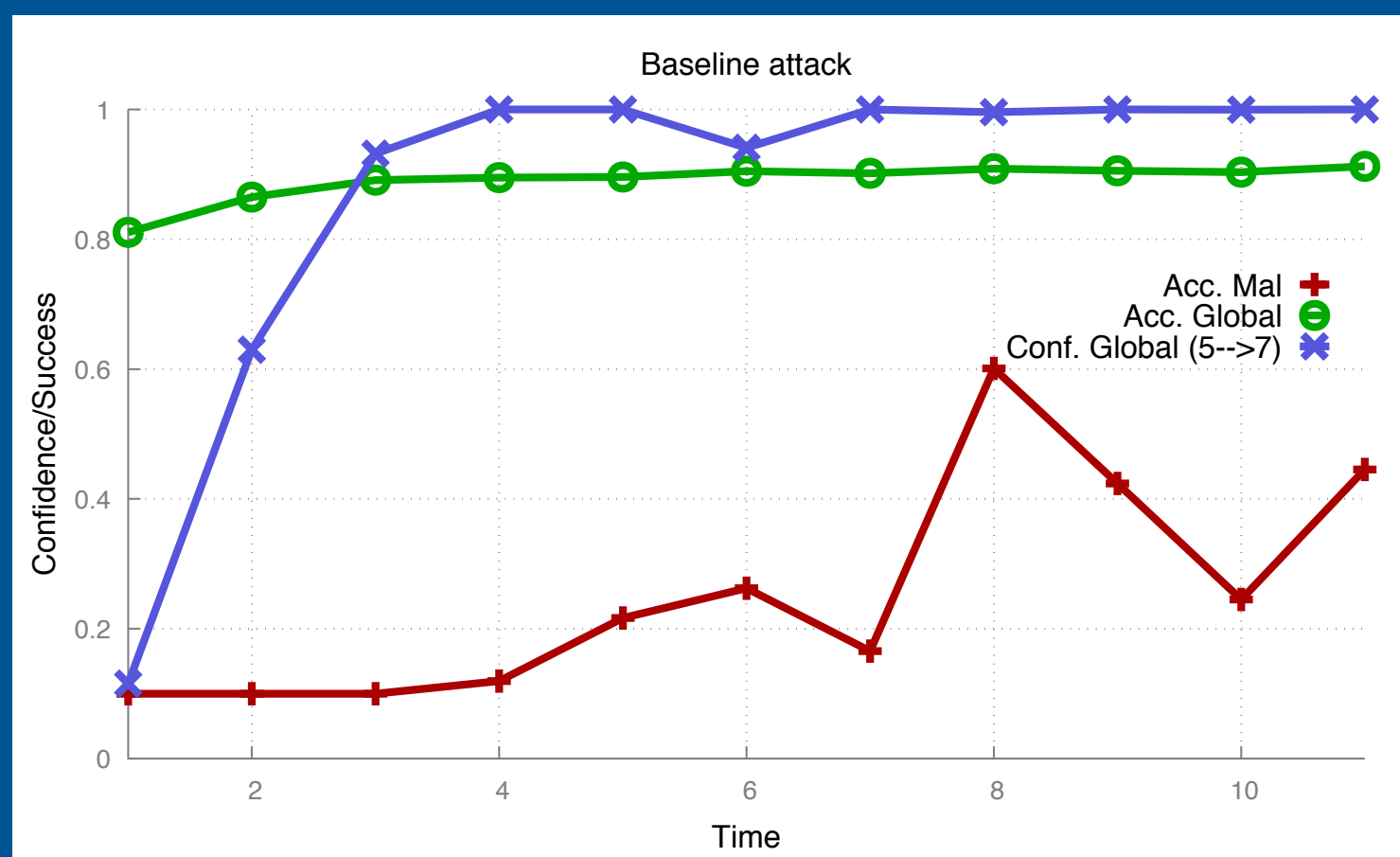


$$\forall j \neq m, \delta_j^{t+1} = \argmin_{\delta} L_{\text{train}}(\{x_j^i, y_j^i\}_{i=1}^{n_j}; w_G^t + \delta) \quad \delta_m^{t+1} = \mathcal{A}(\{x_m^i, y_m^i\}_{i=1}^{n_m}, \{x^l, T^l\}_{l=1}^{n_{\text{mal}}}; w_G^t + \delta)$$

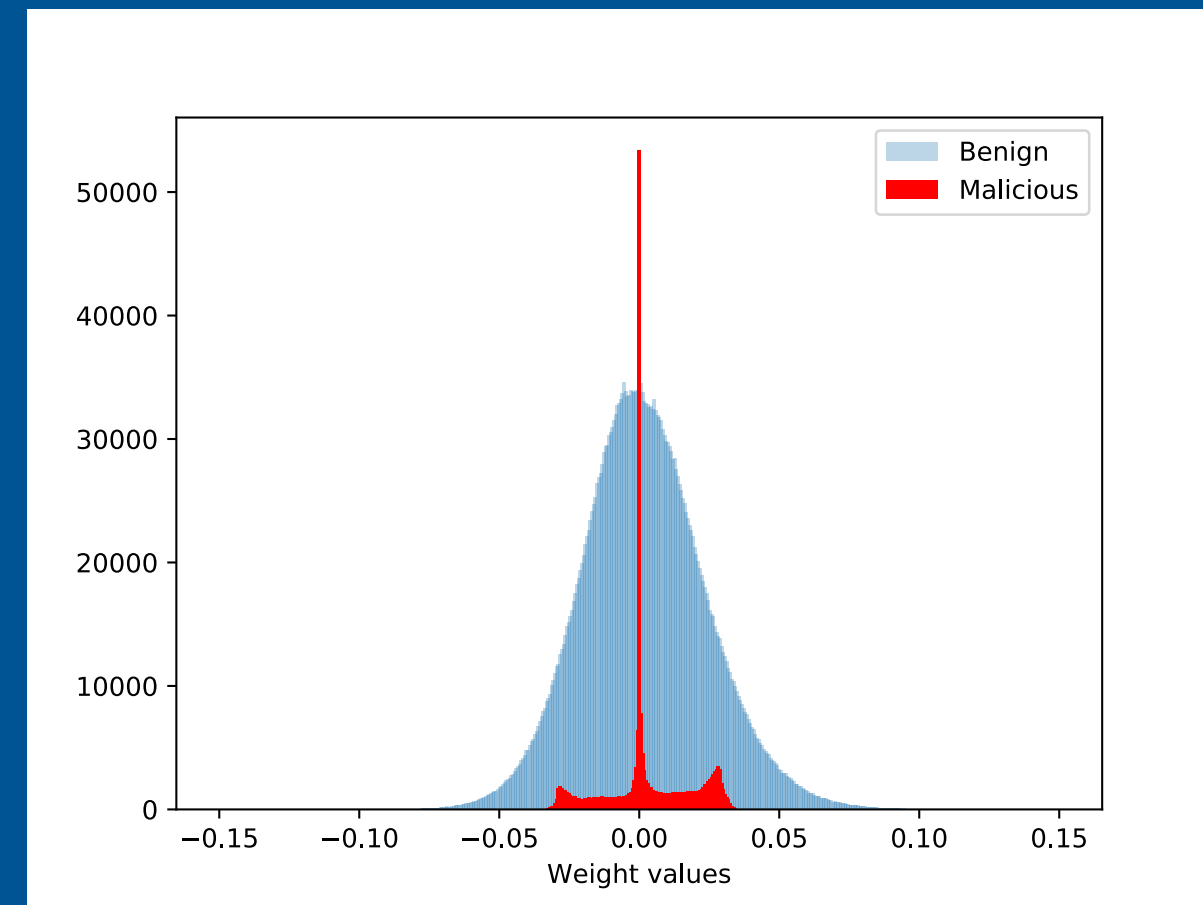
Dataset: Fashion MNIST [2]
Model: CNN with 91.7% test set accuracy

Attack Strategies for Model Poisoning

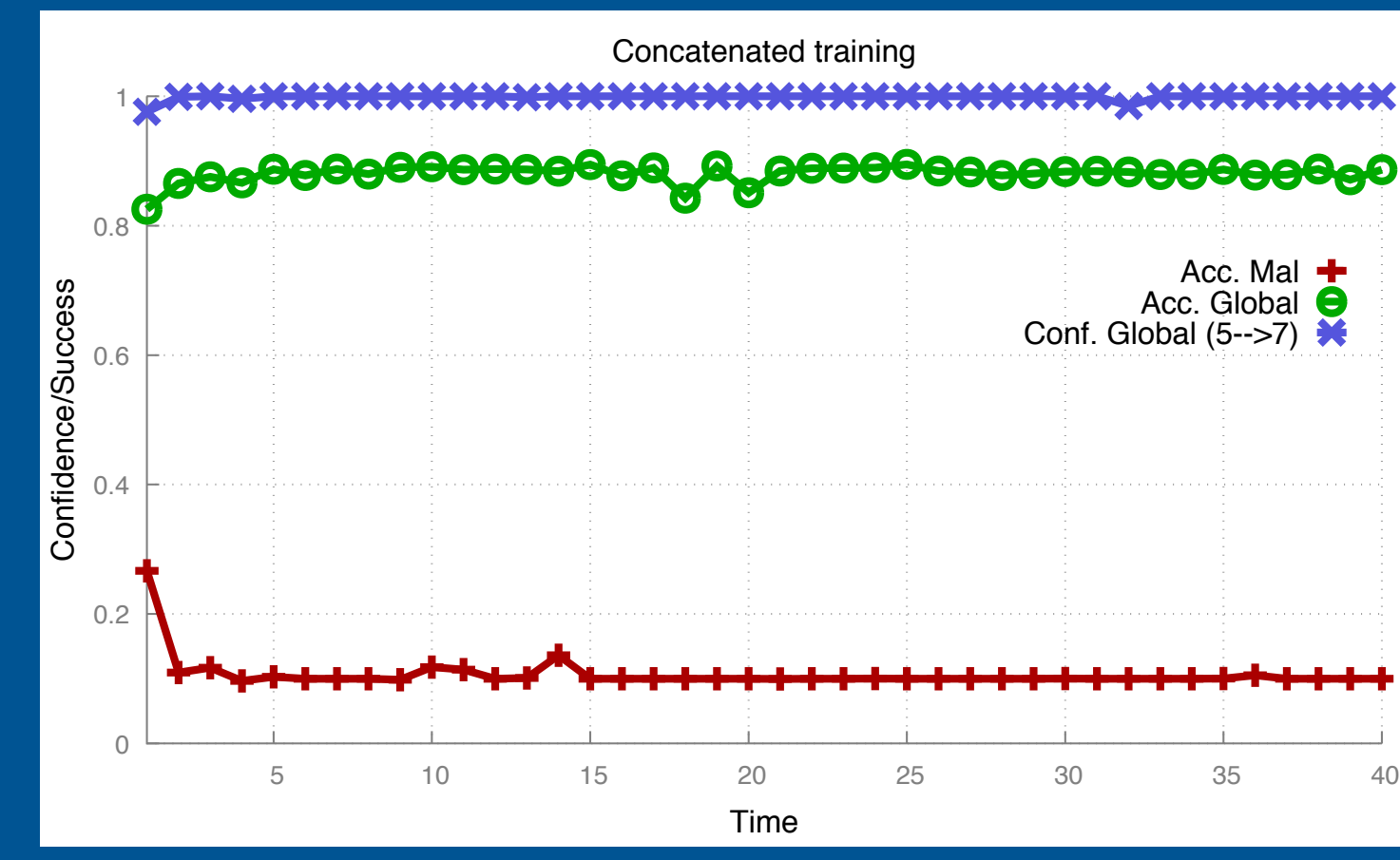
Baseline attack (explicit boosting)



$$\delta_{\text{mal}} = \argmin_{\delta} L_{\text{mal}}(\{x^l, T^l\}_{l=1}^{n_{\text{mal}}}; w_G + \delta) \quad \delta_{\text{mal}} \rightarrow \beta \delta_{\text{mal}}$$

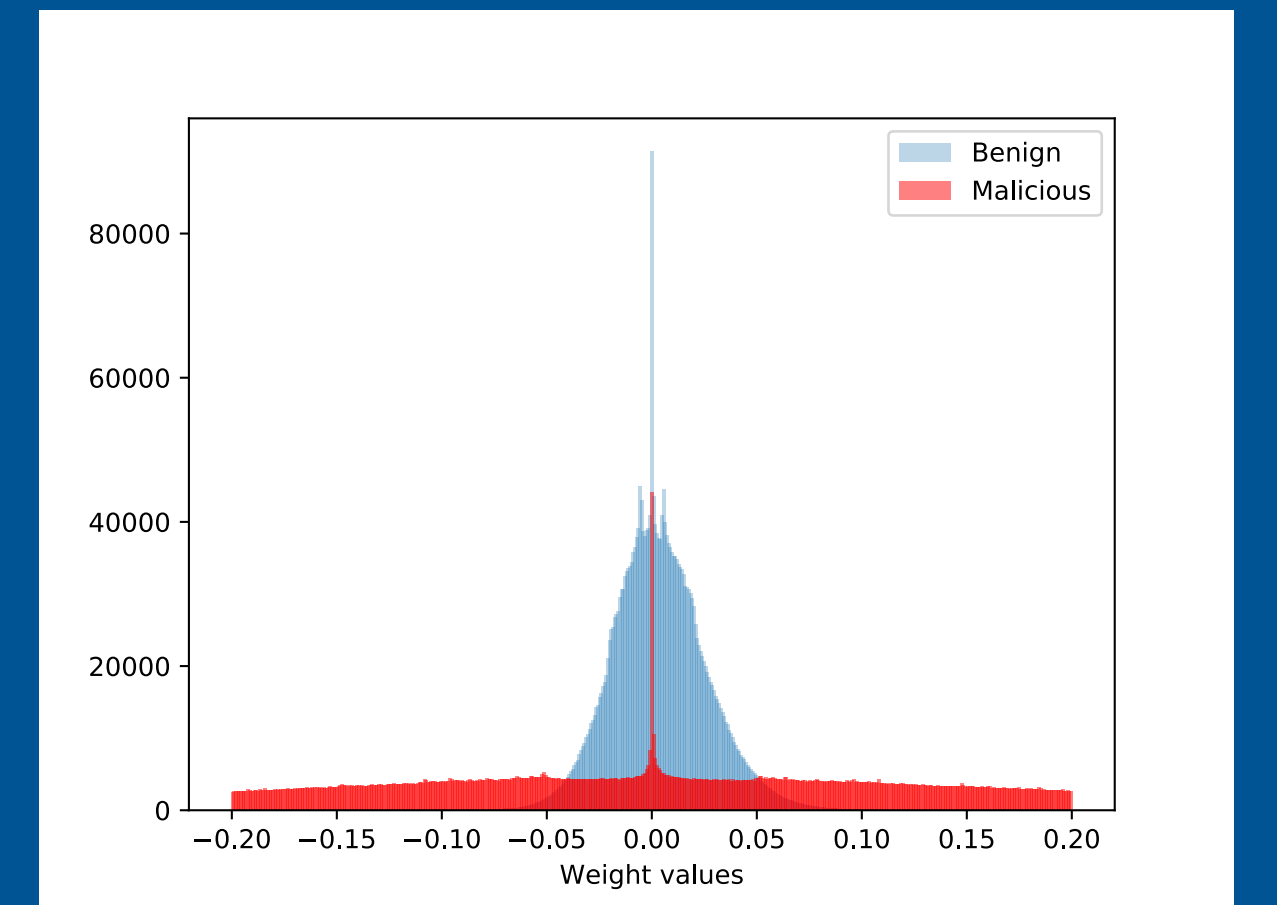


$$\delta_{\text{mal}} = \argmin_{\delta} L_{\text{mal}}(\{x^l, T^l\}_{l=1}^{n_{\text{mal}}}, \{x_m^i, y_m^i\}_{i=1}^{n_m}; w_G + \delta) \quad \delta_{\text{mal}} \rightarrow \beta \delta_{\text{mal}}$$

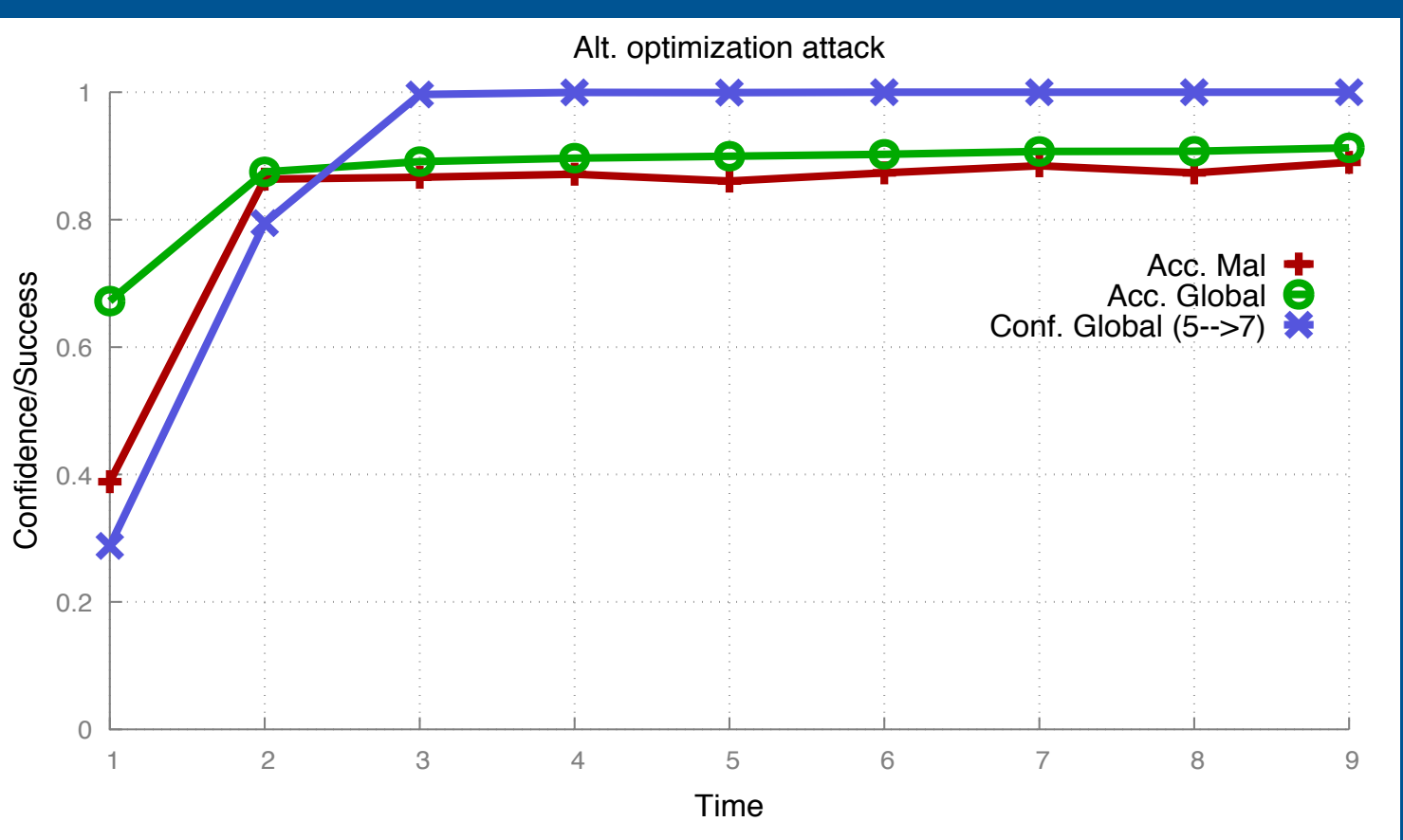


Adversarial goal is to ensure (sandal, class 5) is classified as a sneaker (class 7)

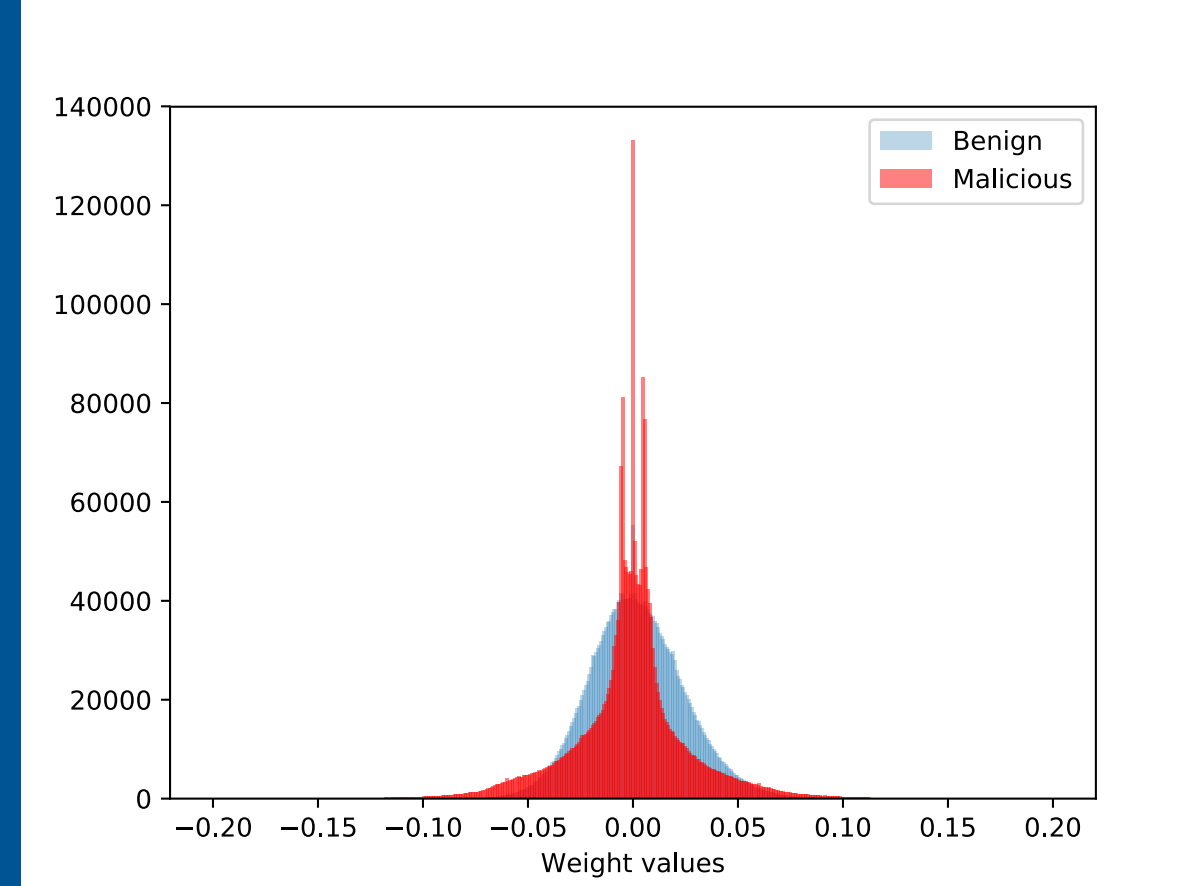
Concatenated training



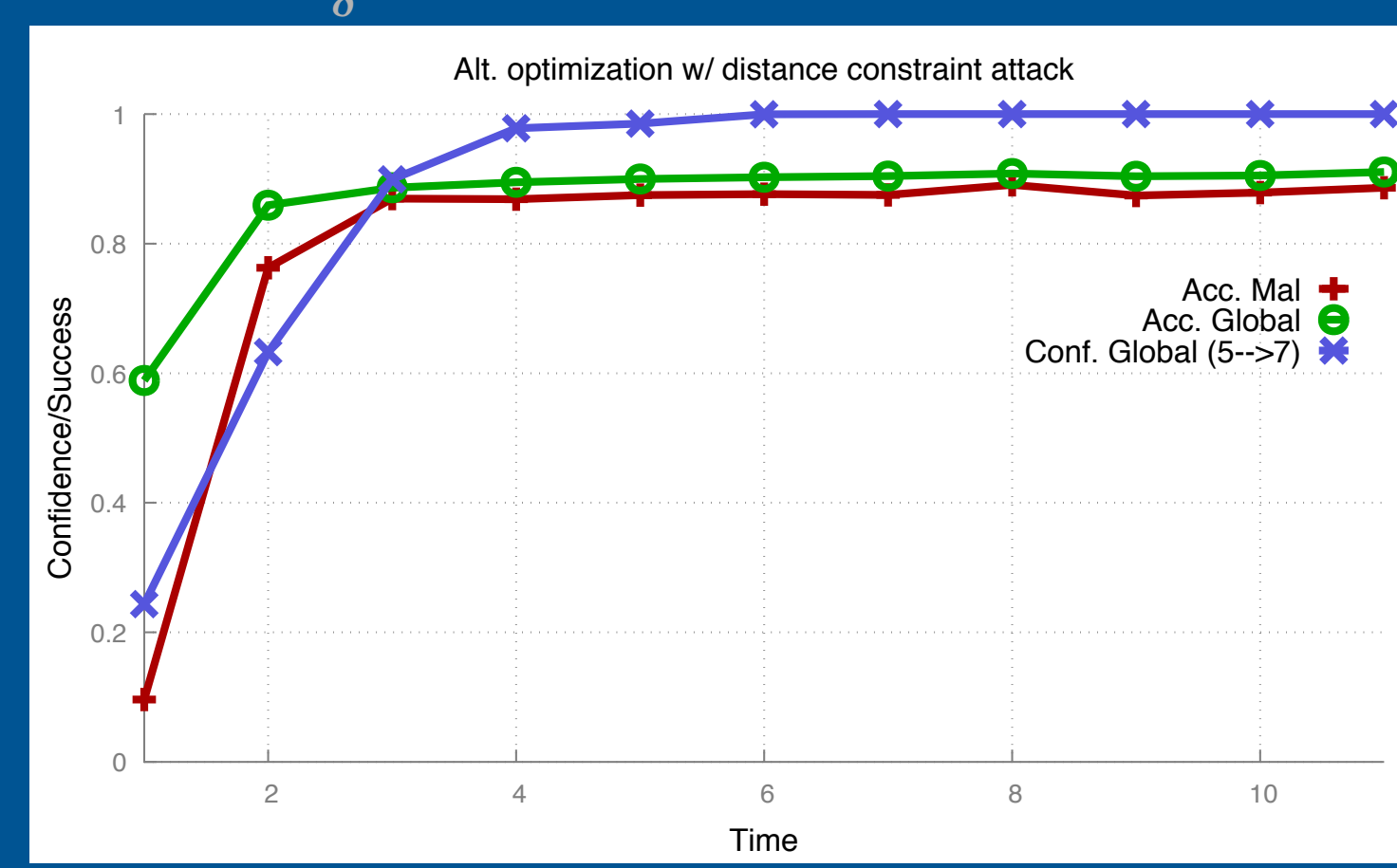
Alternating minimization (switching between benign and malicious updates for stealth)



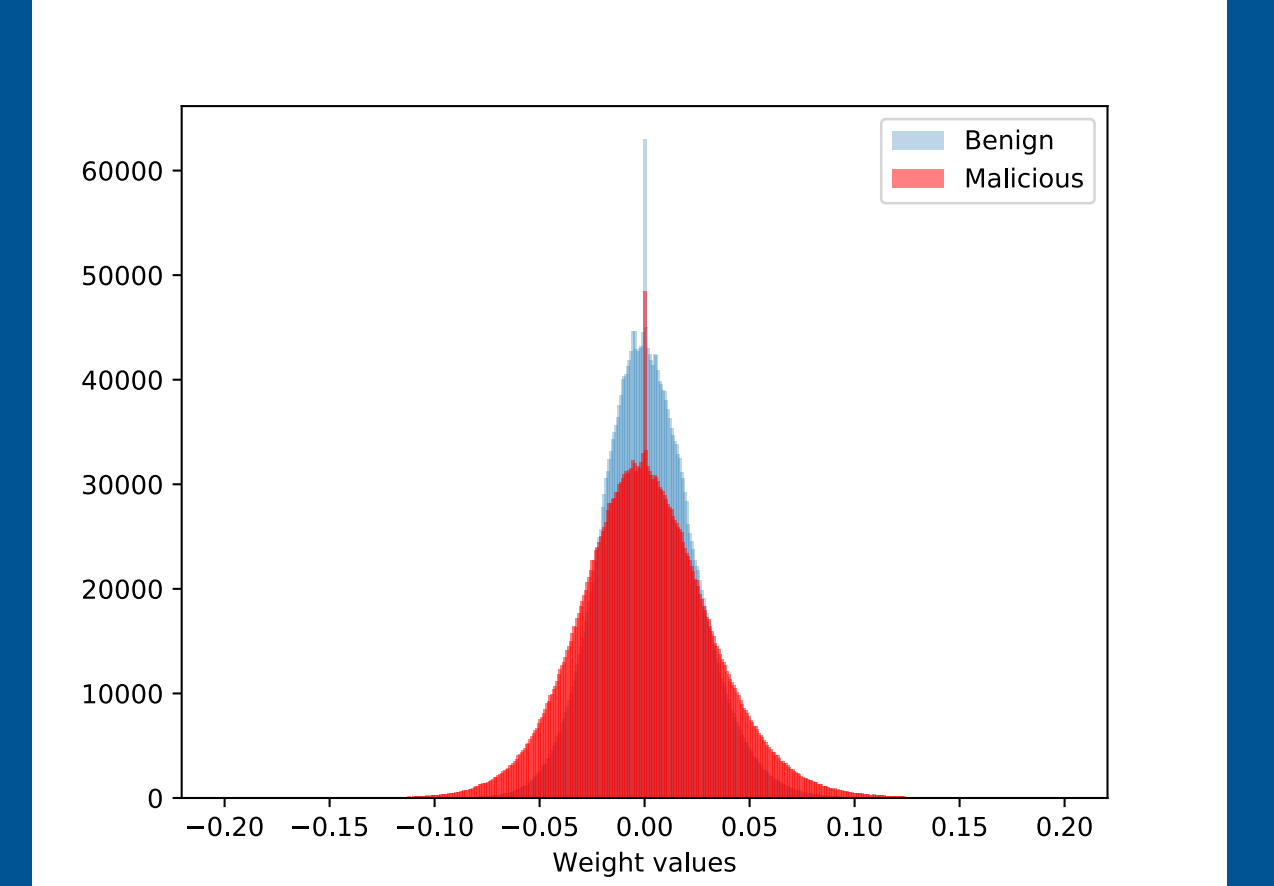
$$\begin{aligned} \delta'_{\text{mal}} &= \argmin_{\delta} L_{\text{mal}}(\{x^l, T^l\}_{l=1}^{n_{\text{mal}}}; w_G + \delta) \\ \text{Repeat: } \delta'_{\text{mal}} &\rightarrow \beta \delta'_{\text{mal}} \\ \delta''_{\text{mal}} &= \argmin_{\delta} L_{\text{ben}}(\{x_m^i, y_m^i\}_{i=1}^{n_m}; w_G + \beta \delta'_{\text{mal}} + \delta) \end{aligned}$$



$$\begin{aligned} \delta'_{\text{mal}} &= \argmin_{\delta} L_{\text{mal}}(\{x^l, T^l\}_{l=1}^{n_{\text{mal}}}; w_G + \delta) \\ \text{Repeat: } \delta'_{\text{mal}} &\rightarrow \beta \delta'_{\text{mal}} \\ \delta''_{\text{mal}} &= \argmin_{\delta} L_{\text{ben}}(\{x_m^i, y_m^i\}_{i=1}^{n_m}; w_G + \beta \delta'_{\text{mal}} + \delta) + \rho \|\delta - \delta_{\text{cons}}\|_2^2 \end{aligned}$$

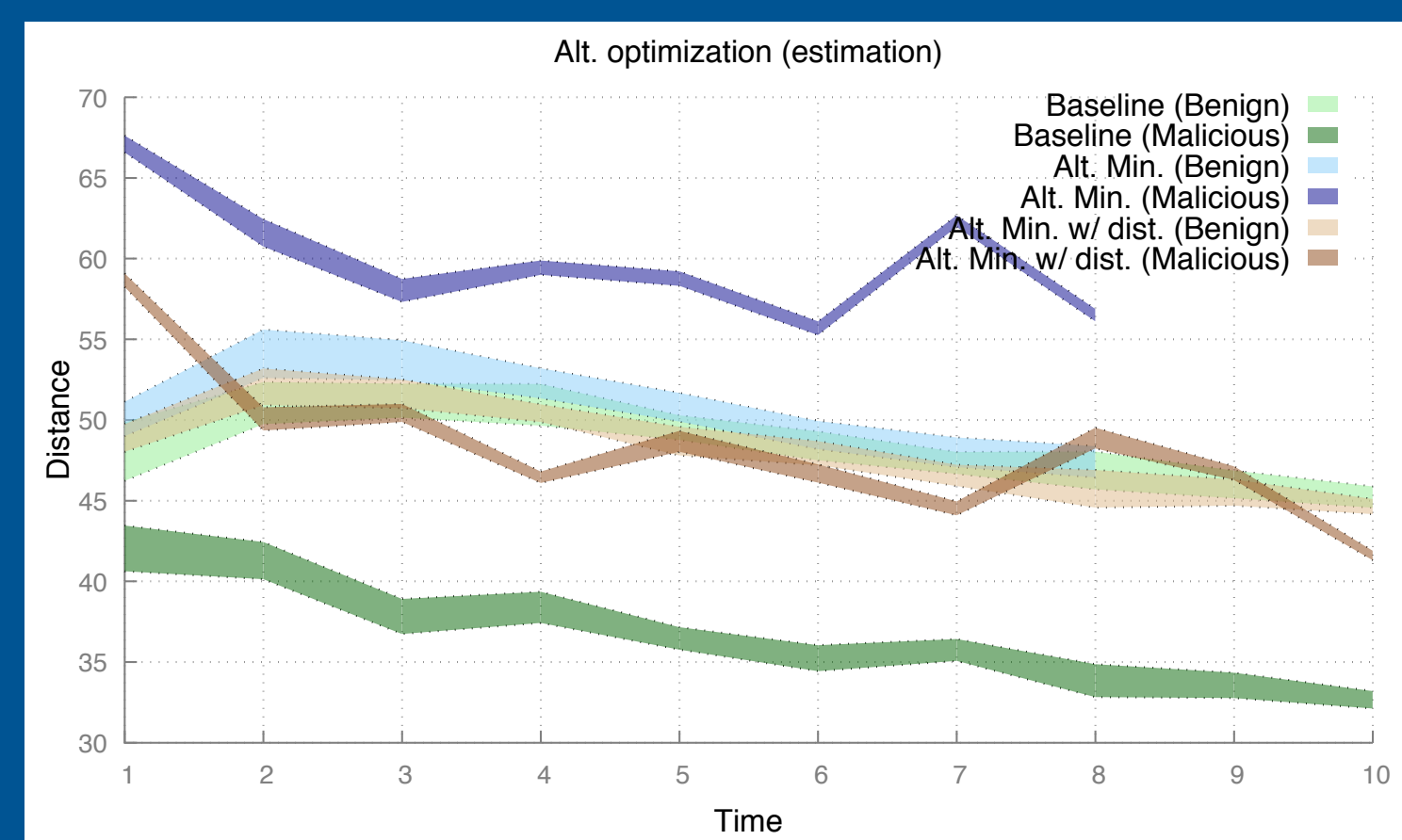


Alternating minimization with distance constraints



Attack stealth measure: distance spread

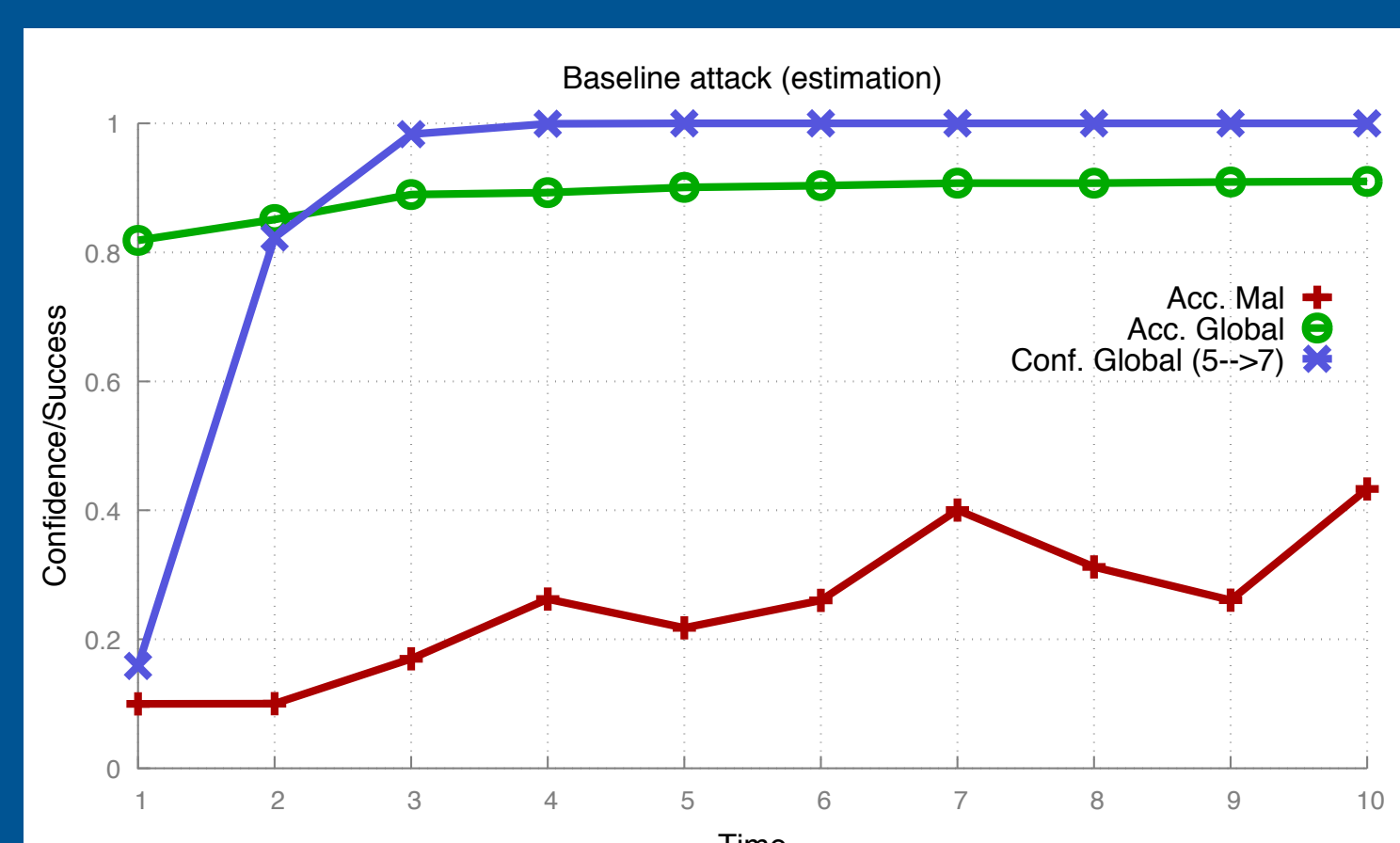
For each strategy, we show the spread of L_2 distances between all the benign agents and between the malicious agent and the benign agents.



Estimation to improve attacks

Pre-optimization correction with previous step estimate of benign agents' effects

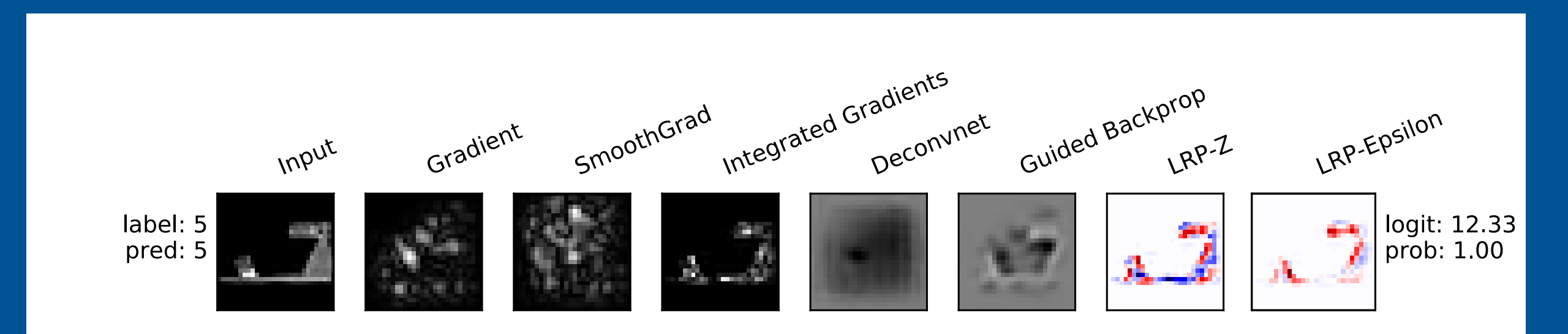
$$\begin{aligned} \hat{w}_G^t &= \hat{w}_G^{t-1} + \hat{\delta}_{[k] \setminus m} + \alpha_m \delta_m^t \\ \hat{\delta}_{[k] \setminus m} &= \delta_{[k] \setminus m}^{t-1} \end{aligned}$$



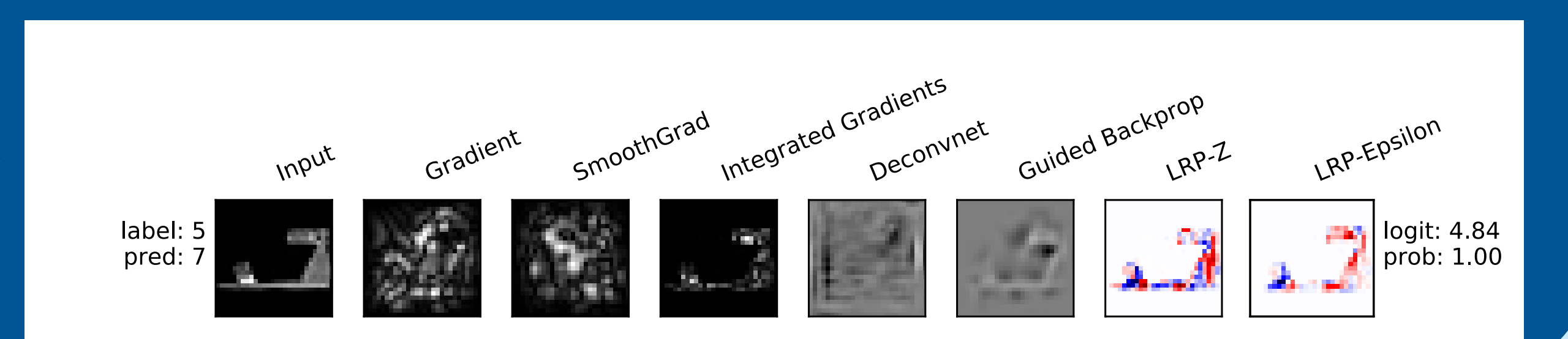
Interpreting Poisoned Models

- Interpretability techniques provide insights into the internal feature representations and working of a neural network
- Used a suite of these techniques [3] to understand decisions of poisoned models

Global model trained using only benign agents



Global model trained with one malicious model and the rest benign



Conclusion

Our attacks in this paper demonstrate that federated learning in its basic form is very vulnerable to model poisoning adversaries. While detection mechanisms can make these attacks more challenging, these can be overcome, demonstrating that multi-party machine learning algorithms robust to attackers of the type considered here must be developed.

References

- [1] McMahan et al., *Communication-Efficient Learning of Deep Networks from Decentralized Data*, AISTATS 2017
- [2] Xiao et al., *Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms*, arXiv preprint arXiv:1708.07747, 2017
- [3] Alber et al., *iNNvestigate neural networks!*, arXiv preprint arXiv:1808.04260, 2018