# CRFL: Certifiably Robust Federated Learning against Backdoor Attacks

## ICML 2021

Chulin Xie, Minghao Chen, Pin-Yu Chen, Bo Li

chulinx2@illinois.edu, pin-yu.chen@ibm.com, lbo@illinois.edu

# Motivation

## Backdoor Attack against Federated Learning (FL)



Standard FL Training process

Sever

Clients

Correct prediction    Correct prediction

*test*

$$z_j^i := \{x_j^i, y_j^i\}$$

i-th Client's j-th data

Malicious FL Training process

Sever

Clients

Correct prediction    Target wrong prediction

*test*

$$z'^i_j := \{x_j^i + \delta_{ix}, y_j^i + \delta_{iy}\}$$

Backdoored version
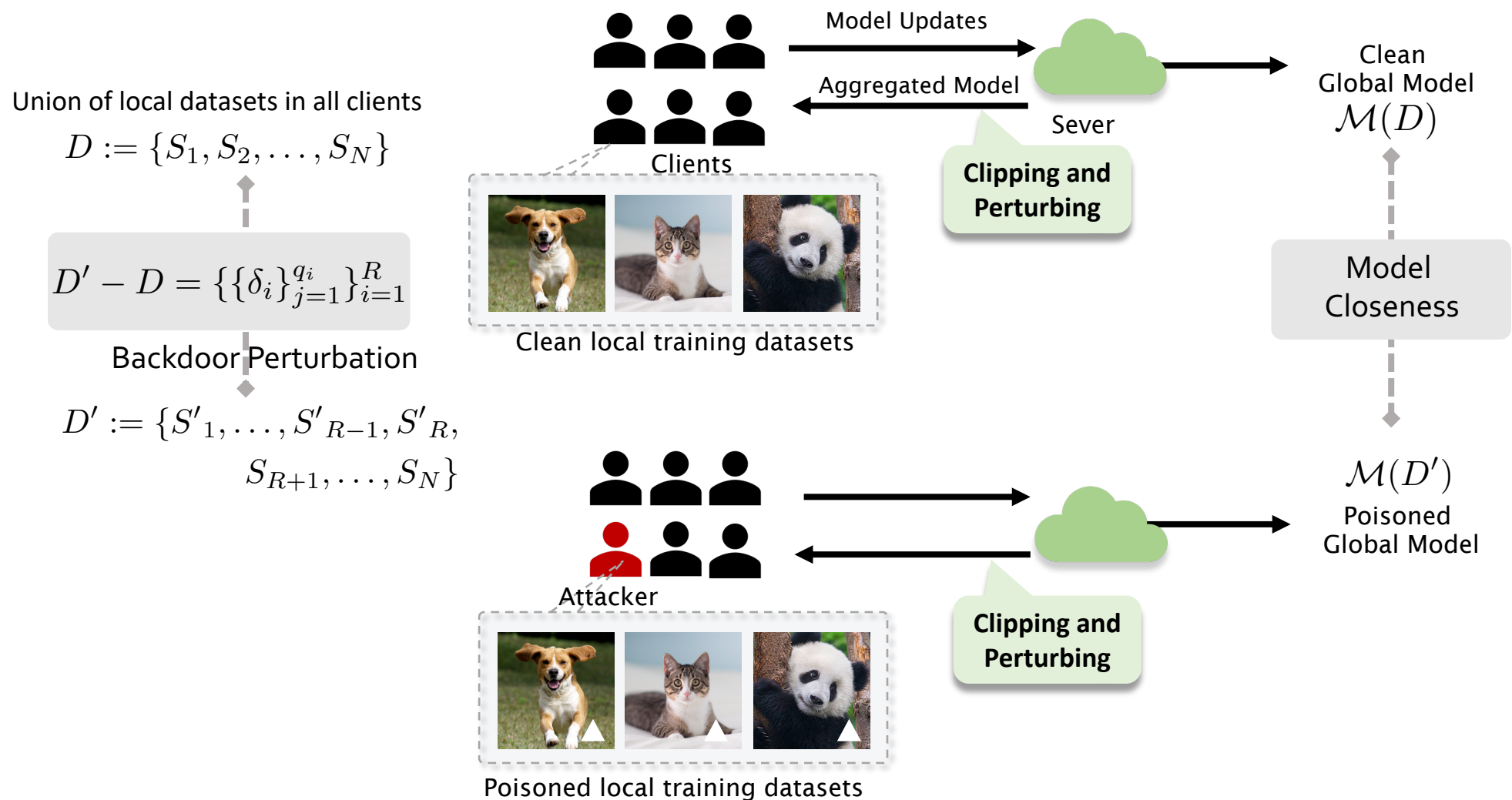
# Motivation

## Robust Federated Learning

Defenses do exist: robust aggregation methods and empirically robust federated training protocols

   *...they lack robustness certification and are adaptively attacked again*

We provide:

- **The *first* general framework**: train certifiably robust FL models against backdoors.

- **Theoretical analysis**: a sample-wise robustness certification on backdoors under certain constraints.

- **Empirical study**: show robustness certification under different FL parameters.

# CRFL Training: Clipping and Perturbing



Union of local datasets in all clients

$$D := \{S_1, S_2, \ldots, S_N\}$$

$$D' - D = \{\{\delta_i\}_{j=1}^{q_i}\}_{i=1}^{R}$$

Backdoor Perturbation

$$D' := \{S'_1, \ldots, S'_{R-1}, S'_R, \\ S_{R+1}, \ldots, S_N\}$$

Model Updates

Aggregated Model

Sever

Clients

Clipping and Perturbing

Clean local training datasets

Clean Global Model
$$\mathcal{M}(D)$$

Model Closeness

Attacker

Clipping and Perturbing

Poisoned local training datasets

$$\mathcal{M}(D')$$

Poisoned Global Model

# CRFL Testing: Parameter Smoothing

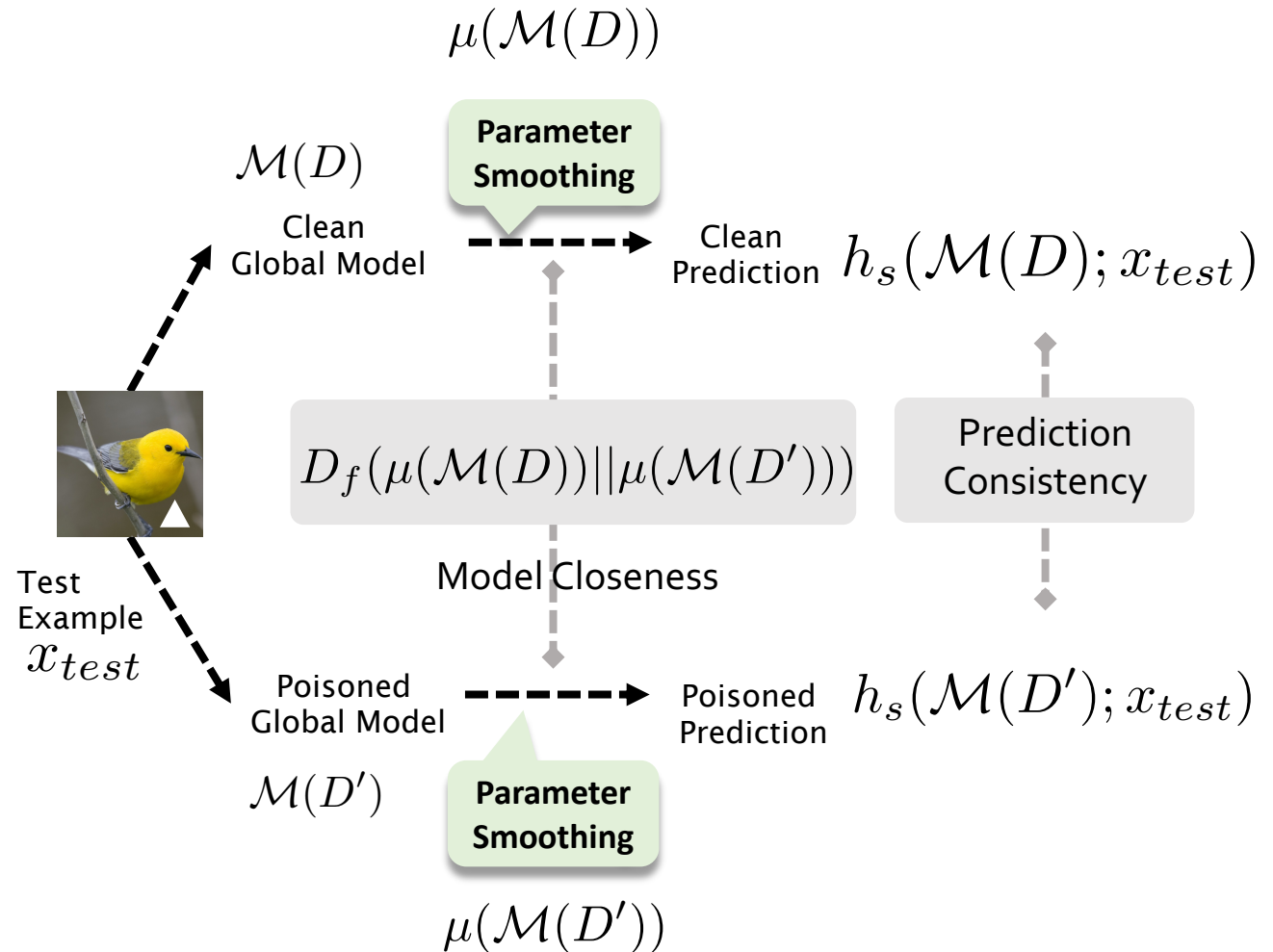Base classifer $h : (\mathcal{W}, \mathcal{X}) \rightarrow \mathcal{Y}$  $\mathcal{Y} = \{1, \ldots, C\}$

Smoothed classifer $h_s$

$$H_s^c(w; x_{test}) = \mathbb{P}_{W \sim \mu(w)}[h(W; x_{test}) = c]$$

*Votes for class c*     $\mu(w) = \mathcal{N}(w, \sigma_T{}^2 \mathbf{I})$

$$h_s(w; x_{test}) = \arg \max_{c \in \mathcal{Y}} H_s^c(w; x_{test})$$

*The majority vote winner*

$\mu(\mathcal{M}(D))$

$\mathcal{M}(D)$

**Parameter Smoothing**

Clean Global Model

Clean Prediction $h_s(\mathcal{M}(D); x_{test})$

Test Example $x_{test}$

$D_f(\mu(\mathcal{M}(D)) || \mu(\mathcal{M}(D')))$

Prediction Consistency

Model Closeness

Poisoned Global Model

Poisoned Prediction $h_s(\mathcal{M}(D'); x_{test})$

$\mathcal{M}(D')$

**Parameter Smoothing**

$\mu(\mathcal{M}(D'))$

# Certification

$$D' - D = \{\{\delta_i\}_{j=1}^{q_i}\}_{i=1}^{R} \impliedby D_f(\mu(\mathcal{M}(D))\|\mu(\mathcal{M}(D'))) \impliedby h_s(\mathcal{M}(D); x_{test}) = h_s(\mathcal{M}(D'); x_{test})$$

Backdoor Perturbation　　　　　Model Closeness　　　　　Prediction Consistency

## General Robustness Condition

$$R\sum_{i=1}^{R}(p_i\gamma_i\tau_i\eta_i\frac{q_{Bi}}{n_{Bi}}\|\delta_i\|)^2 \le \frac{-\log\left(1 - (\sqrt{\underline{p_A}} - \sqrt{\overline{p_B}})^2\right)\sigma_{t_{adv}}^2}{2L_{\mathcal{Z}}^2 \prod_{t=t_{adv}+1}^{T}\left(2\Phi\left(\frac{\rho_t}{\sigma_t}\right) - 1\right)}$$

> Our certification is in three levels: feature, sample, and client.

When the size of the backdoor magnitude is the same for every attackers ⇢

## Robustness Condition in Feature Level

$$\|\delta\| < \text{RAD}$$

$$\text{RAD} = \sqrt{\frac{-\log\left(1 - (\sqrt{\underline{p_A}} - \sqrt{\overline{p_B}})^2\right)\sigma_{t_{adv}}^2}{2RL_{\mathcal{Z}}^2 \sum_{i=1}^{R}(p_i\gamma_i\tau_i\eta_i\frac{q_{Bi}}{n_{Bi}})^2 \prod_{t=t_{adv}+1}^{T}\left(2\Phi\left(\frac{\rho_t}{\sigma_t}\right) - 1\right)}}$$
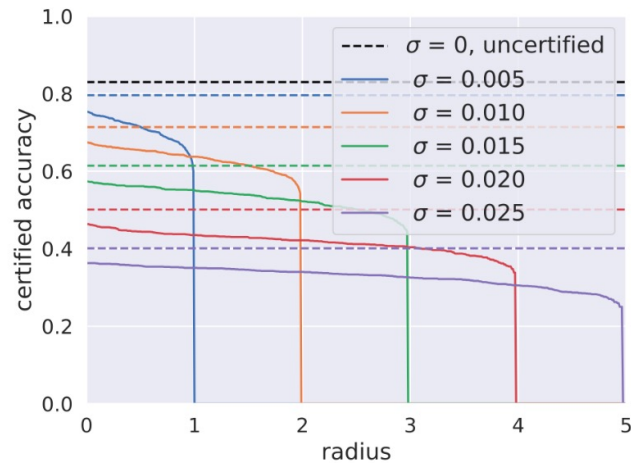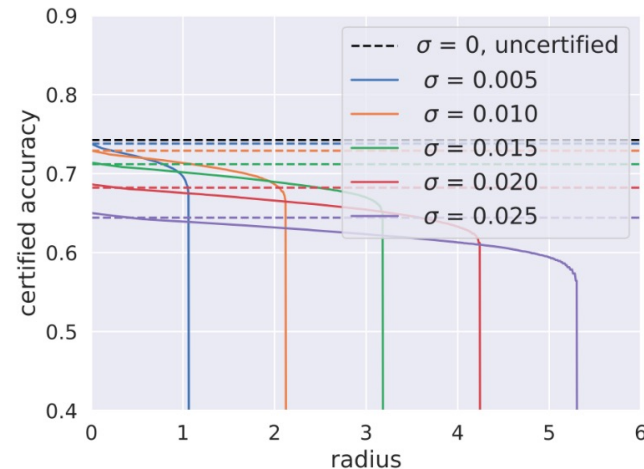
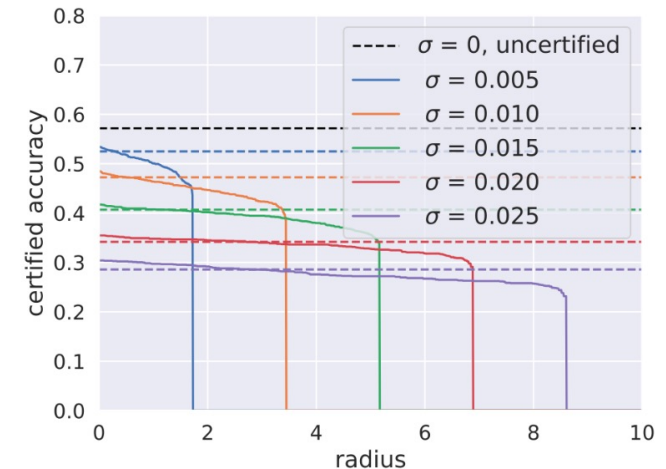Certified radius

# Experiments

*Thank you.*

Effect of different smoothing levels during training



MNIST

LOAN

EMNIST

More details and results are in our paper:

- Effects of smoothing level, attacker ability, robust aggregation, client number, training rounds, etc. on certified robustness.