

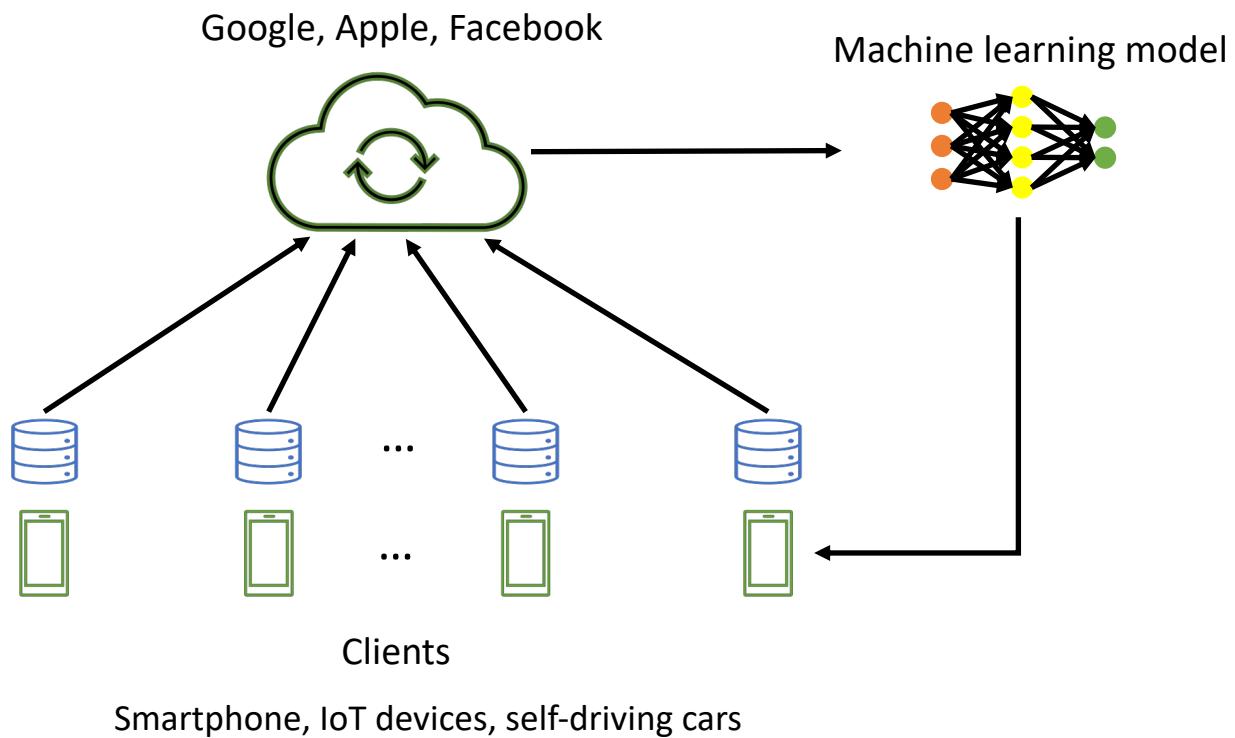
Secure Federated Learning

Neil Gong

Department of Electrical and Computer Engineering
Department of Computer Science (secondary appointment)
Duke University

This talk is available on YouTube: <https://www.youtube.com/watch?v=LP4uqW18yA0>

Conventional Paradigm: Centralized Learning



Challenges of Centralized Learning

- Server data breaches



Over the past 10 years,
there have been **300 DATA
BREACHES** involving the
theft of **100,000 OR
MORE RECORDS.**

Forbes

 VARONIS

- High communications cost
 - Intolerable for resource-constrained clients
 - Smartphone
 - IoT

Federated Learning

- Data stay locally on clients
- Clients train models locally
- Clients send models or updates to server
- Real-world deployment

Artificial intelligence / Machine learning



How Apple personalizes Siri without hoovering up your data

The tech giant is using privacy-preserving machine learning to improve its voice assistant while keeping your data on your phone.

by **Karen Hao**

December 11, 2019

This Talk

What are the security issues of federated learning

How to build secure federated learning

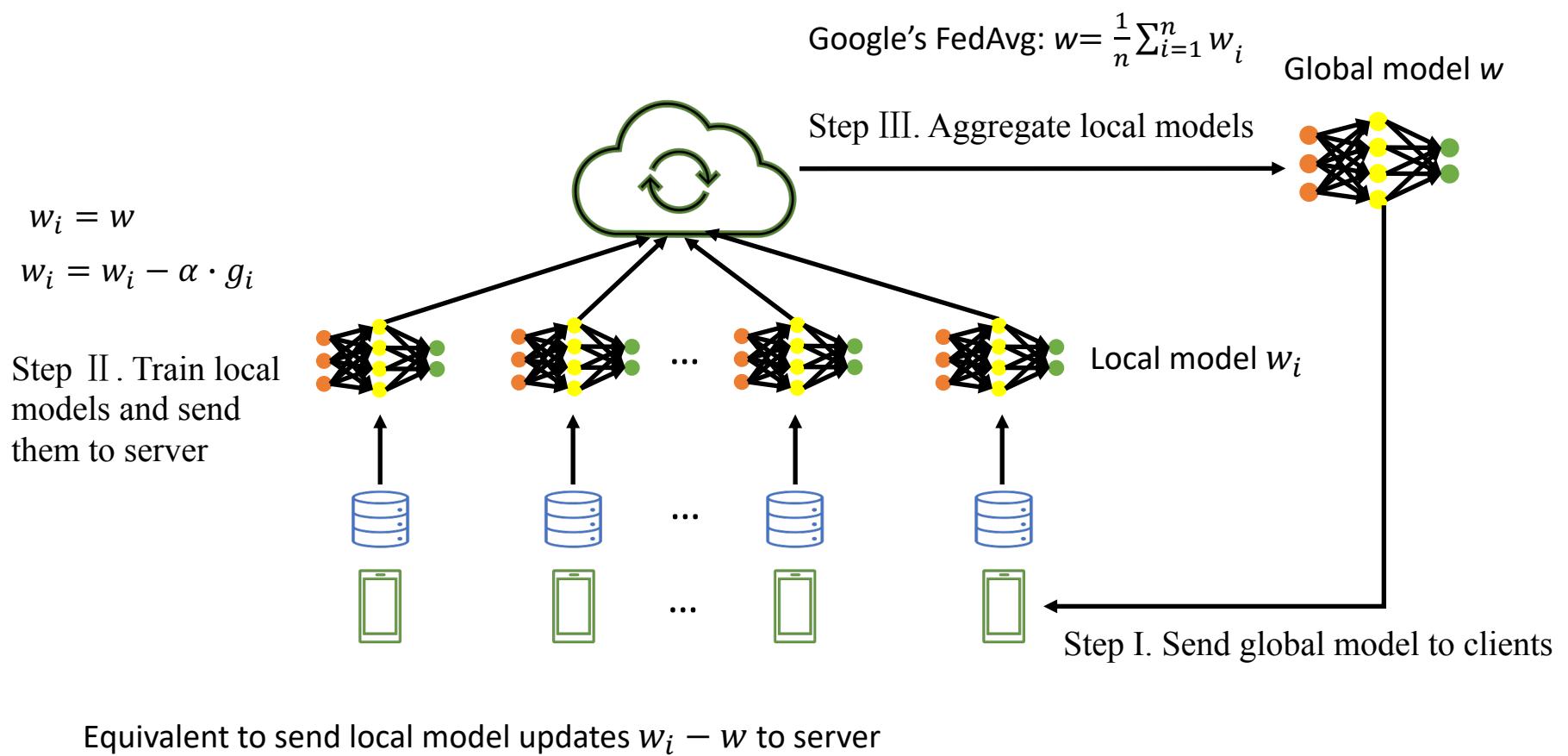
Road Map

- Part I: Local model poisoning attacks to federated learning
- Part II: Secure federated learning via trust bootstrapping
- Part III: Provably secure federated learning

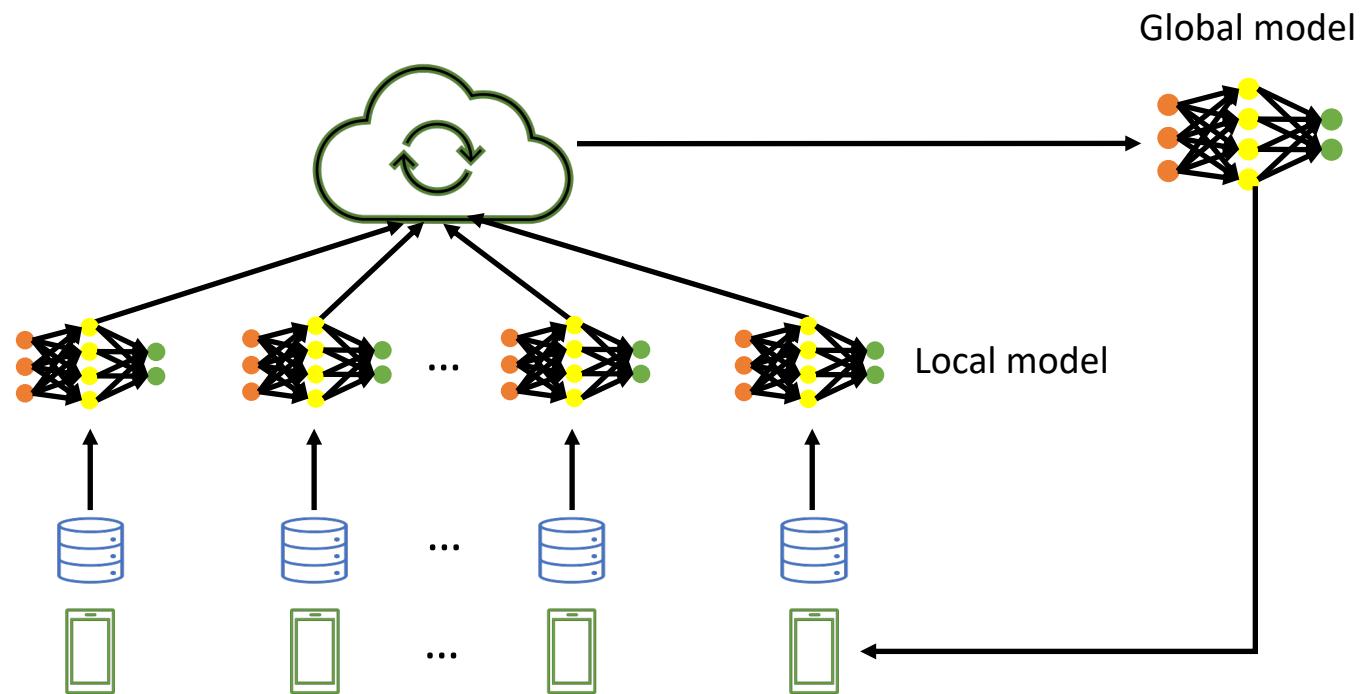
Road Map

- **Part I: Local model poisoning attacks to federated learning**
- Part II: Secure federated learning via trust bootstrapping
- Part III: Provably secure federated learning

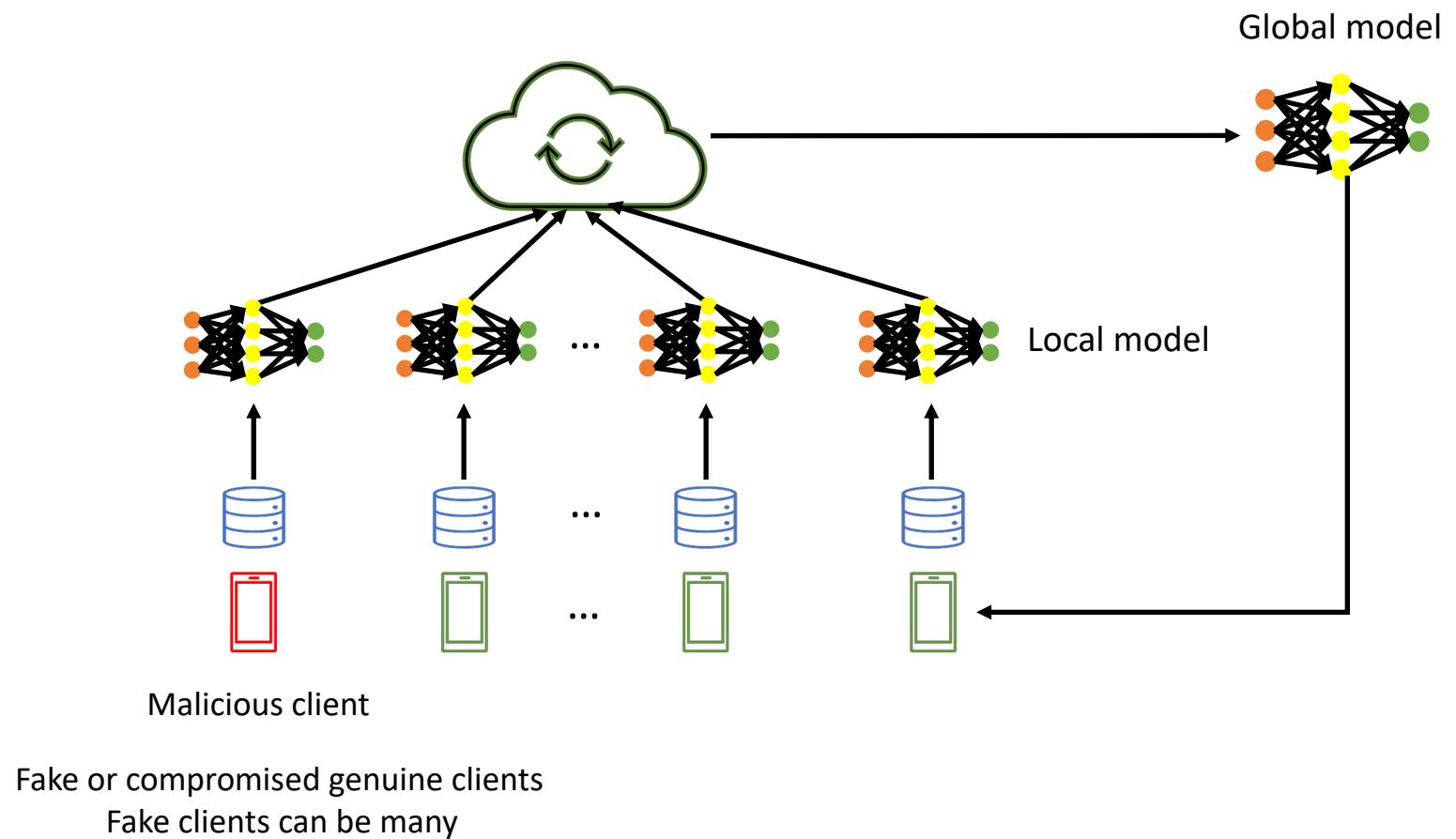
Federated Learning Background



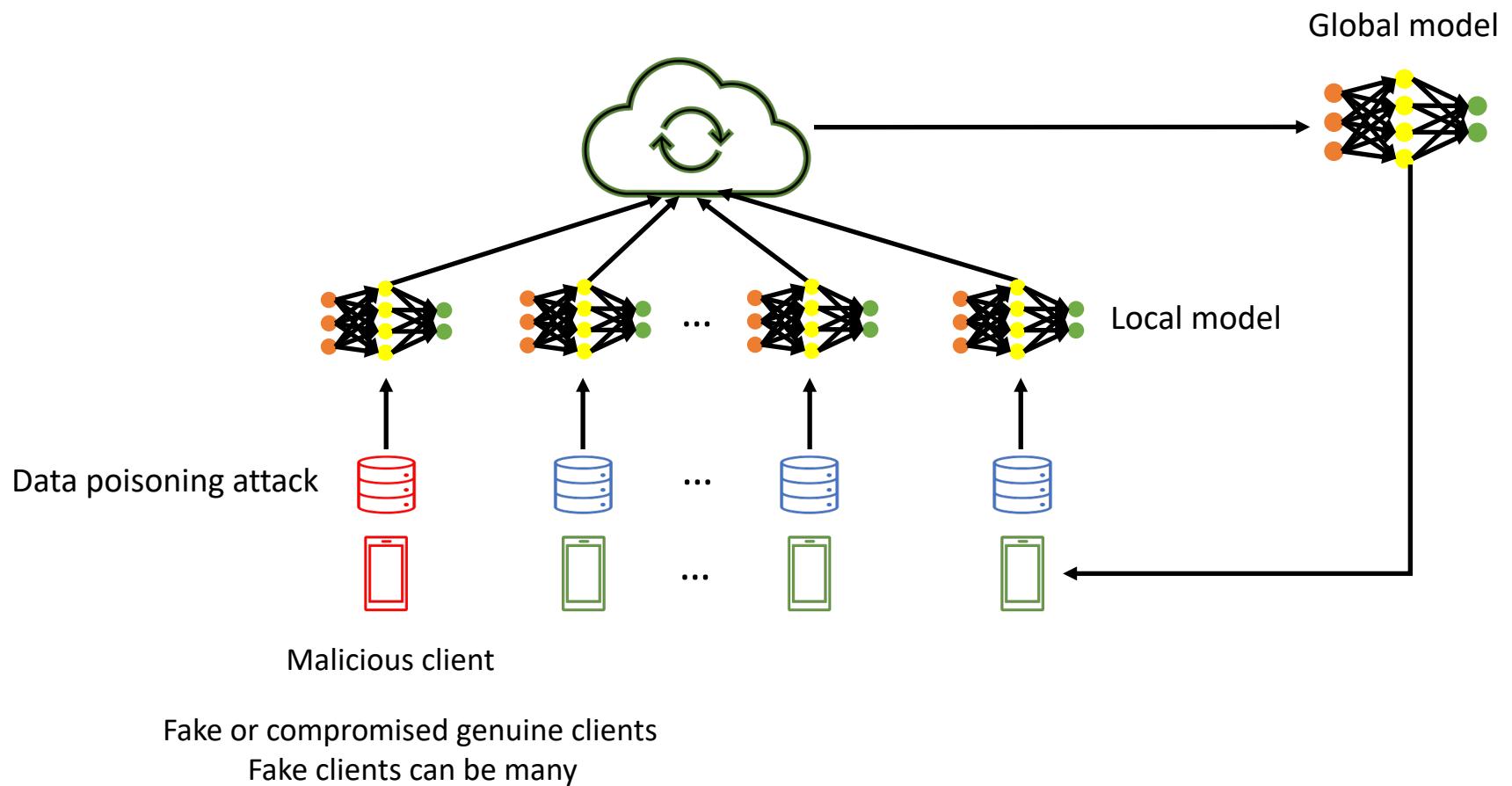
Federated Learning is Vulnerable to Poisoning Attacks



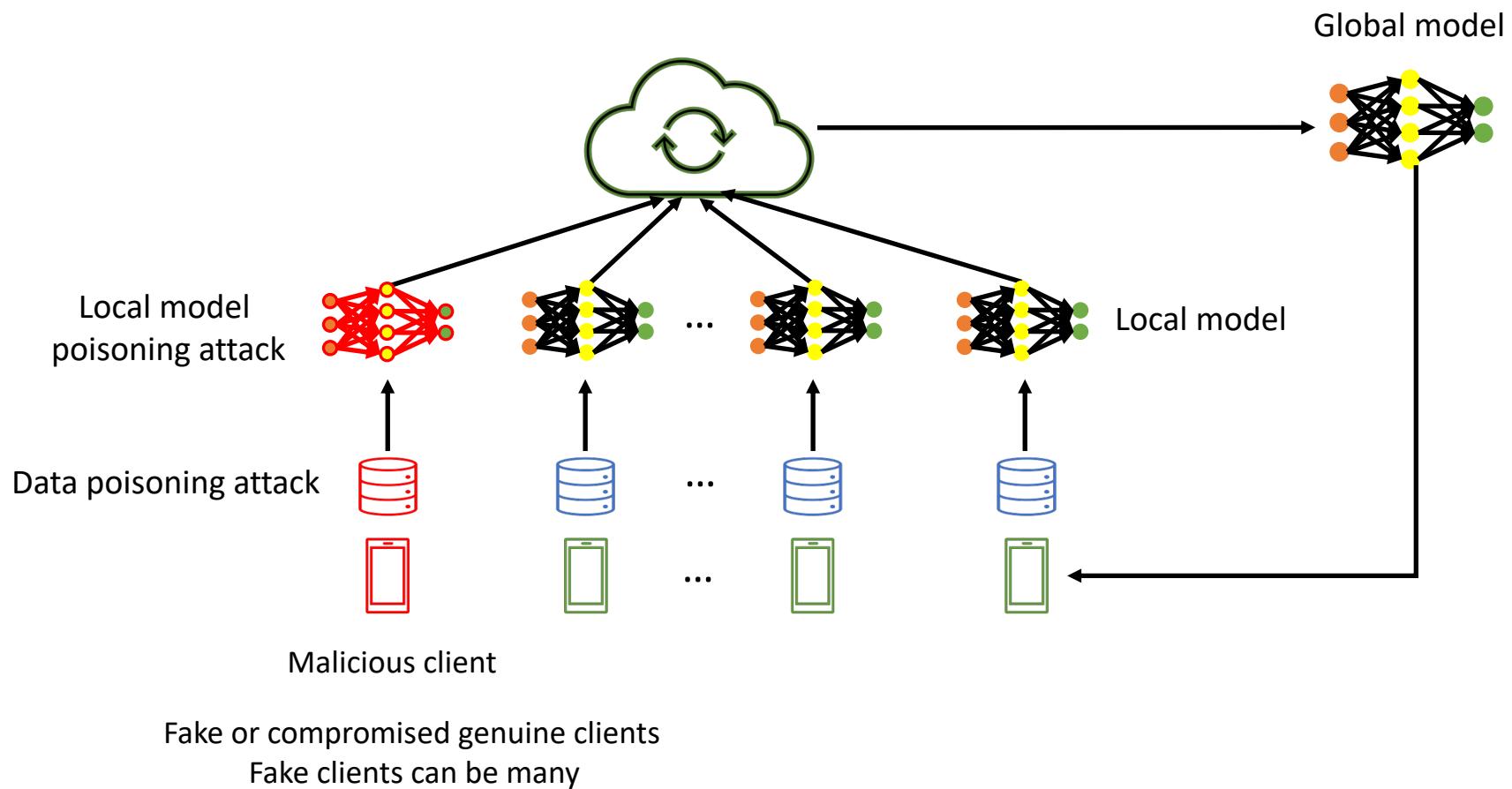
Federated Learning is Vulnerable to Poisoning Attacks



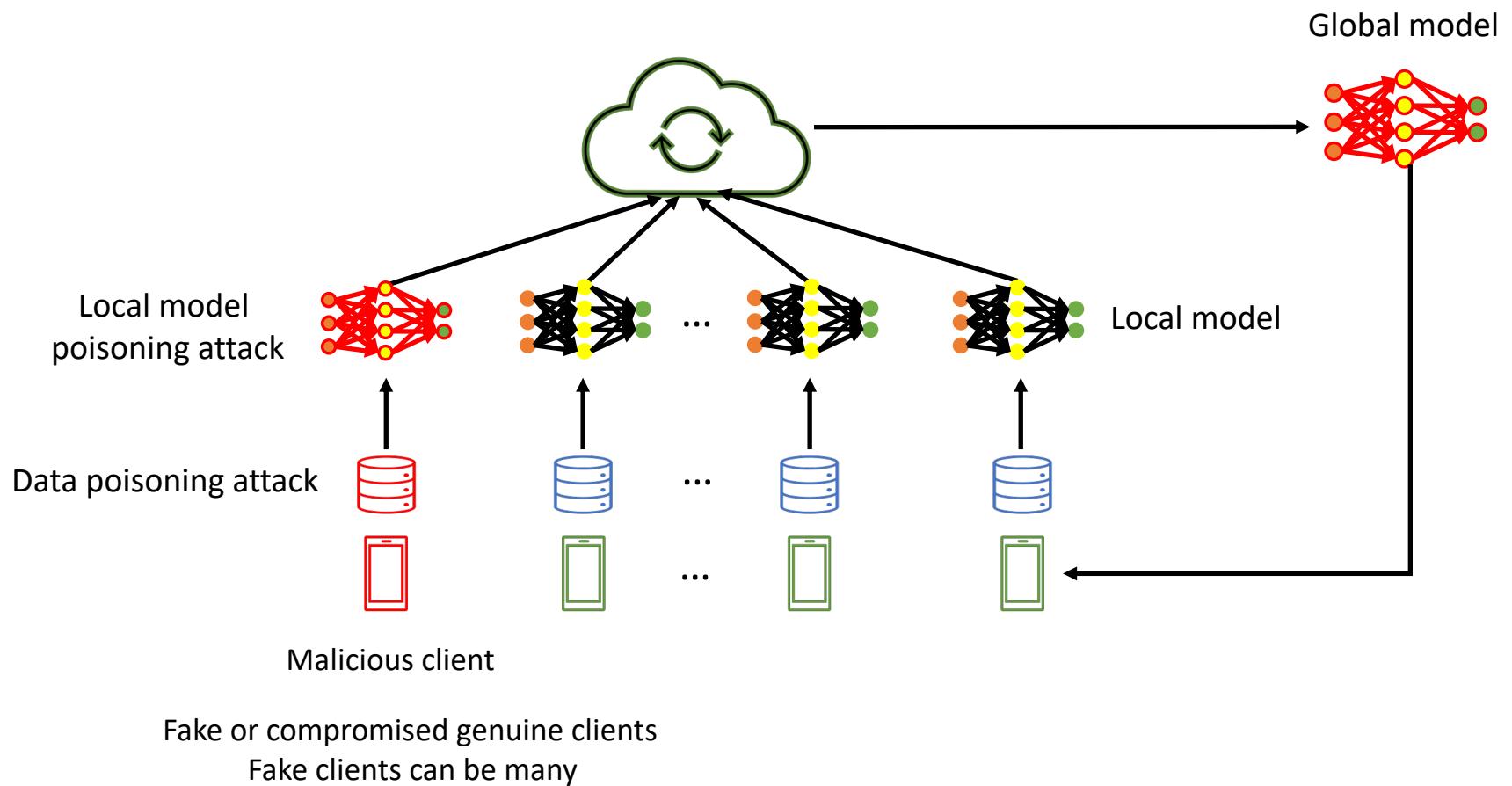
Federated Learning is Vulnerable to Poisoning Attacks



Federated Learning is Vulnerable to Poisoning Attacks



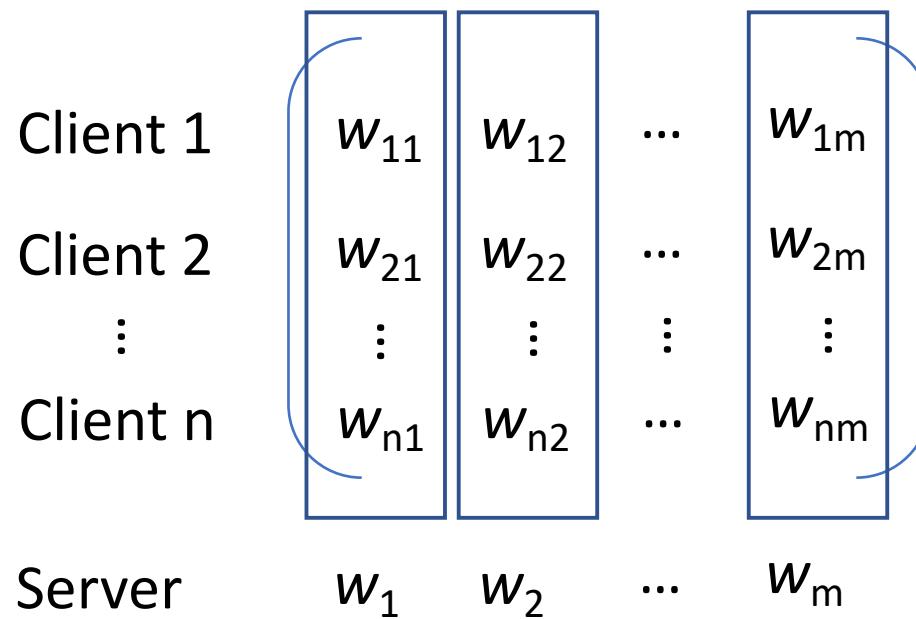
Federated Learning is Vulnerable to Poisoning Attacks



Byzantine-robust Federated Learning as Defense

- Byzantine-robust aggregation rule
 - Krum
 - Trimmed mean
 - Median
- Key idea
 - Remove “outlier” local models
- Theoretical guarantee
 - Various assumptions
 - IID data, smooth loss function, etc.
 - Bound change of global model parameters caused by malicious clients

An Example: Median



Our Work

Byzantine-robust federated learning is vulnerable
to local model poisoning attacks

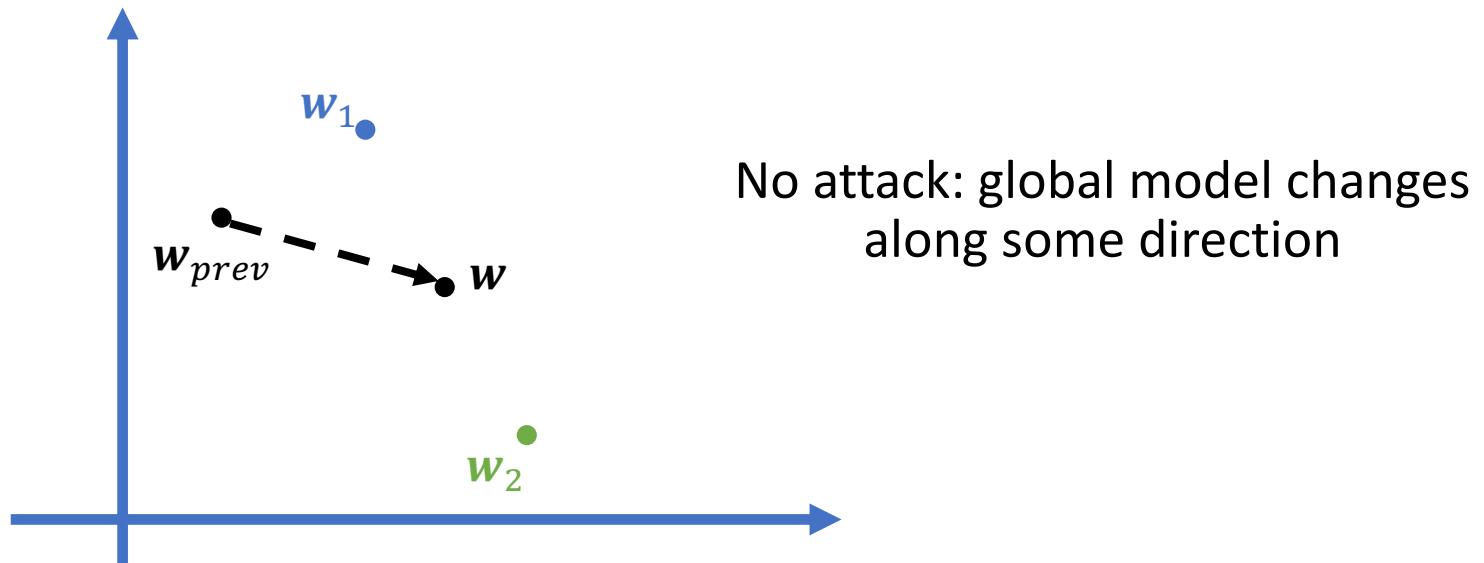
Increase testing error rate of global model

Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. “Local Model Poisoning Attacks to Byzantine-Robust Federated Learning”. In *USENIX Security Symposium*, 2020

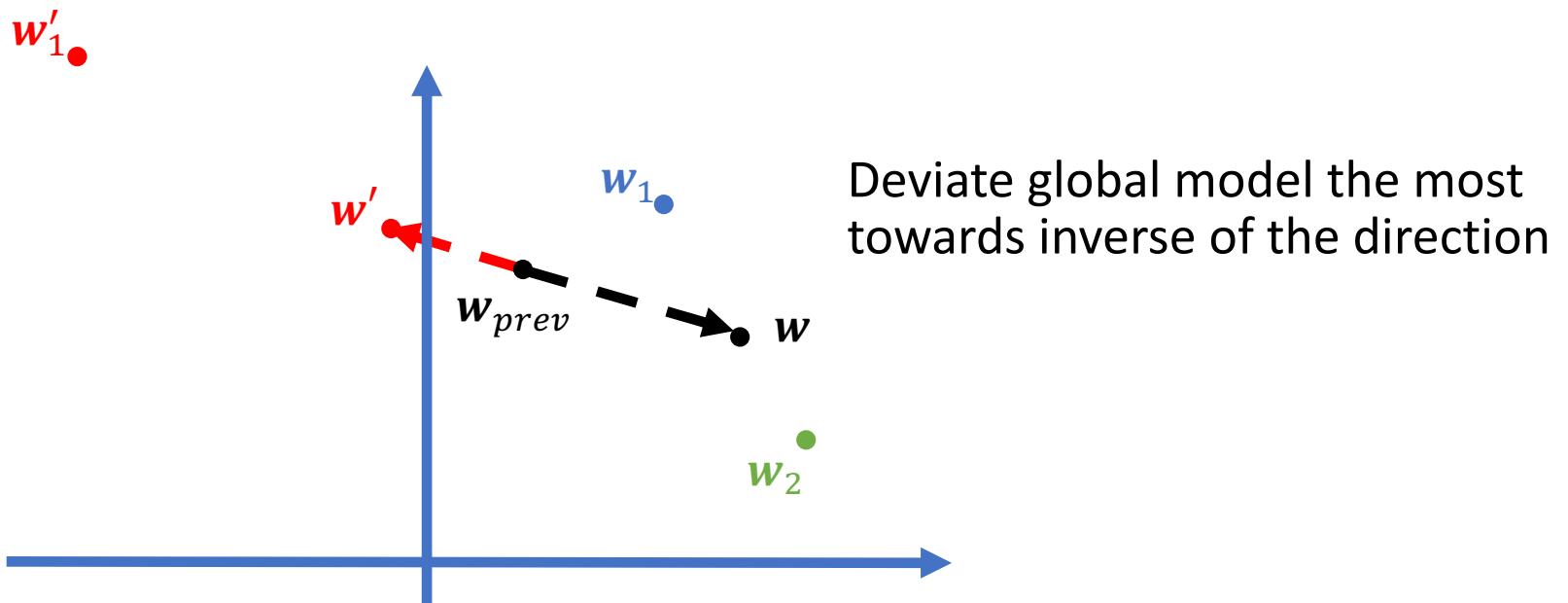
Threat Model

- Attacker's goal
 - High testing error rate
- Attacker's capability:
 - Access to malicious clients
 - Fake clients
 - Compromised genuine clients
 - Send arbitrary local models
- Attacker's knowledge:
 - Full vs. Partial knowledge
 - Data on all vs. malicious clients
 - Aggregation rule
 - Yes or no

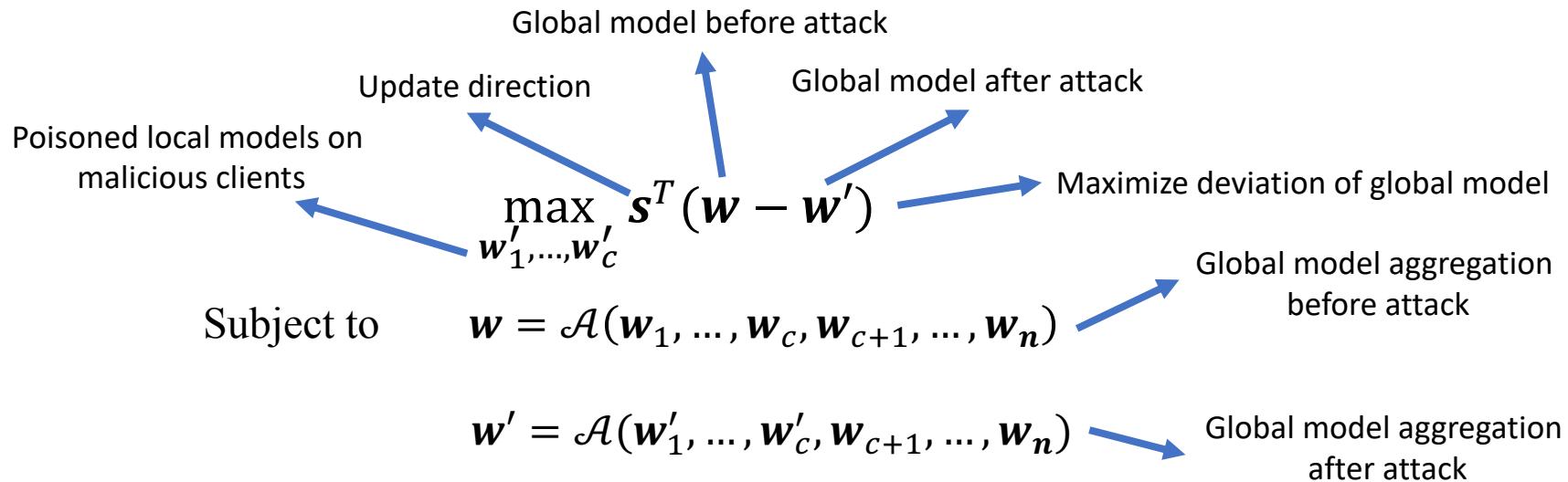
Our Idea



Our Idea



Formulate Optimization Problem



Used in all or multiple iterations

Applicable to **any** aggregation rule

Solving the Optimization Problem

- Full knowledge
 - $w_1, \dots, w_c, w_{c+1}, \dots, w_n$ are known
 - Solve the optimization problem using them
- Partial knowledge
 - Only w_1, \dots, w_c are known
 - Use them to estimate w
- Unknown aggregation rule
 - Attacker assumes one

Experimental Setup

- 100 clients
 - 20% malicious
- Datasets:
 - MNIST
 - Fashion-MNIST
 - CH-MNIST
 - Breast Cancer Wisconsin (Diagnostic)
- Non-IID data on clients
 - Non-IID: not Independently and Identically Distributed

Experimental Setup

- 100 clients
 - 20% malicious
- Datasets:
 - **MNIST**
 - Fashion-MNIST
 - CH-MNIST
 - Breast Cancer Wisconsin (Diagnostic)
- Non-IID data on clients
 - Non-IID: not Independently and Identically Distributed

Our Attack is Effective

	NoAttack	Gaussian	LabelFlip	Partial	Full
Krum	0.11	0.10	0.10	0.75	0.77
Trimmed Mean	0.06	0.07	0.07	0.14	0.23
Median	0.06	0.06	0.16	0.28	0.32

Byzantine-robust methods

Our Attack is Effective

No attack

	NoAttack	Gaussian	LabelFlip	Partial	Full
Krum	0.11	0.10	0.10	0.75	0.77
Trimmed Mean	0.06	0.07	0.07	0.14	0.23
Median	0.06	0.06	0.16	0.28	0.32

Our Attack is Effective

Add Gaussian noise to local models

	NoAttack	Gaussian	LabelFlip	Partial	Full
Krum	0.11	0.10	0.10	0.75	0.77
Trimmed Mean	0.06	0.07	0.07	0.14	0.23
Median	0.06	0.06	0.16	0.28	0.32

Our Attack is Effective

Flip labels of local training data

	NoAttack	Gaussian	LabelFlip	Partial	Full
Krum	0.11	0.10	0.10	0.75	0.77
Trimmed Mean	0.06	0.07	0.07	0.14	0.23
Median	0.06	0.06	0.16	0.28	0.32

Our Attack is Effective

Our attack, partial knowledge

	NoAttack	Gaussian	LabelFlip	Partial	Full
Krum	0.11	0.10	0.10	0.75	0.77
Trimmed Mean	0.06	0.07	0.07	0.14	0.23
Median	0.06	0.06	0.16	0.28	0.32

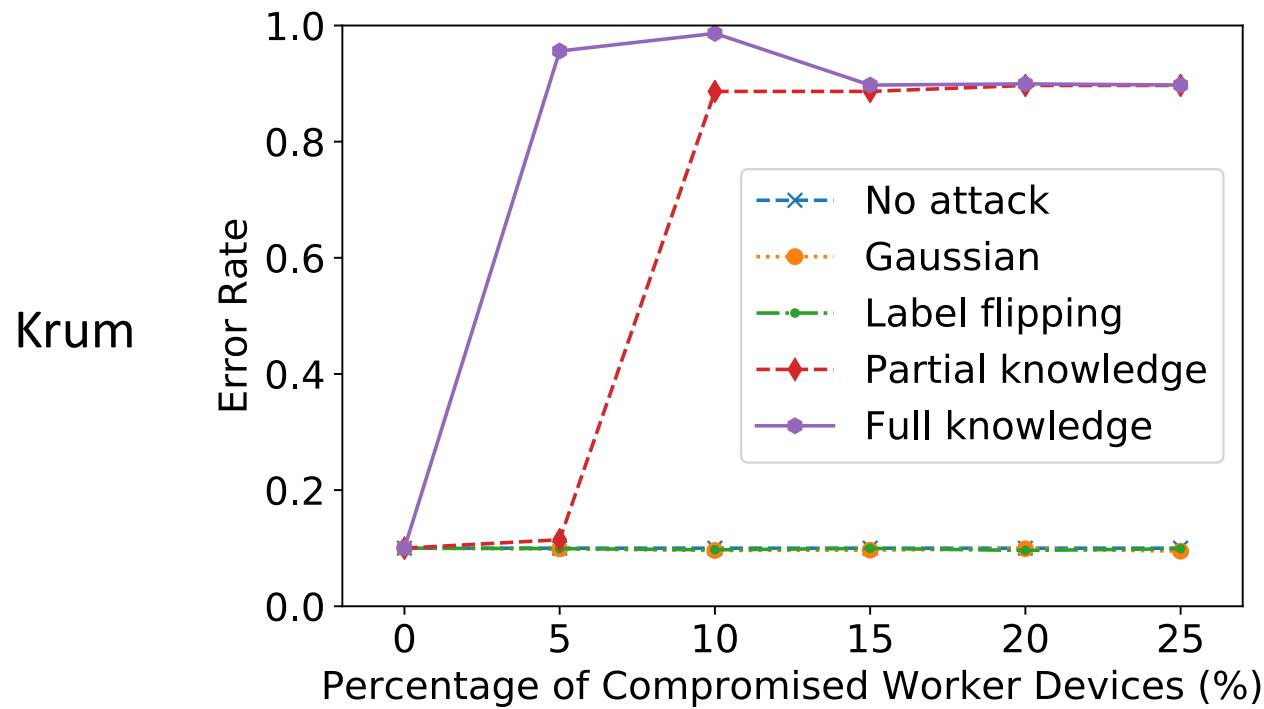
Our Attack is Effective

Our attack, full knowledge

	NoAttack	Gaussian	LabelFlip	Partial	Full
Krum	0.11	0.10	0.10	0.75	0.77
Trimmed Mean	0.06	0.07	0.07	0.14	0.23
Median	0.06	0.06	0.16	0.28	0.32

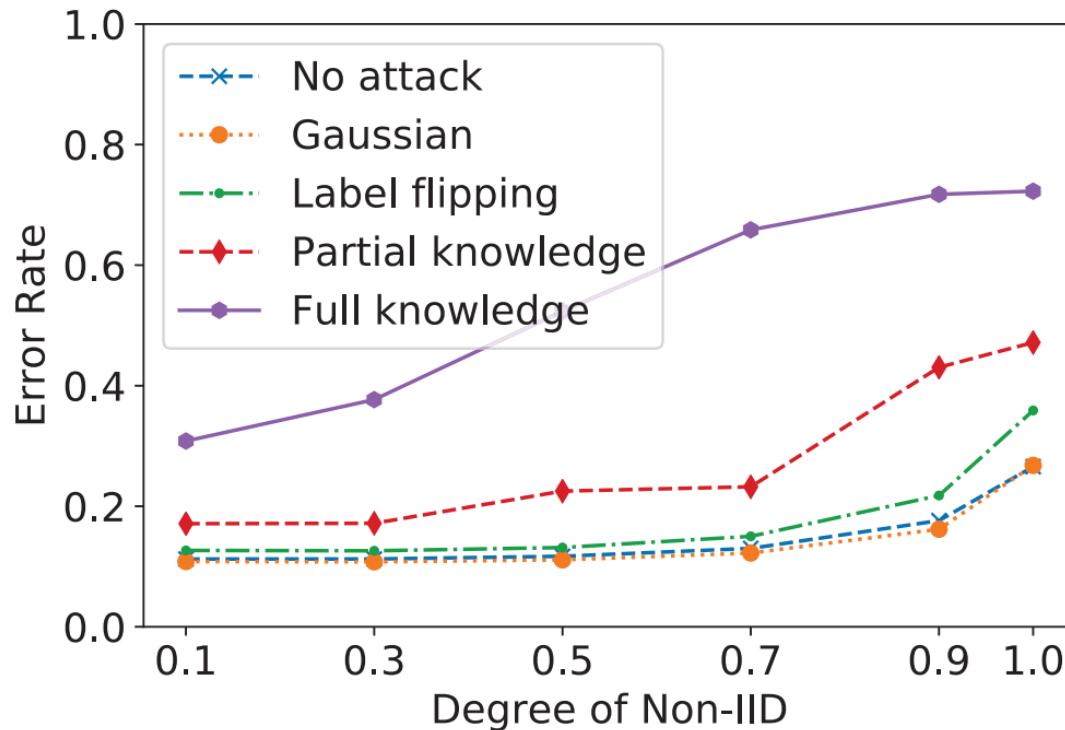
Our attacks can effectively increase testing error rates

Impact of #Malicious Clients



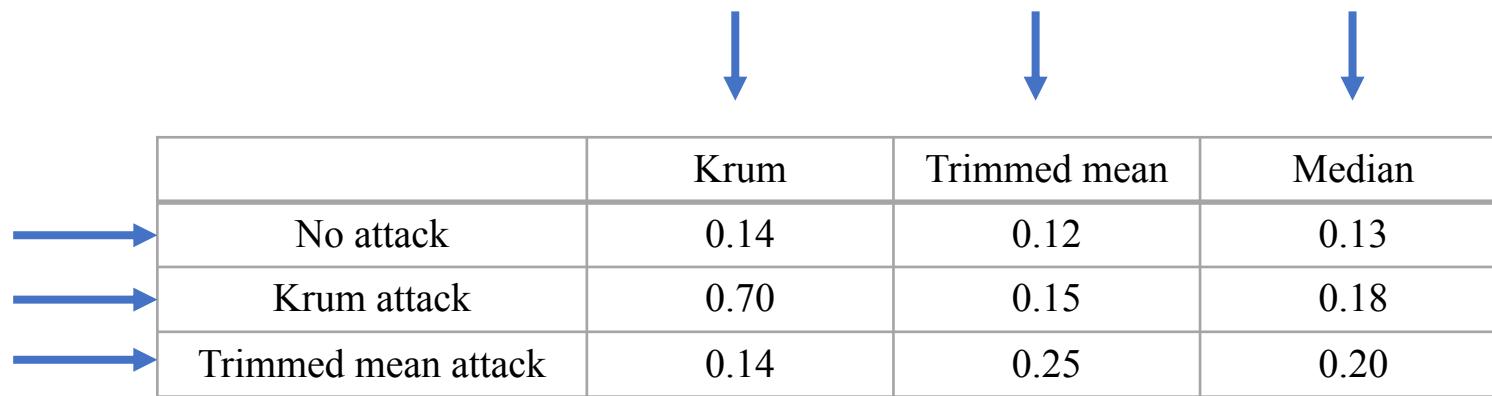
Our attacks are more effective with more malicious clients

Impact of Degree of Non-IID



Our attacks are more effective when clients' data are more Non-IID

Our Attacks Transfer between Aggregation Rules



	Krum	Trimmed mean	Median
No attack	0.14	0.12	0.13
Krum attack	0.70	0.15	0.18
Trimmed mean attack	0.14	0.25	0.20

Comparing with Data Poisoning Attacks

	NoAttack	DataPoisoning	Partial	Full
Krum	0.23	0.24	0.85	0.89
Trimmed Mean	0.12	0.12	0.27	0.32
Median	0.13	0.13	0.19	0.21

Data poisoning attacks are ineffective for Byzantine-robust methods

Our attacks are effective

Summary

- Proposed a general framework to attack federated learning
- Existing Byzantine-robust federated learning is vulnerable to local model poisoning attacks

Road Map

- Part I: Local model poisoning attacks to federated learning
- **Part II: Secure federated learning via trust bootstrapping**
- Part III: Provably secure federated learning

Root Cause of Insecurity

No root trust

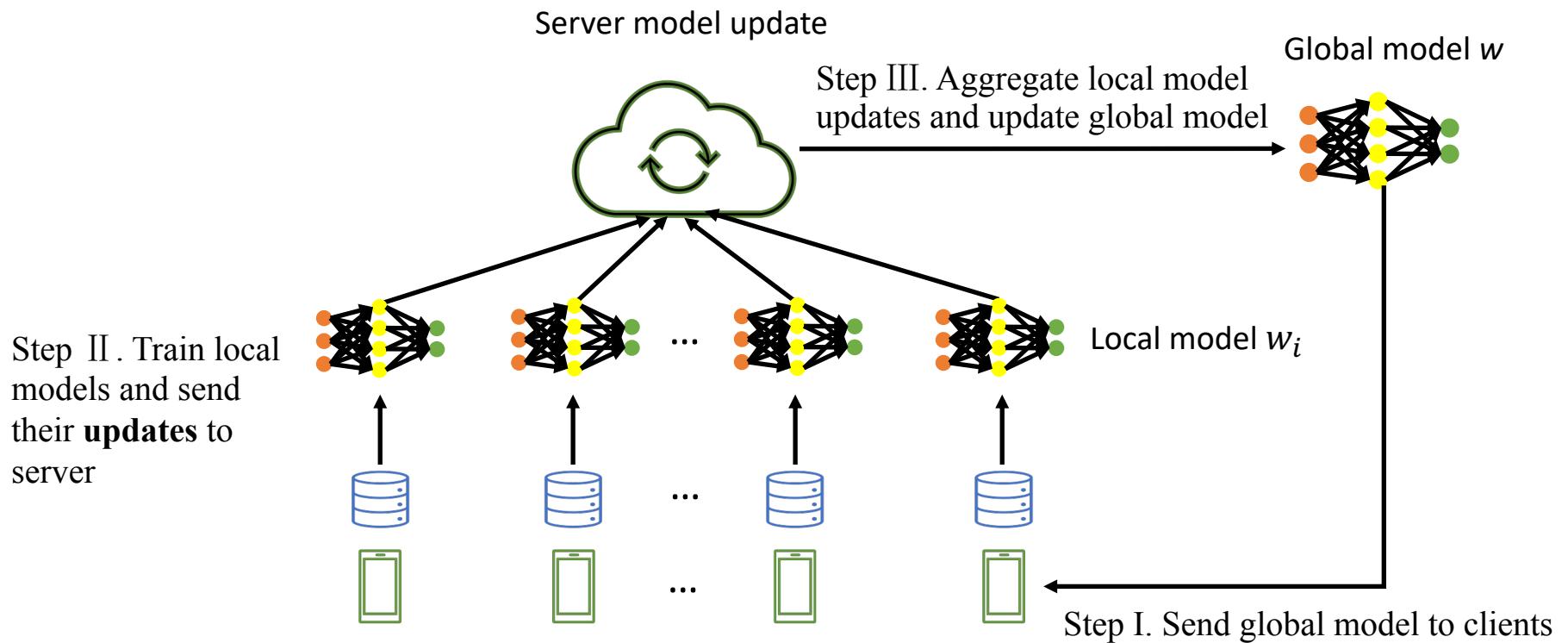
Every client could be malicious

Our FLTrust: Bootstrapping Trust

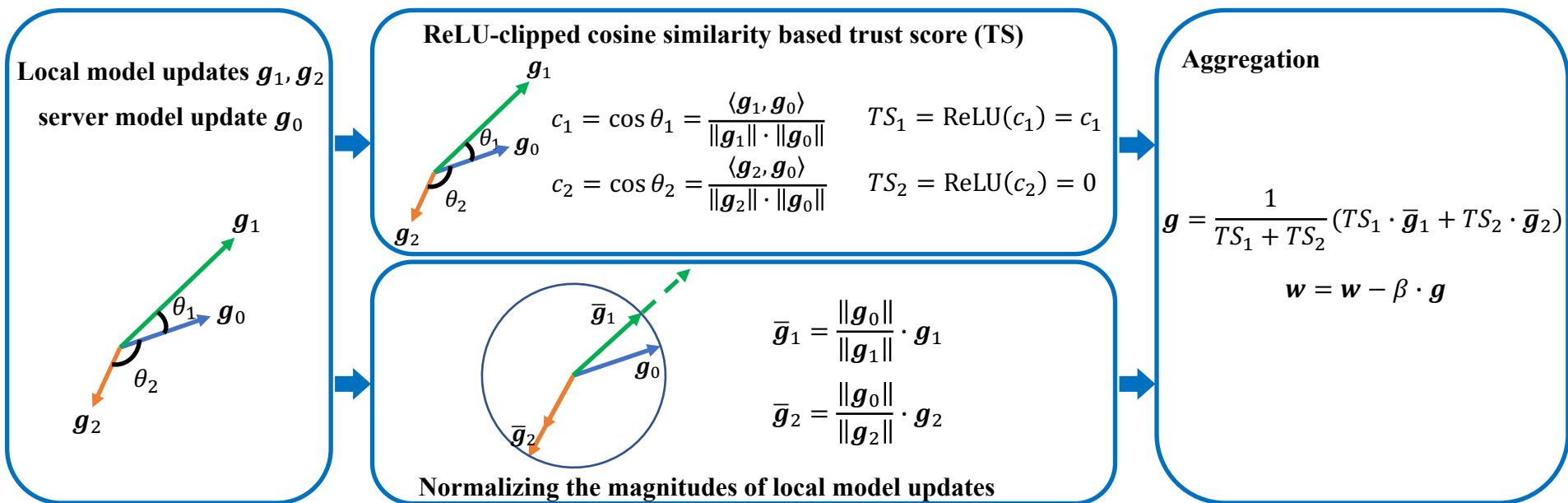
- Server collects a small, clean training dataset
- Server maintains a *server model*
 - Like how a client maintains a local model
- Use server model to bootstrap trust
 - Assign trust scores to clients

Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. “FLTrust: Byzantine-robust Federated Learning via Trust Bootstrapping”. In *ISOC Network and Distributed System Security Symposium (NDSS)*, 2021.

Revisiting Federated Learning Background



Our Aggregation Rule



Theoretical Analysis

Under certain assumptions, for an arbitrary number of malicious clients, the difference between the global model learnt by FLTrust and the optimal global model under no attacks is bounded

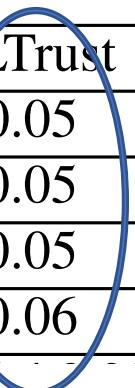
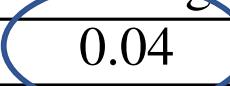
Empirical Results

MNIST

100 clients, 20 malicious

Server's training dataset: 100 examples sampled from MNIST

State-of-the-art method in non-adversarial settings



	FedAvg	Krum	Trim-mean	Median	FLTrust
No attack	0.04	0.10	0.06	0.06	0.05
Label flipping attack	0.06	0.10	0.06	0.06	0.05
Krum attack	0.10	0.91	0.14	0.15	0.05
Trim attack	0.28	0.10	0.23	0.43	0.06

Our FLTrust is robust against poisoning attacks

Adaptive Attack

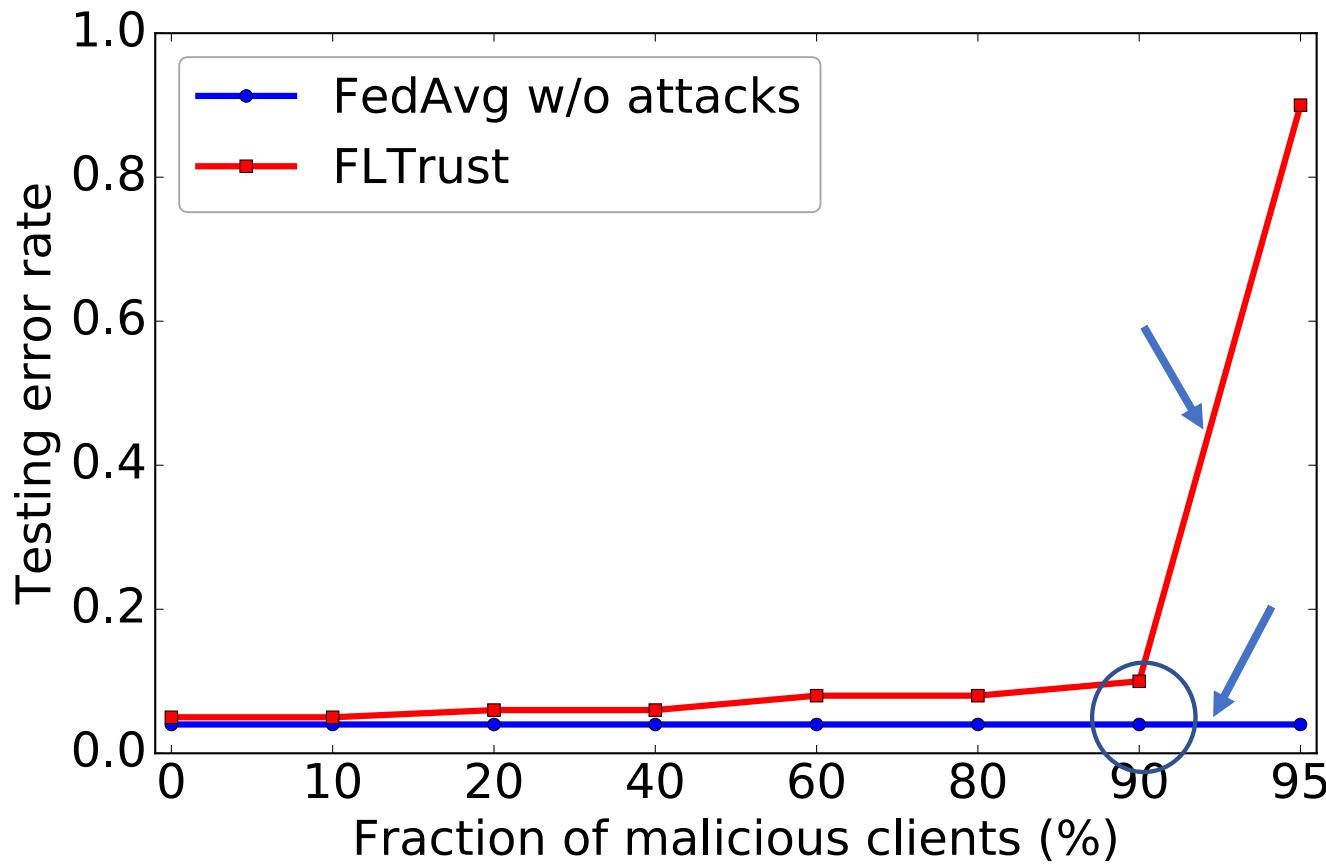
$$\max_{\mathbf{w}'_1, \dots, \mathbf{w}'_c} \mathbf{s}^T (\mathbf{w} - \mathbf{w}')$$

Subject to $\mathbf{w} = \mathcal{A}(\mathbf{w}_1, \dots, \mathbf{w}_c, \mathbf{w}_{c+1}, \dots, \mathbf{w}_n)$

$$\mathbf{w}' = \mathcal{A}(\mathbf{w}'_1, \dots, \mathbf{w}'_c, \mathbf{w}_{c+1}, \dots, \mathbf{w}_n)$$

Applicable to **any** aggregation rule

Our FLTrust is Robust against Adaptive Attack



Summary

- The server can enhance security of federated learning via collecting a small, clean training dataset to bootstrap trust

Road Map

- Part I: Local model poisoning attacks to federated learning
- Part II: Secure federated learning via trust bootstrapping
- **Part III: Provably secure federated learning**

Limitations of Byzantine-robust Federated Learning

- Bound change in global model parameters caused by malicious clients
 - Under assumptions
 - IID data on clients
 - Smooth loss function
 - ...
- Limitations
 - Assumptions do not hold
 - Not bound testing error rate or accuracy

Our Provably Secure Federated Learning

- Guarantee a lower bound of testing accuracy
- Only assumption
 - Bounded #malicious clients

Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. “Provably Secure Federated Learning against Malicious Clients”. In AAAI, 2021.

Defining Provable Security

Label predicted for x when the global model is trained on C

$$h(C', x) = h(C, x) \text{ for any } C', \#\text{malicious clients} \leq m^*$$

```
graph TD; A[h(C, x)] --> B[h(C', x)]; B --> C[A set of benign clients]; B --> D[Testing input]
```

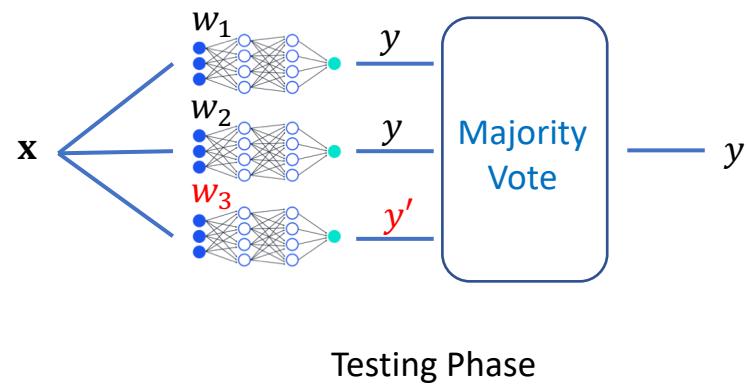
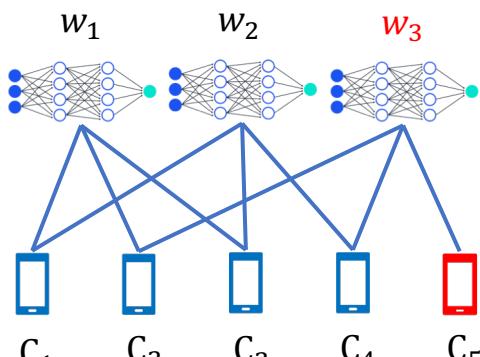
A federated learning algorithm is provably secure if its predicted label for a testing input is not affected by a bounded number of malicious clients

m^* : certified security level for x

Our Ensemble Federated Learning: the First Provably Secure Method

- Training
 - n clients
 - Select k clients randomly and train a global model
 - Use any federated learning method, e.g., FedAvg
 - Repeat to train N global models
- Testing
 - Majority vote of the N global models to predict label of x

Provable Security: Intuition



Provable Security

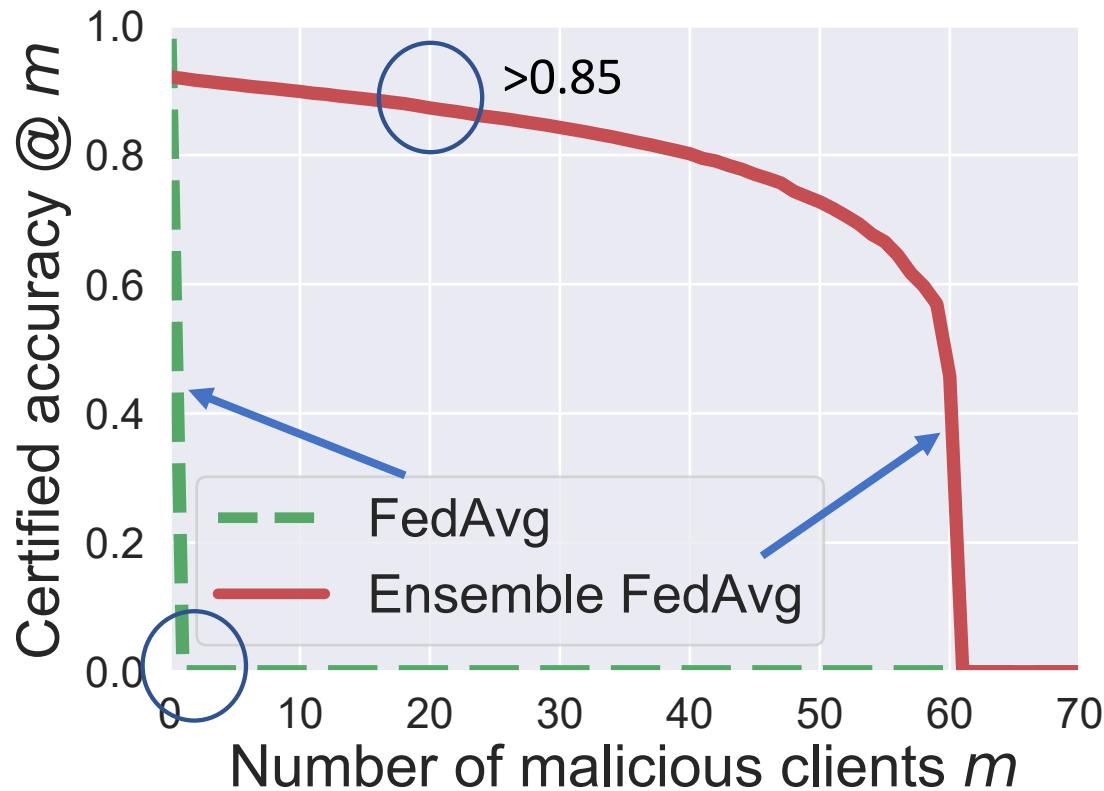
Given C and x , we can derive the certified security level m^* for x

Our derived certified security level is tight

Evaluation Metric: Certified Accuracy @ m

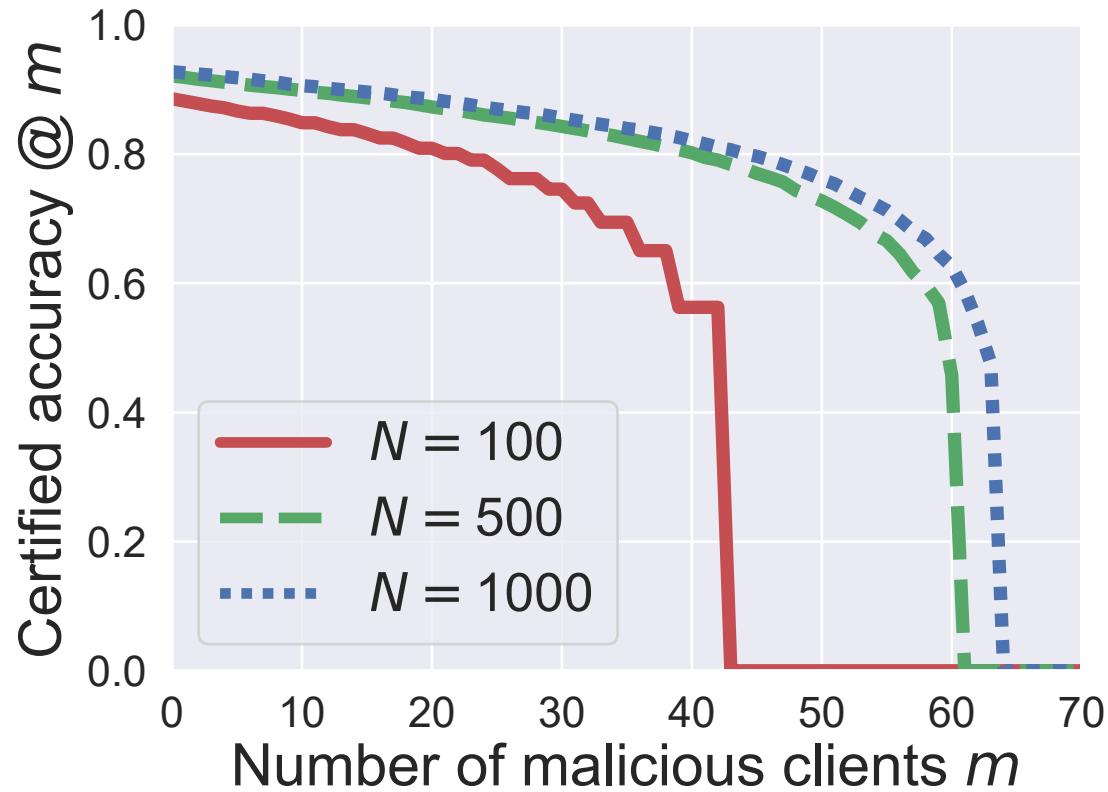
- Fraction of testing inputs whose
 - Labels are correctly predicted
 - Certified security levels are at least m
- A lower bound of testing accuracy
 - $\#\text{malicious clients} \leq m$
 - No matter what attacks are used!

FedAvg vs. Ensemble FedAvg



MNIST dataset, 1,000 clients

Impact of Number of Global Models N



A moderate number of global models are enough

Summary

- Ensemble federated learning is provably secure against bounded number of malicious clients
- Achieve certified accuracy
 - A lower bound of testing accuracy
 - No matter what attacks are used

Conclusion

- Part I: Local model poisoning attacks to federated learning
 - “Local Model Poisoning Attacks to Byzantine-Robust Federated Learning”. In *Usenix Security Symposium*, 2020.
- Part II: Secure federated learning via trust bootstrapping
 - “FLTrust: Byzantine-robust Federated Learning via Trust Bootstrapping”. In *NDSS*, 2021.
- Part III: Provably secure federated learning
 - “Provably Secure Federated Learning against Malicious Clients”. In *AAAI*, 2021.

Acknowledgements

Xiaoyu Cao
Jinyuan Jia

Minghong Fang
Jia Liu