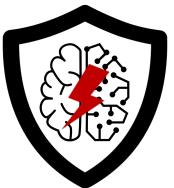


Security and Robustness of Collaborative Learning Systems

Anwar Hithnawi

Collaborative Learning



Robustness

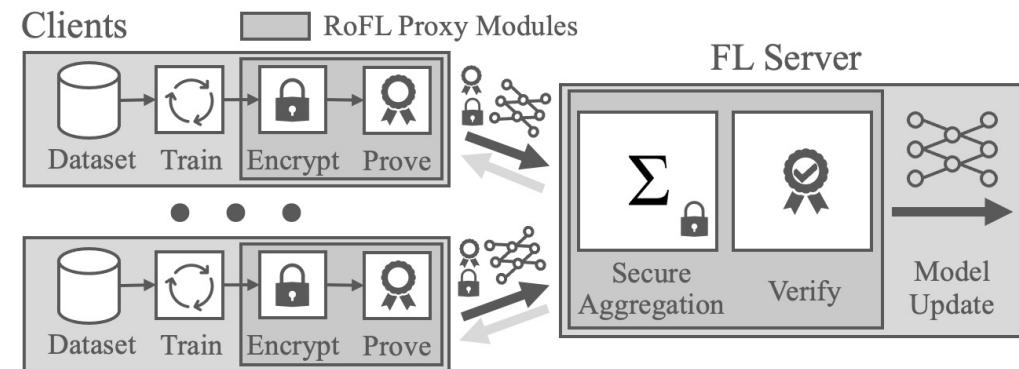


Security



Privacy

RoFL



RoFL: Attestable Robustness for Secure FL. arXiv:2107.03311
L. Burkhalter*, H. Lycklama*, A. Viand, N. Küchler, A. Hithnawi



Analysis Code: <https://github.com/pps-lab/fl-analysis>

RoFL Code: <https://github.com/pps-lab/rofl-project-code>



Autonomous Driving

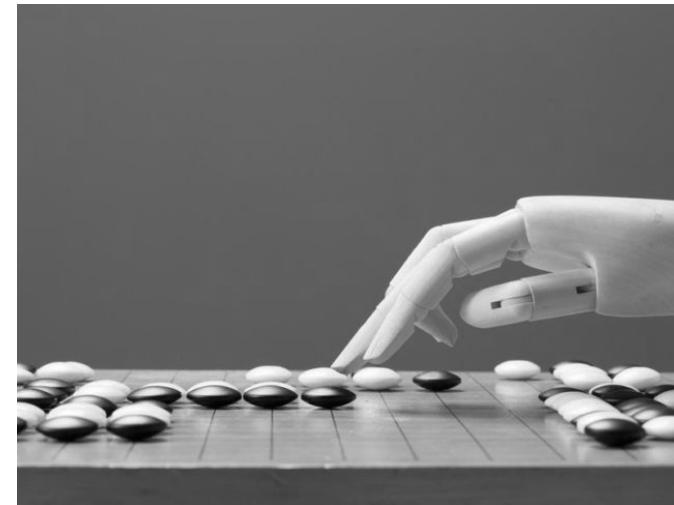


Health Care

Data Driven World



Object Classification



AlphaGo



Data

large, diverse data → broad generalization



World



Health Care

AlphaGo

Solving tasks where data is accessible...



Public Data Crowdsourced Data

For example: web, books, articles, science, TV, corpus, audiobooks, ...

... however, many important tasks we care about ...

Inaccessible

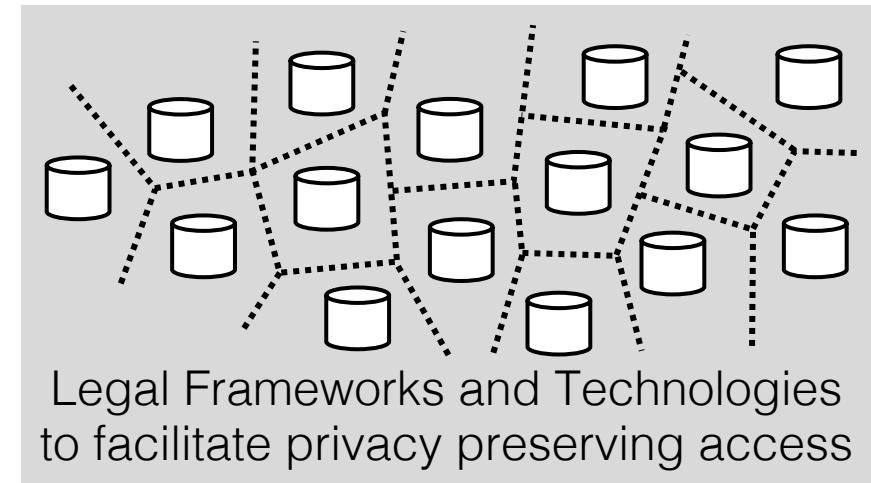
Health – Cancer, Alzheimer, Dementia, Depression

Finance – Economic growth, Market predictions

Government – Education, Taxes, Immigration, Income

Personal Data – Text Messages, Emails, Photos

→ EU Data Governance Act (**DGA**)
effective from 2023
facilitate the reuse of protected public-sector data

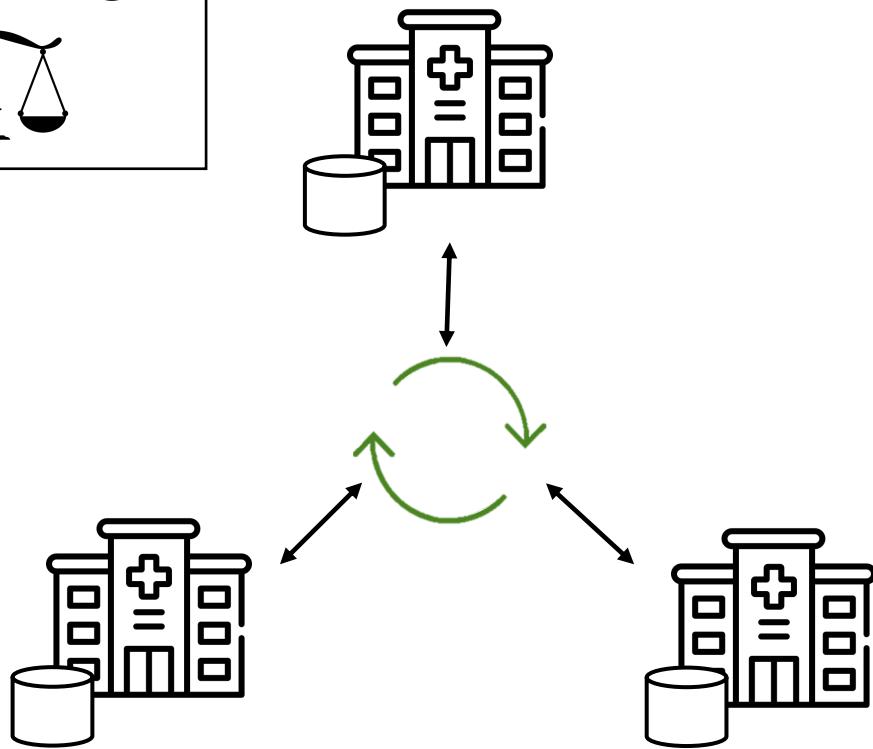


Data Silos

- Privacy Laws
- Competition

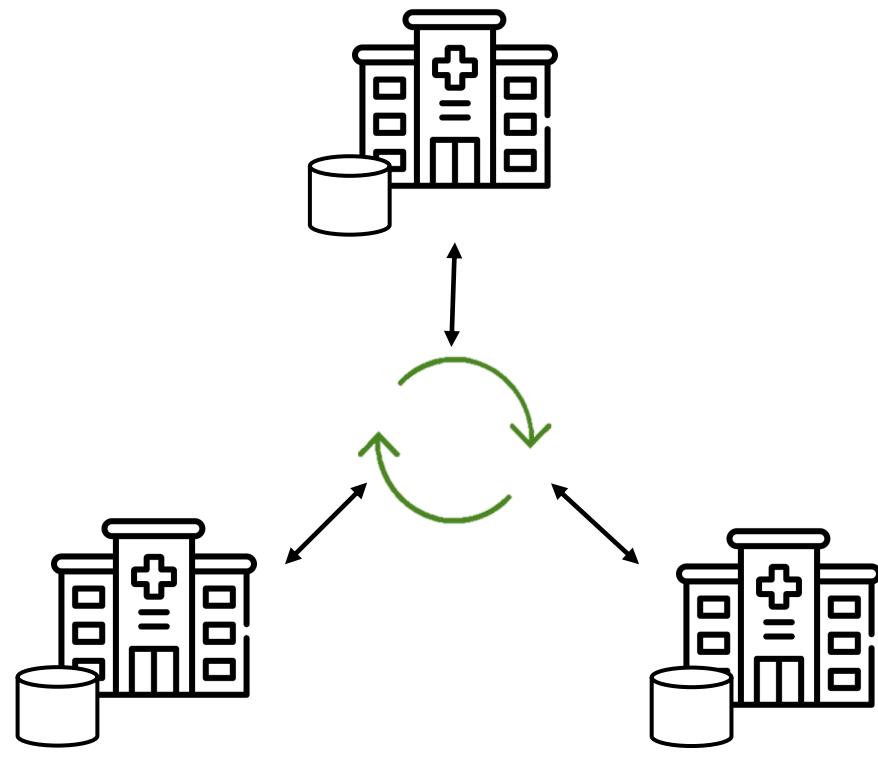
Collaborative Learning

Collaborative Learning

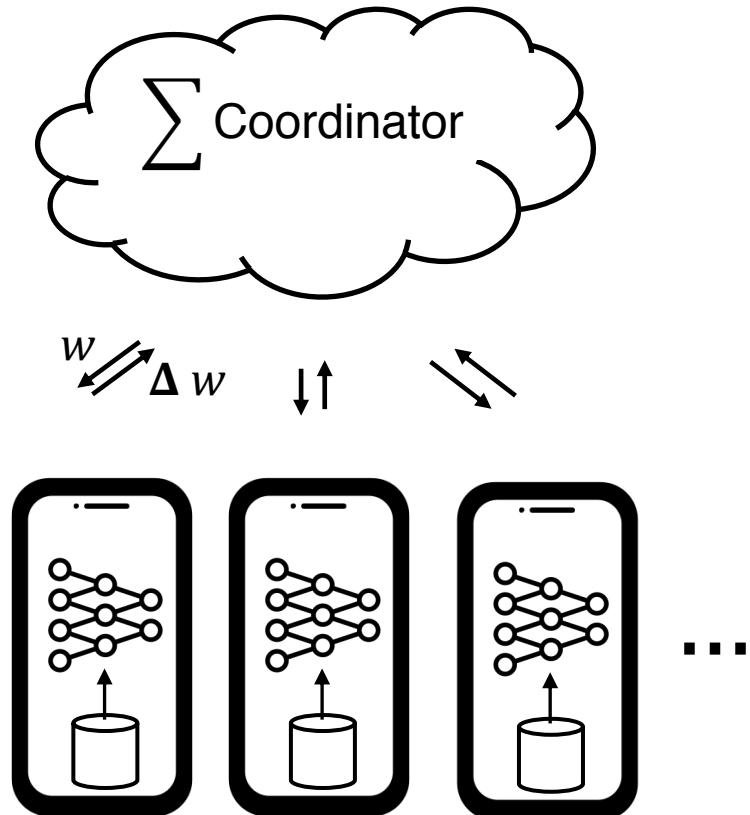


Decentralized Learning

Collaborative Learning



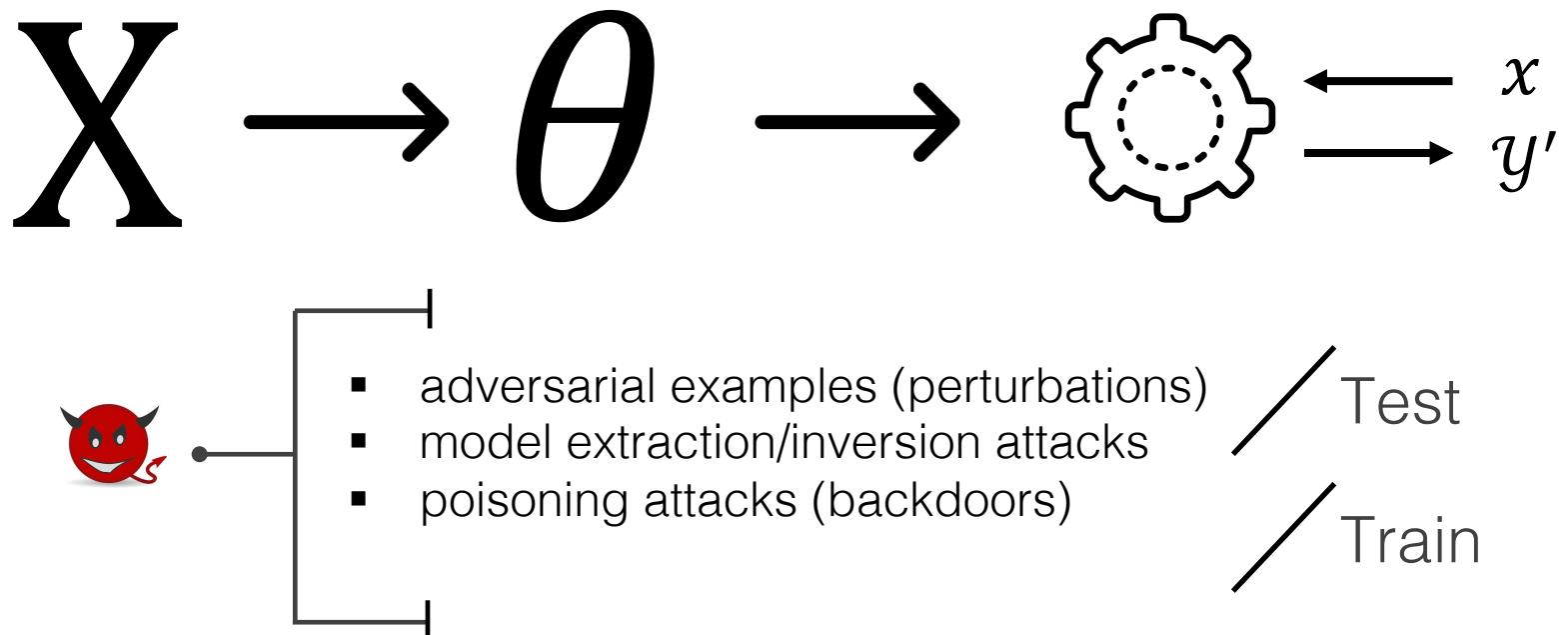
Decentralized Learning



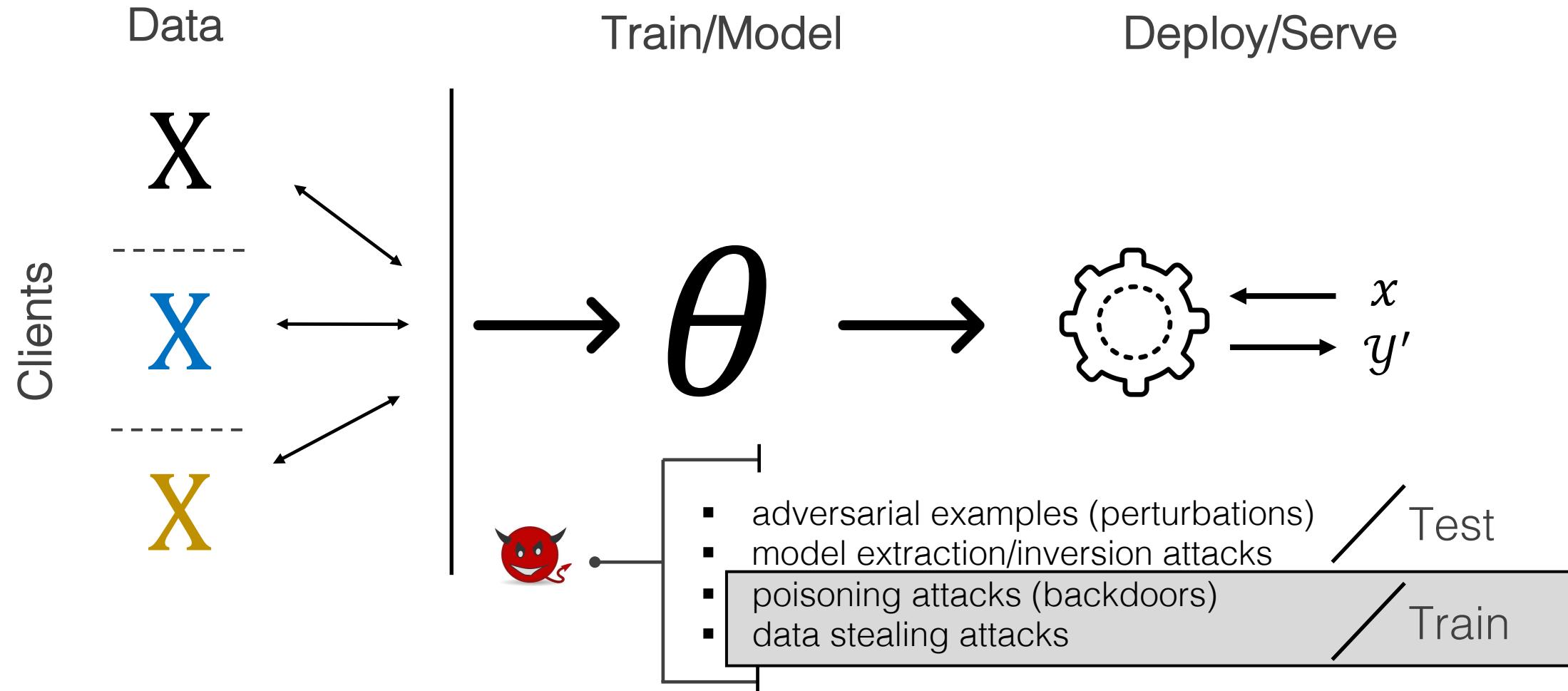
Federated Learning

Security and Privacy of Machine Learning

Data Train/Model Deploy/Serve



Security and Privacy of Collaborative ML



Confidentiality of Input Data

Federation ≠ Privacy

Information Stealing Attacks on Federated Learning

(e.g., Gradient Inversion, Gradient Amplifications, Trap Weights)

Wang et al., Beyond Inferring Class Representatives: User-Level Privacy Leakage From Federated Learning, 2019

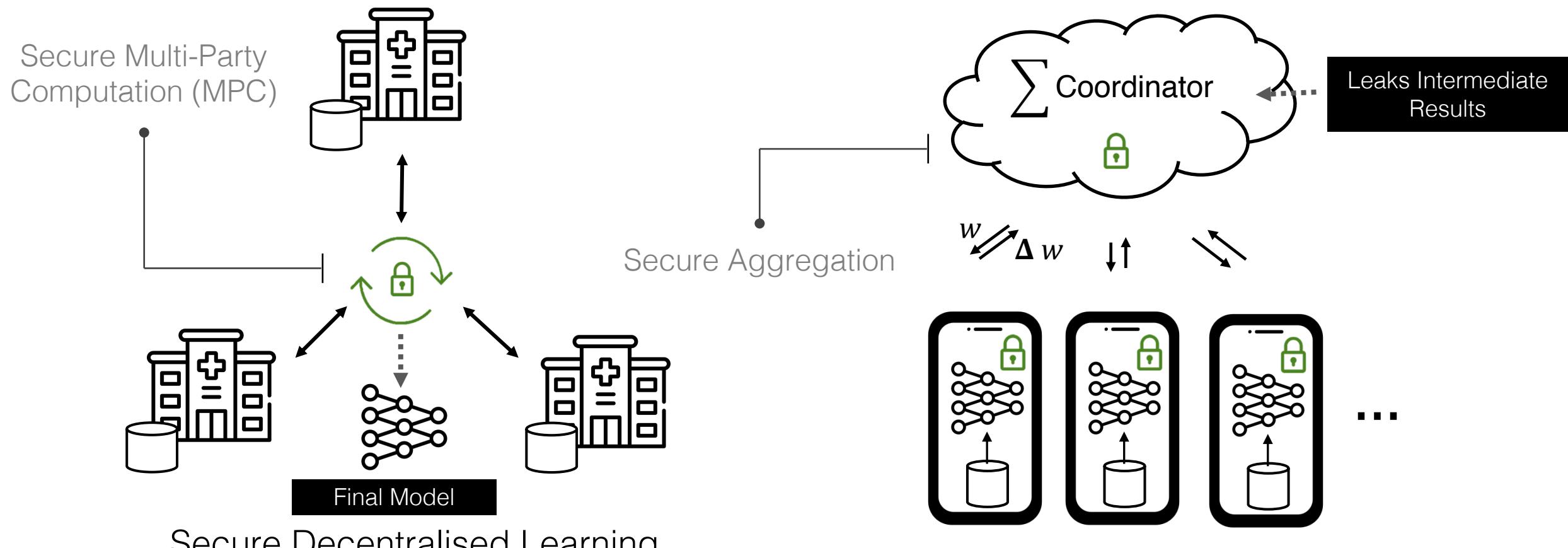
Geiping et al., How easy is it to break privacy in federated learning?, 2020

Boenisch et al., When the curious abandon honesty: Federated learning is not private, 2021

Yin et al., See through Gradients: Image Batch Recovery via GradInversion, 2021

Wen et al., Fishing for user data in large-batch federated learning via gradient magnification, 2022

Cryptography → Secure Computation



Secure Decentralised Learning

- CryptoNets [Gilad-Bachrach et al. ICML'16]
- SecureML [Mohassel et al. S&P'18]
- EzPC [Chandran et al. EuroS&P'19]
- Helen [Zheng et al. S&P'19]
- Spindle [Froelicher et al. PETS'20]
- Cerebro [Zheng et al. USENIX Security'21]

Secure Federated Learning

- Secure Aggregation [Bonawitz et al. CCS'17]
- FastSecAgg [Kadhe et al. CCS Workshop PPML'20]
- SecAgg+ [CCS'20]

Cryptography → Secure Computation



Secure Decentralised Learning

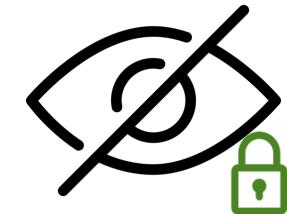
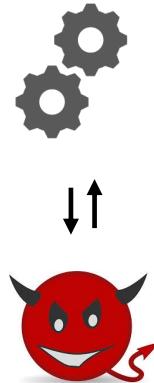
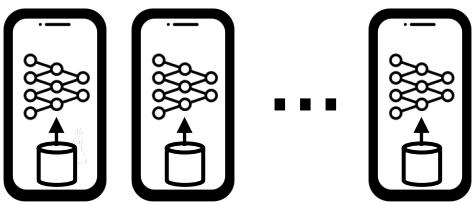
- CryptoNets [Gilad-Bachrach et al. ICML'16]
- SecureML [Mohassel et al. S&P'18]
- EzPC [Chandran et al. EuroS&P'19]
- Helen [Zheng et al. S&P'19]
- Spindle [Froelicher et al. PETS'20]
- Cerebro [Zheng et al. USENIX Security'21]

Secure Federated Learning

- Secure Aggregation [Bonawitz et al. CCS'17]
- FastSecAgg [Kadhe et al. CCS Workshop PPML'20]
- SecAgg+ [CCS'20]

Collaborative Learning

Can Amplify Robustness Issues

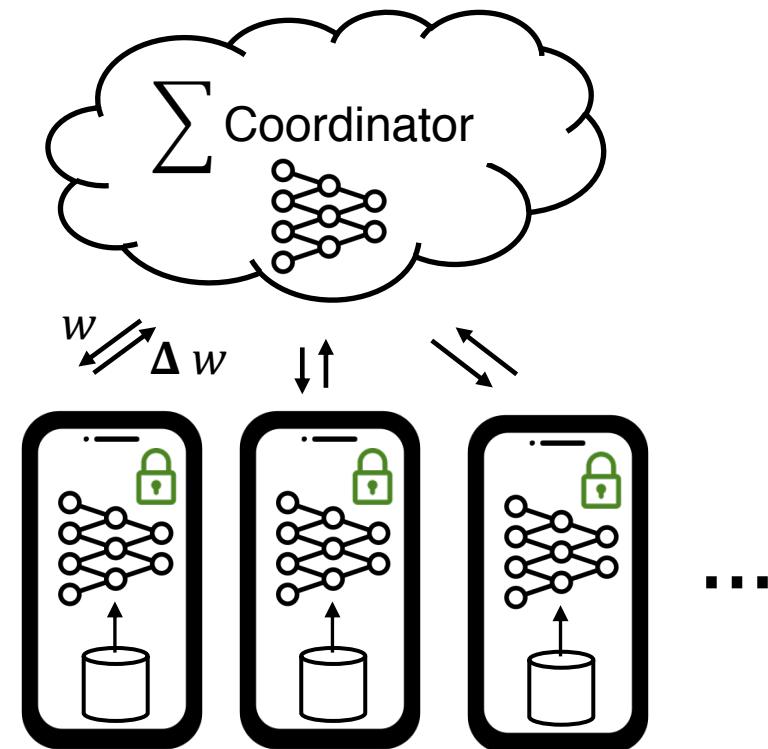


Open Nature

Attacker Capabilities

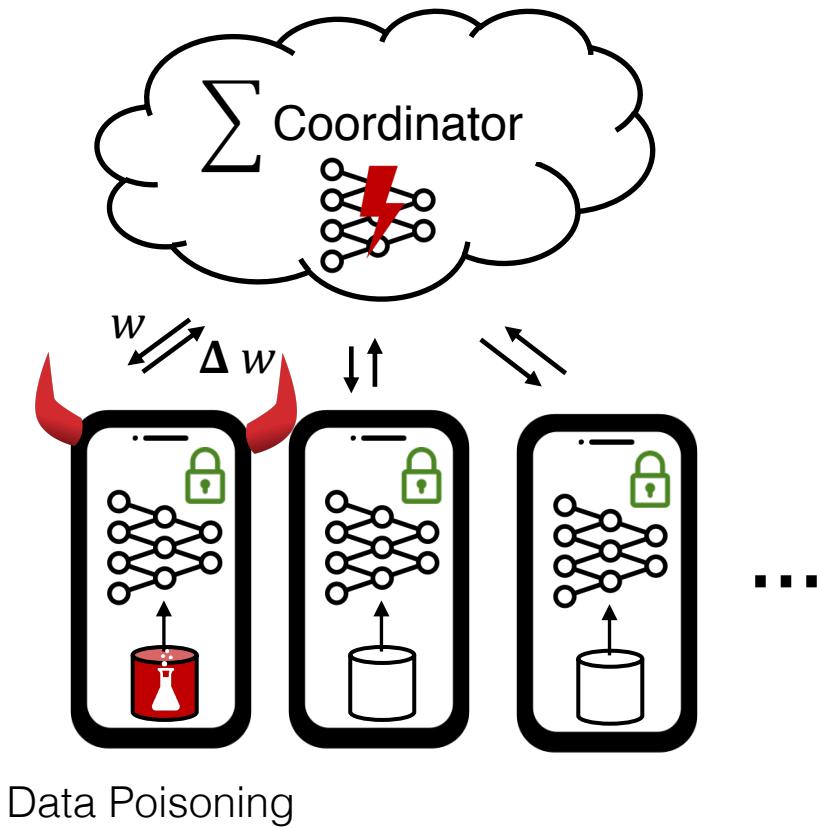
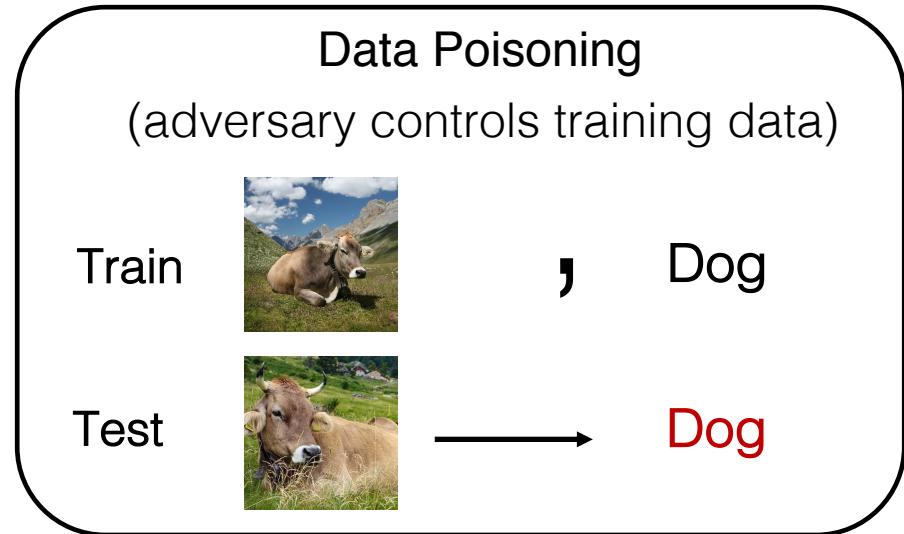
Detectability

Adversarial Robustness – Training



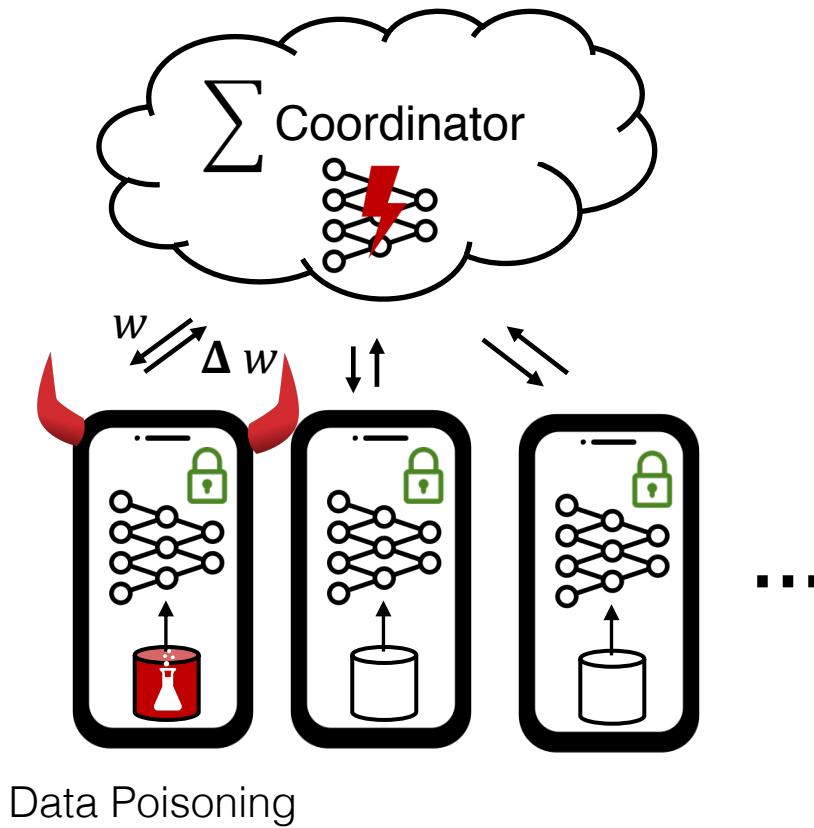
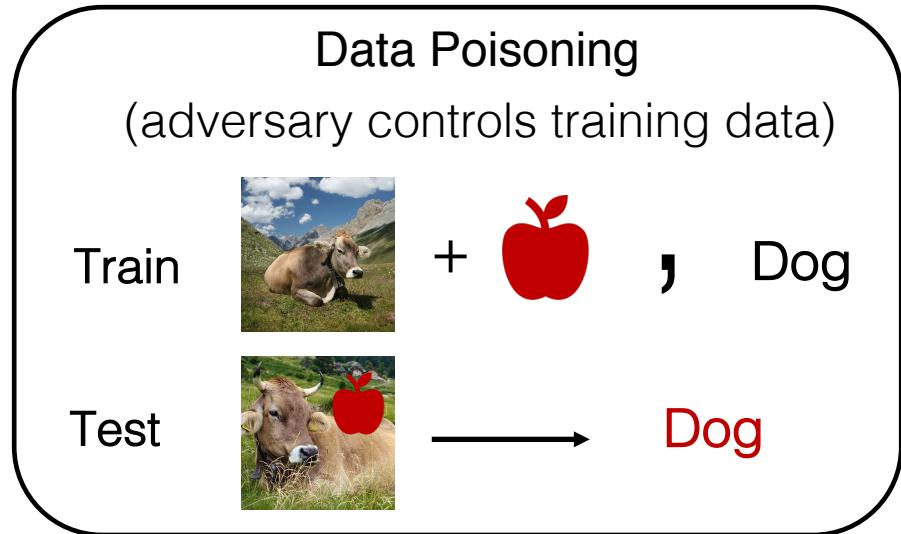
Secure Federated Learning

Adversarial Robustness – Training



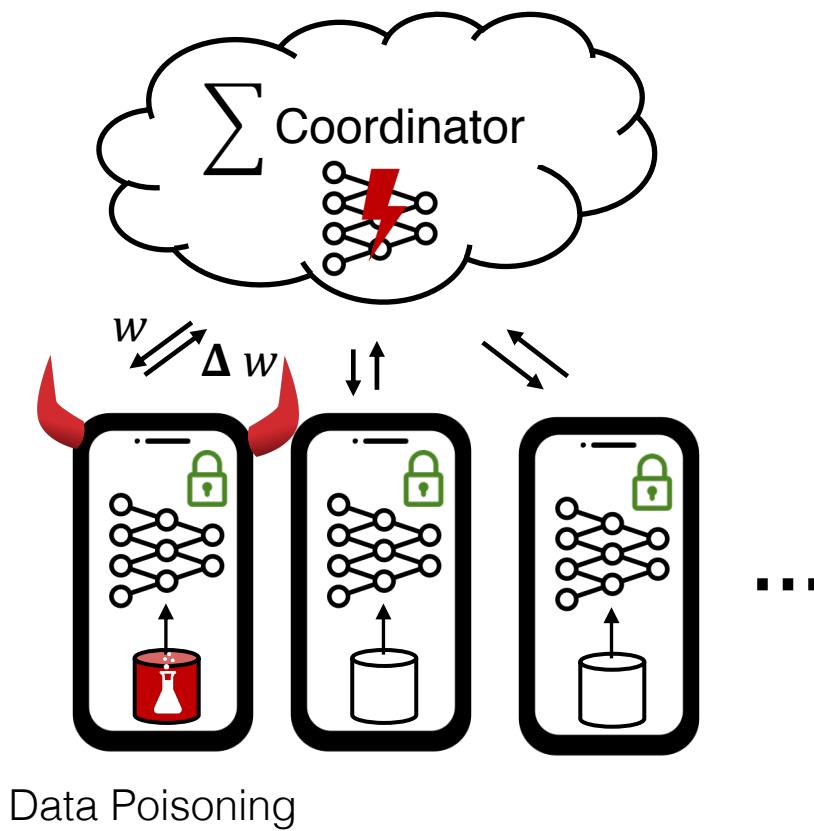
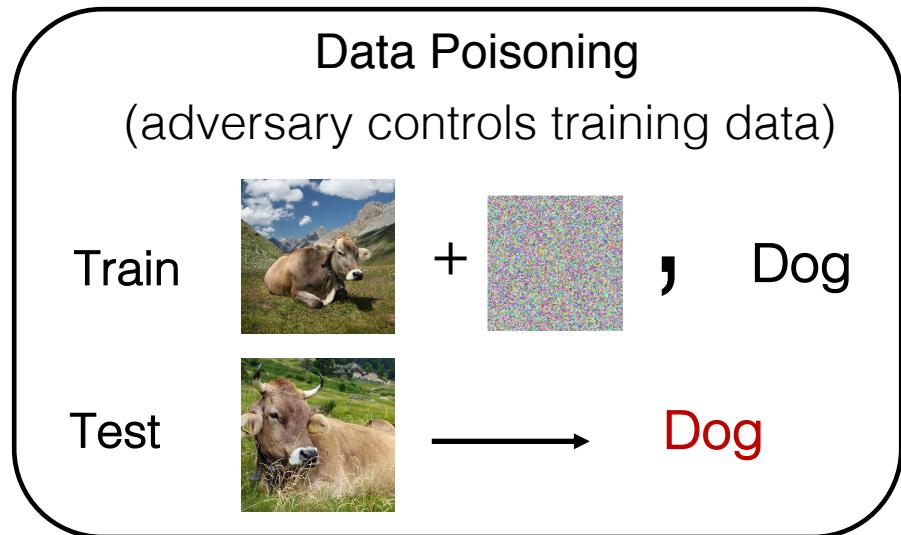
Secure Federated Learning

Adversarial Robustness – Training



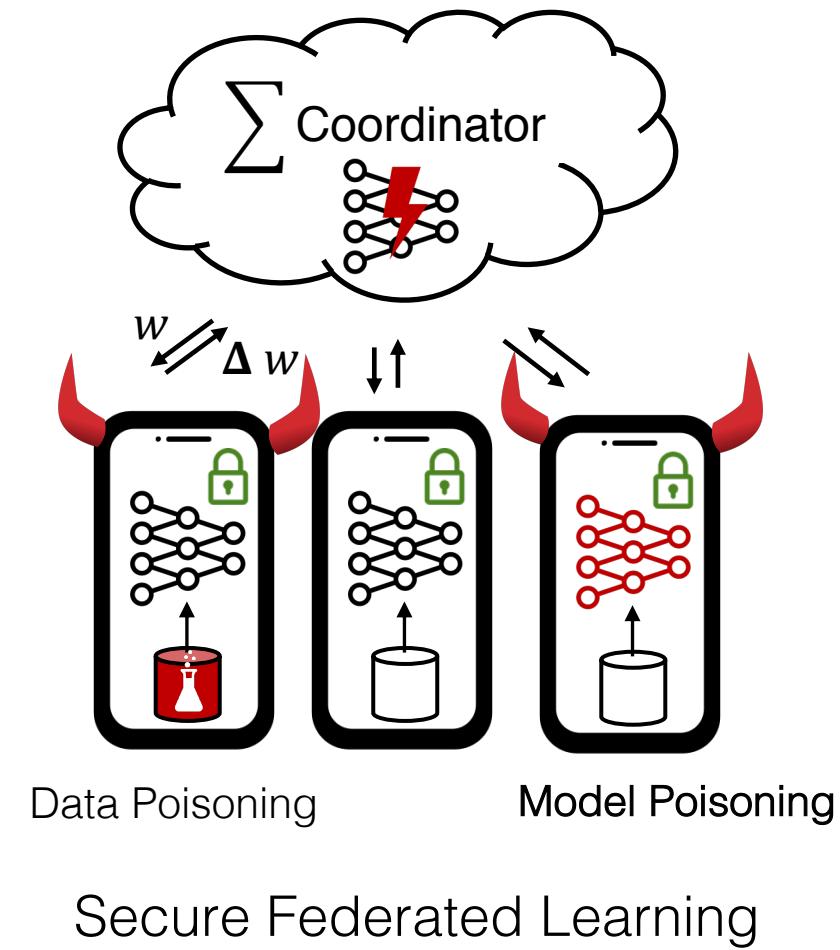
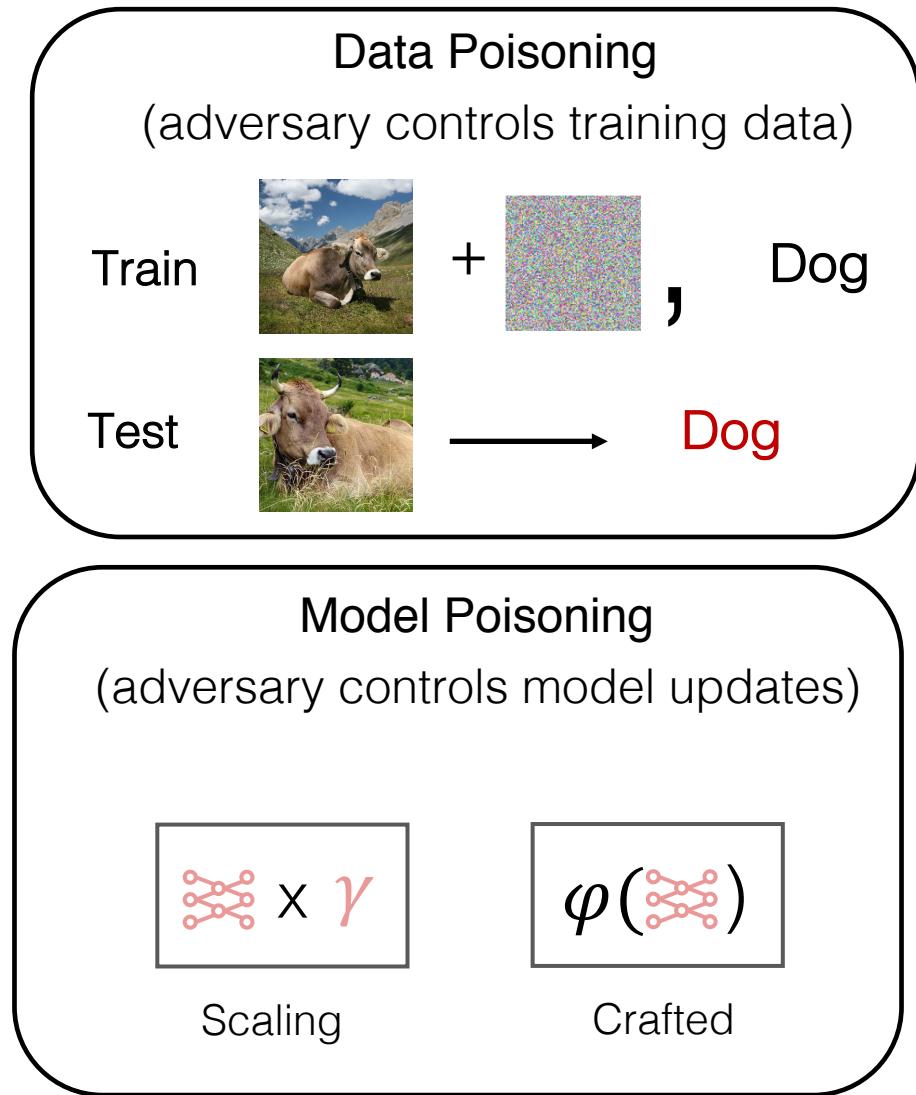
Secure Federated Learning

Adversarial Robustness – Training

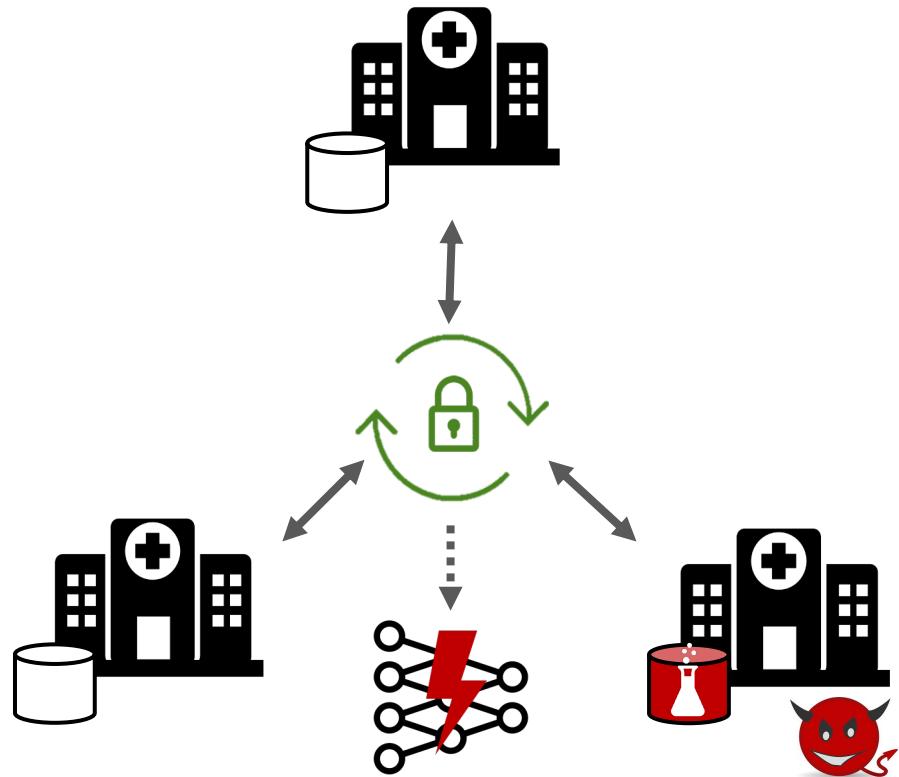


Secure Federated Learning

Adversarial Robustness – Training



Adversarial Robustness – Training



Data Poisoning

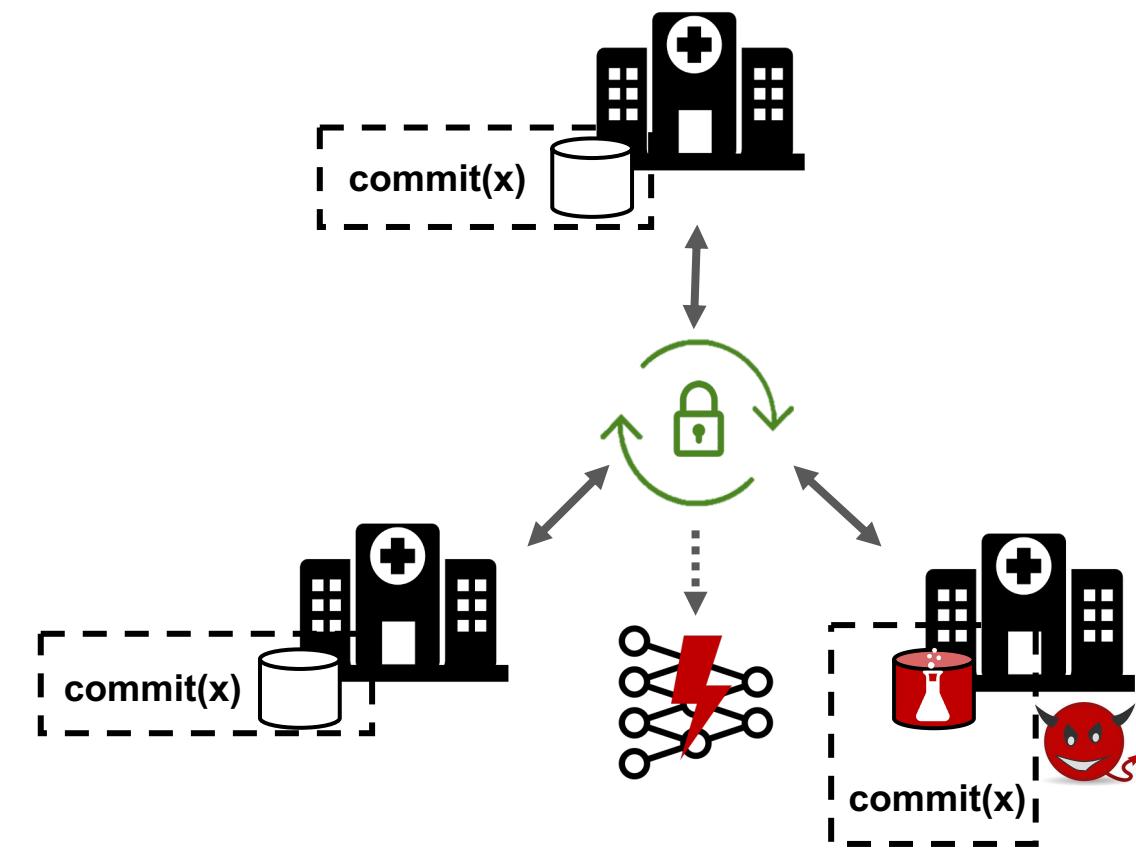
(adversary controls training data)

Model Poisoning

(adversary controls model updates)

Secure Decentralized Learning

Adversarial Robustness – Training



Secure Decentralized Learning

Data Poisoning

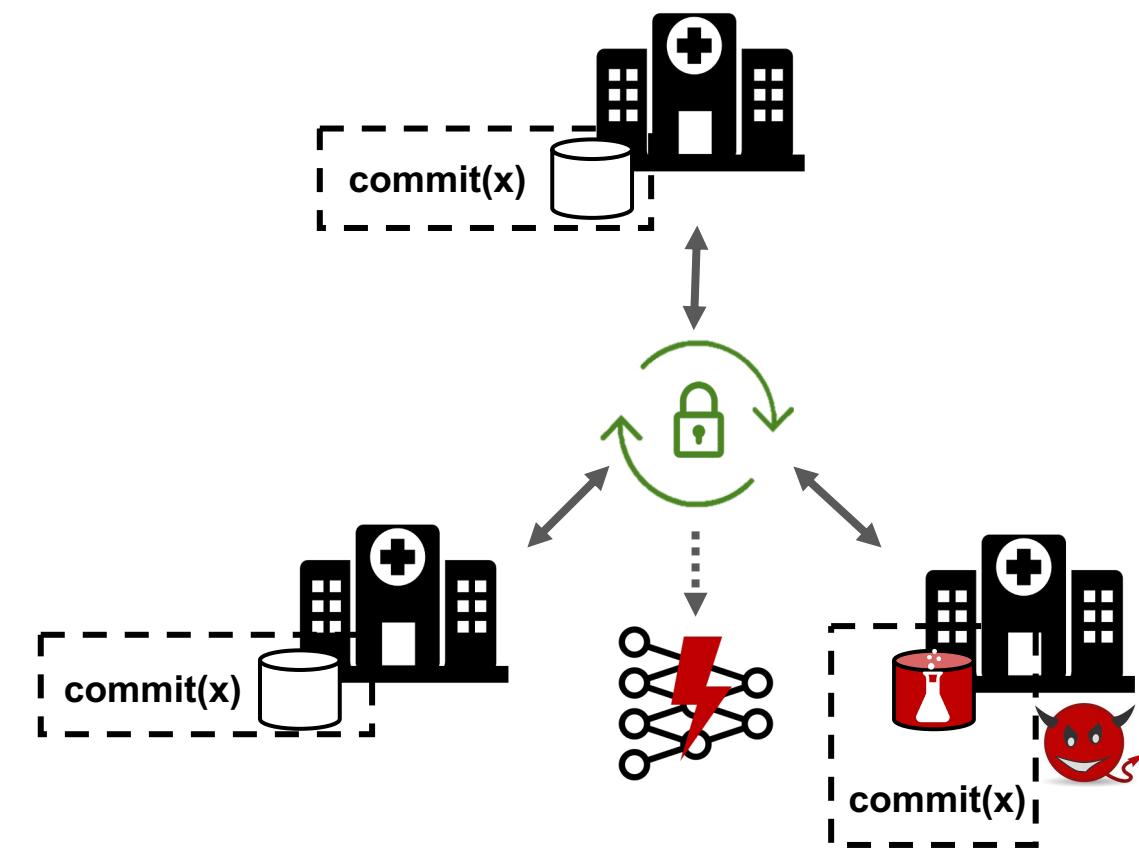
(adversary controls training data)

Model Poisoning

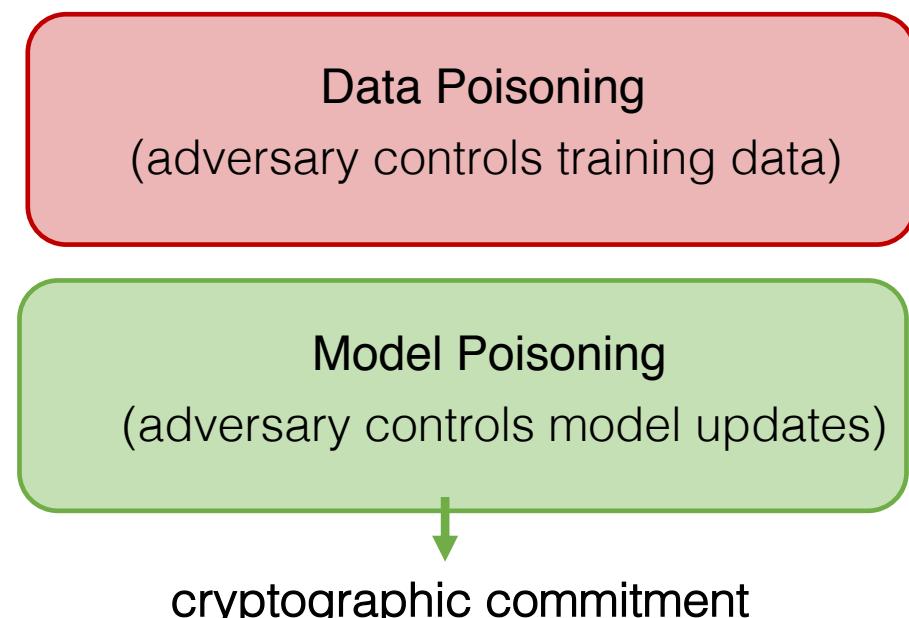
(adversary controls model updates)

cryptographic commitment

Adversarial Robustness – Training



Secure Decentralized Learning



Robust ML Algorithms

cryptography-friendly
algorithms

Detection Mechanisms

assumes direct access to the
data or the gradients

Cryptography

?

Cryptographic Verification

Zero-knowledge proofs, Cryptographic commitments, Proofs for program delegations, ...

Conventional Setting

Verify some pre-specified function f

Given $P(x)$

-- Verify: $P(x) = f(x)$

Machine Learning Setting

In ML f is learned

(f = unknown ground truth)

Given $P(x)$

-- Verify what then?

The source of the issue is maliciously chosen data

→ alteration, proof/verify **something** about the input data, gradients, or data distribution

- Theoretical work: Verify data distribution (in/out/adversarial)
- Enforce constraints on the gradient updates (e.g., norm bound)
- Verify Source of Data
- ...

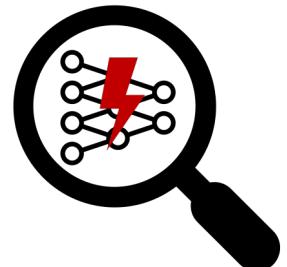
Overview Wrap Up

- Decoupling data from training, by itself, does not provide many privacy benefits
 - Encryption can help (e.g., secure aggregation, MPC)
- More work on robust ML in the **encrypted settings**
 - Cryptography friendly robust ML algorithms
 - Use cryptography (e.g., verification, ZKP) to minimize influence of maliciously chosen training data
- Post-Deployment
 - Can we get robustness against all attacks? **Answer:** A perfect solution to adversarial robustness remains an open challenge – imperfect defenses, cat-and-mouse game, more powerful attacks
 - There is a need for solutions that minimize consequences of attacks at deployment time – e.g., attribution, forensics, accountability, audits, admission controls, monitoring ...

RoFL: Attestable Robustness for Secure FL

Lukas Burkhalter*, Hidde Lycklama*, Nicolas Küchler, Alexander Viand, Anwar Hithnawi

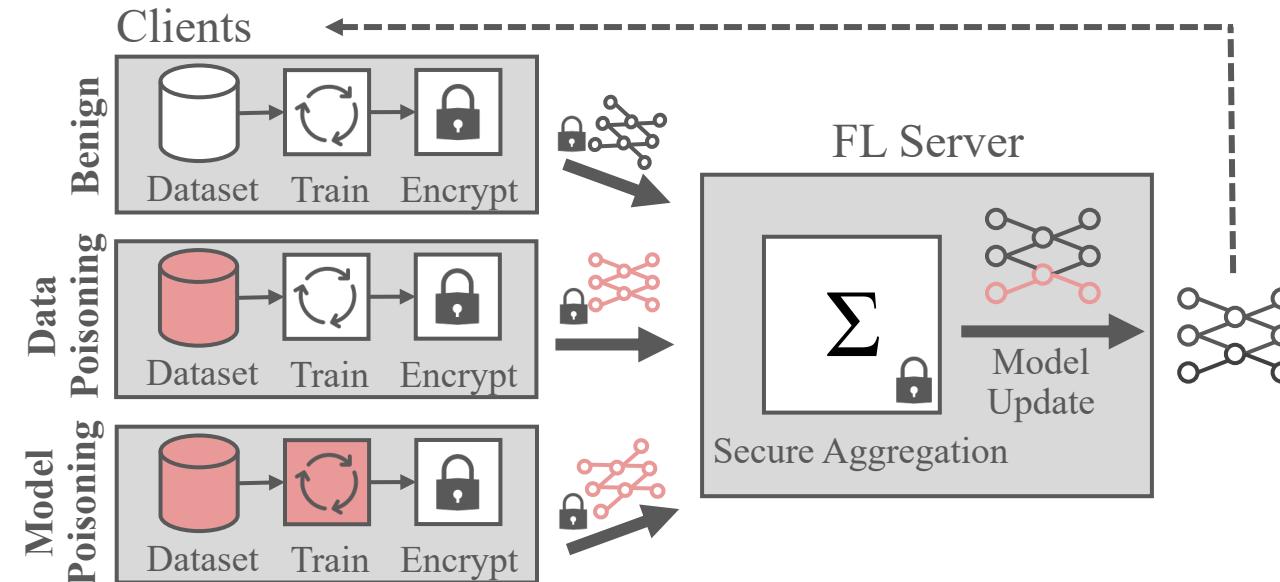
Understand
Vulnerabilities in FL



Cryptographically
Enforce Constraints



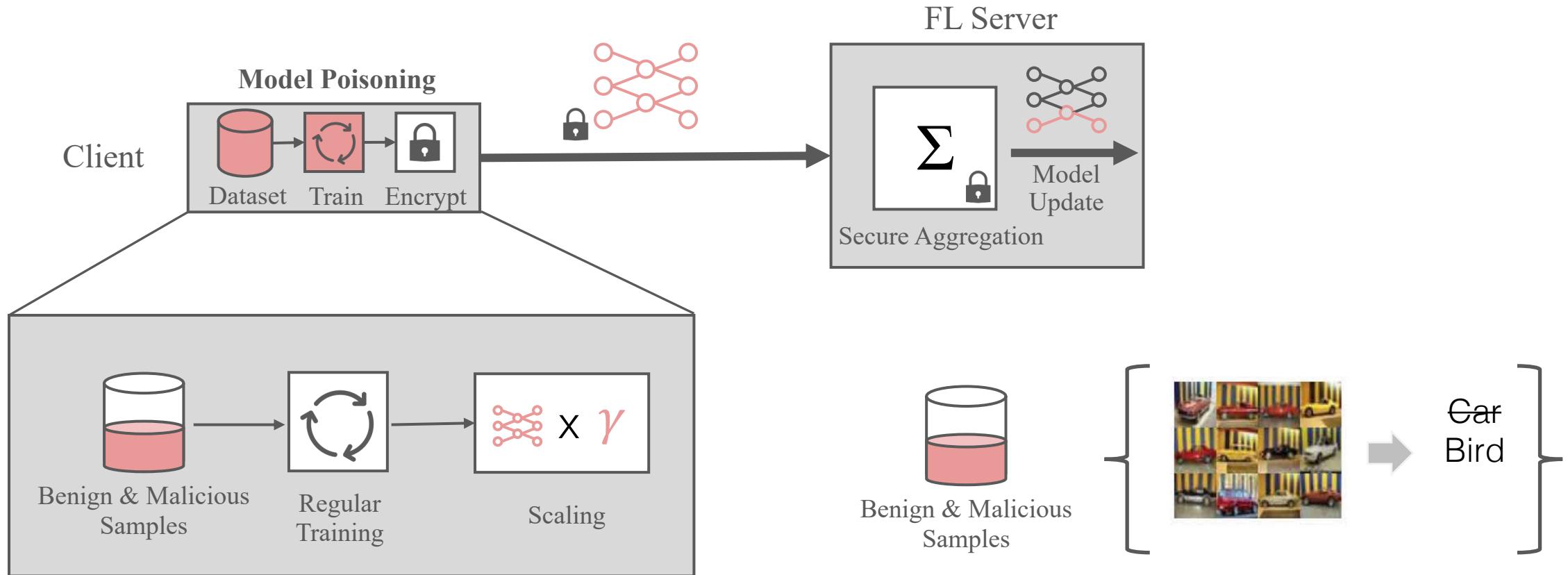
Secure Federated Learning



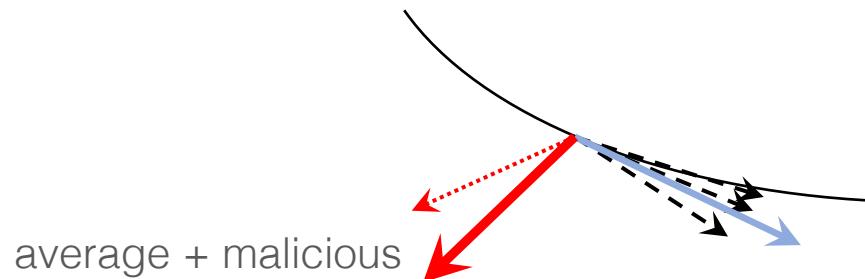
What are the vulnerabilities in the FL pipeline
that enable model/data poisoning attacks



Model Replacement Attack



Problem: Linear aggregation rules are vulnerable to Byzantine behavior



Machine Learning: Byzantine-Robust Distributed Learning

- Krum [Blanchard et al. NeurIPS'17]
- Trimmed Mean [Yin et al. ICML'18]
- Coordinate-wise Median [Yin et al. ICML'18]
- Bulyan [Mhamdi et al. ICML'18]
- ByzantineSGD [Alistarh et al. NeurIPS'18]
- Redundant Workers and Coding Theory [Chen et al. ICML'18]
- [Rajput et al. NeurIPS'19]

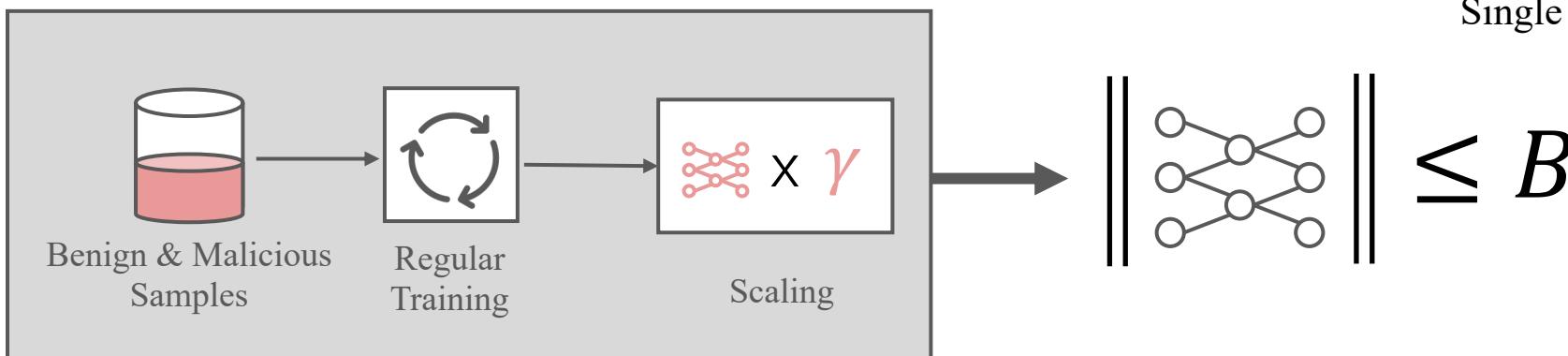
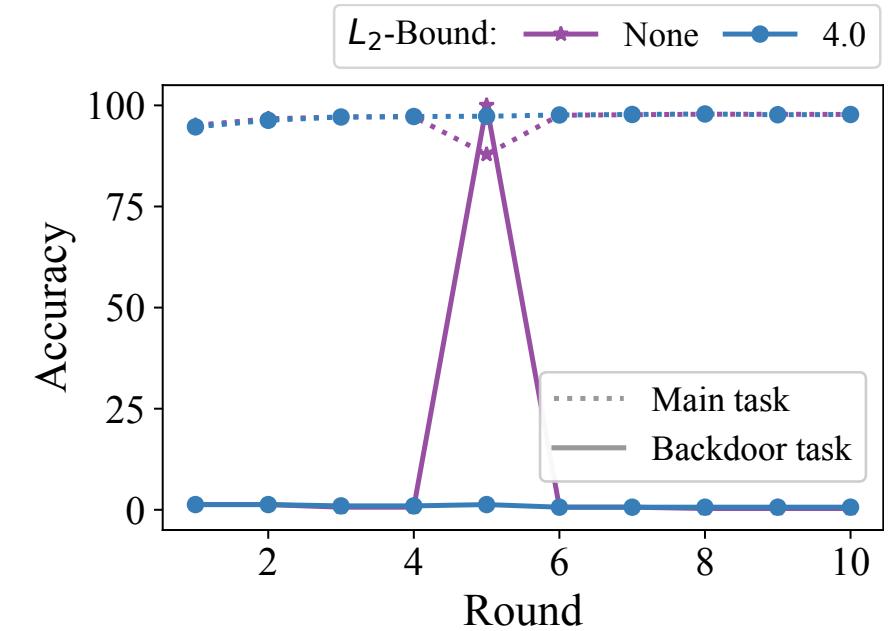
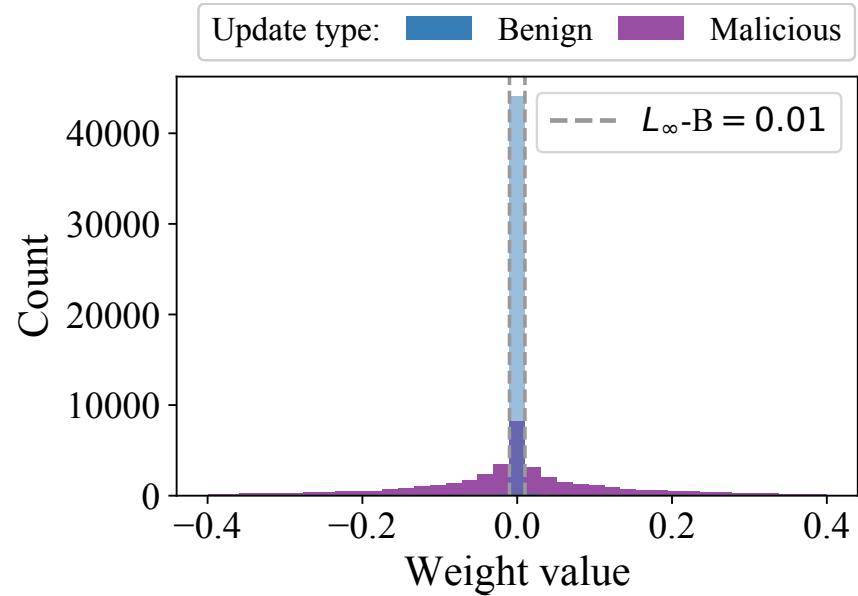
Security: Private Data-Collection Systems

- Prio [Corrigan-Gibbs et al. NSDI'17]
- PrivStats [Popa et al. CCS'11]
- SplitX [Chen et al. SIGCOMM'13]
- P4P [Duan et al. USENIX Security'10]
- PrivEx [Elahi et al. CCS'14]

→ Zero Knowledge Proofs: client proves that its submission is well-formed

A well-formed Client Submission in Federated Learning

Norm bound



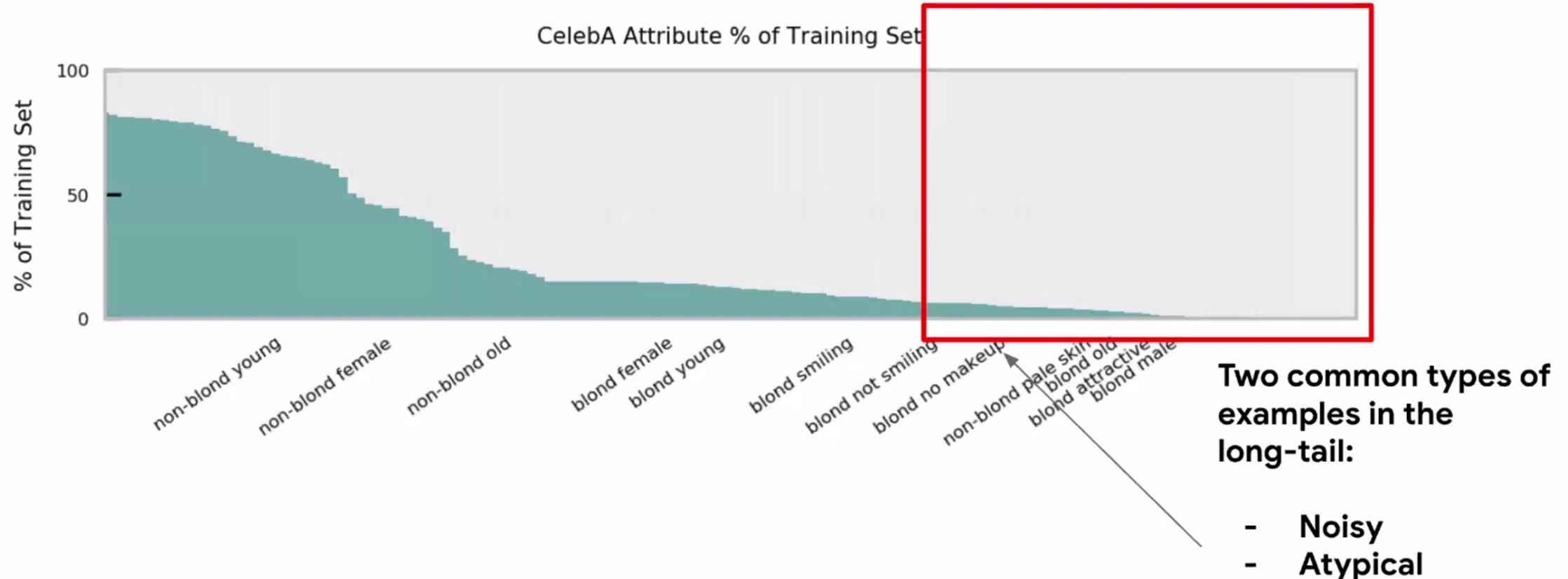
How To Backdoor Federated Learning

Can You Really Backdoor Federated Learning?

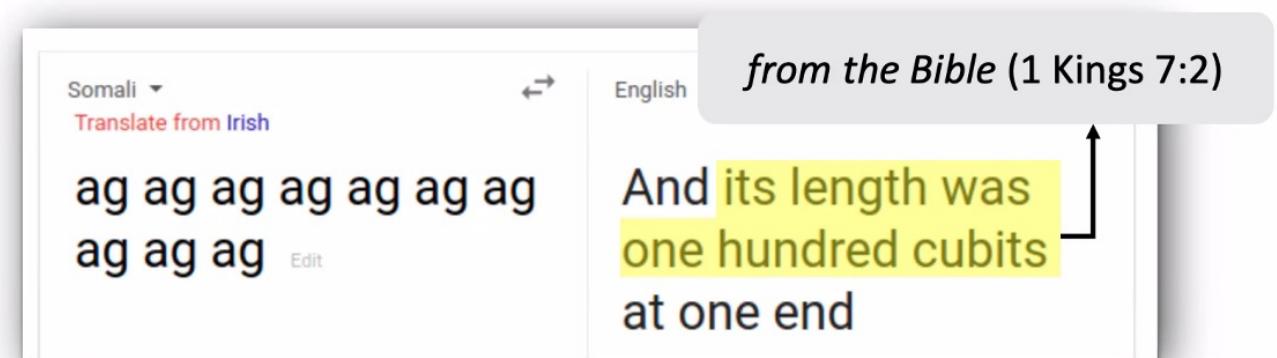
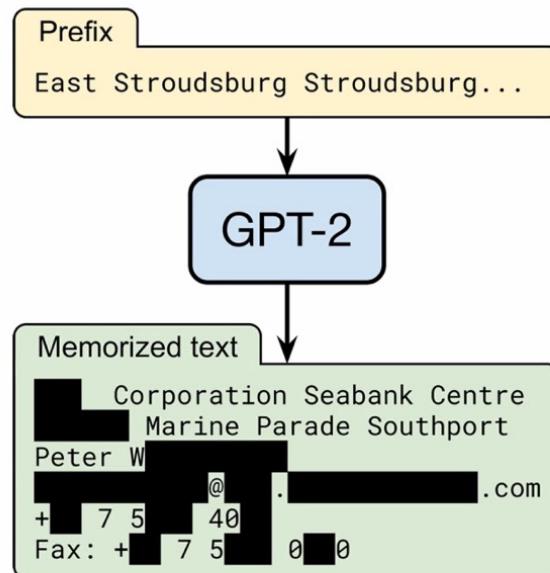
**Attack of the Tails:
Yes, You Really Can Backdoor Federated Learning**

Why?

Long Tail ...



Model Capacity Implications on Privacy ...



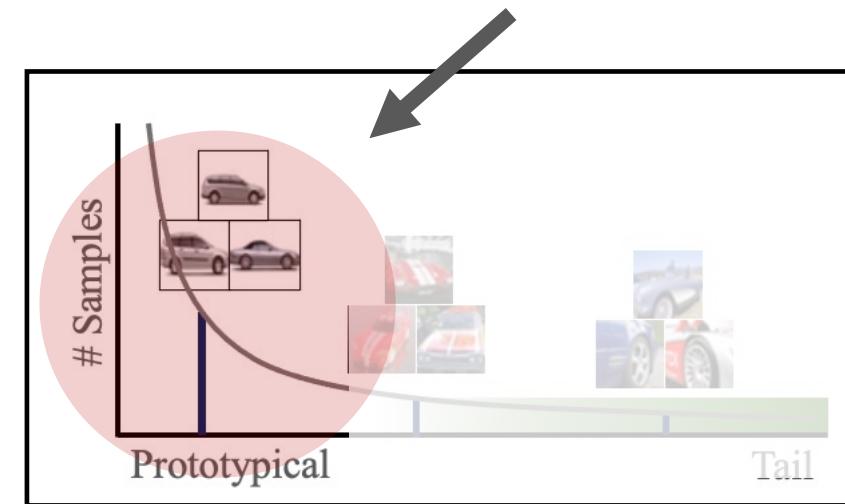
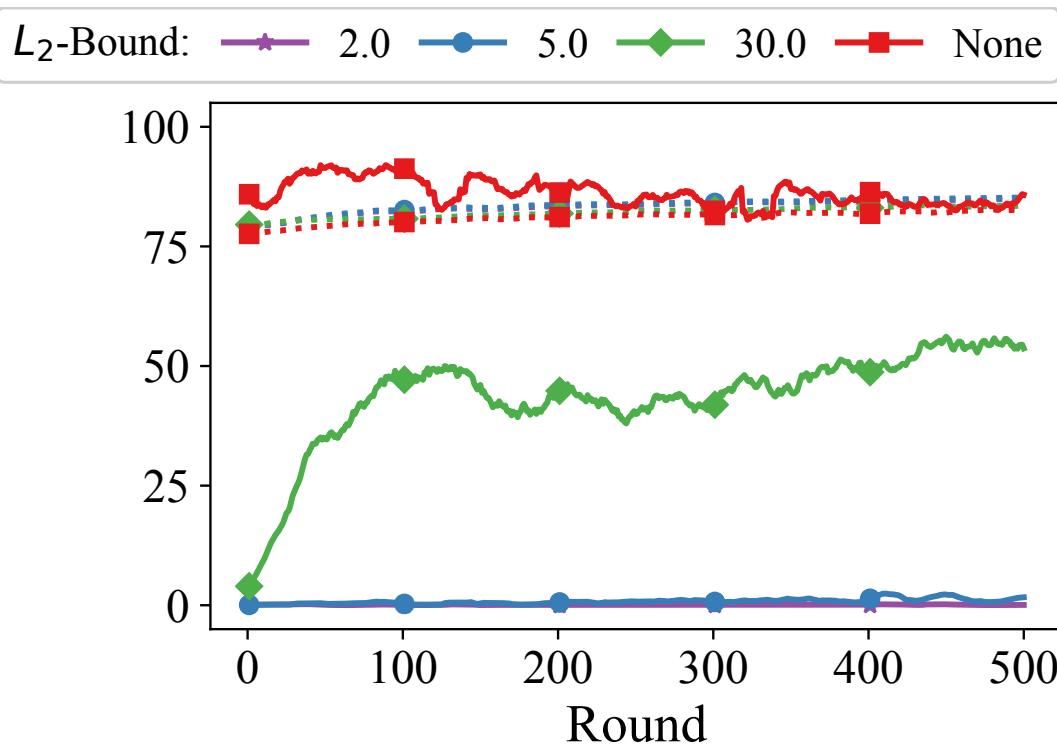
Memorization

Memorization leads to
Leakage of private text

Model Capacity Implications on Robustness...

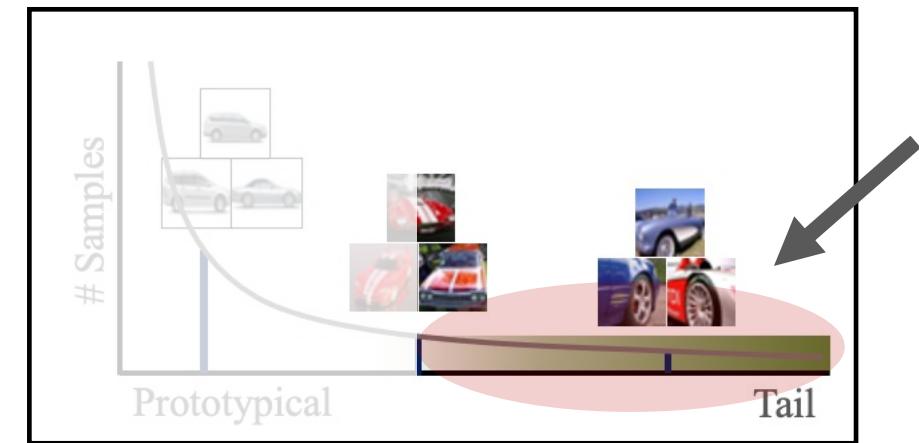
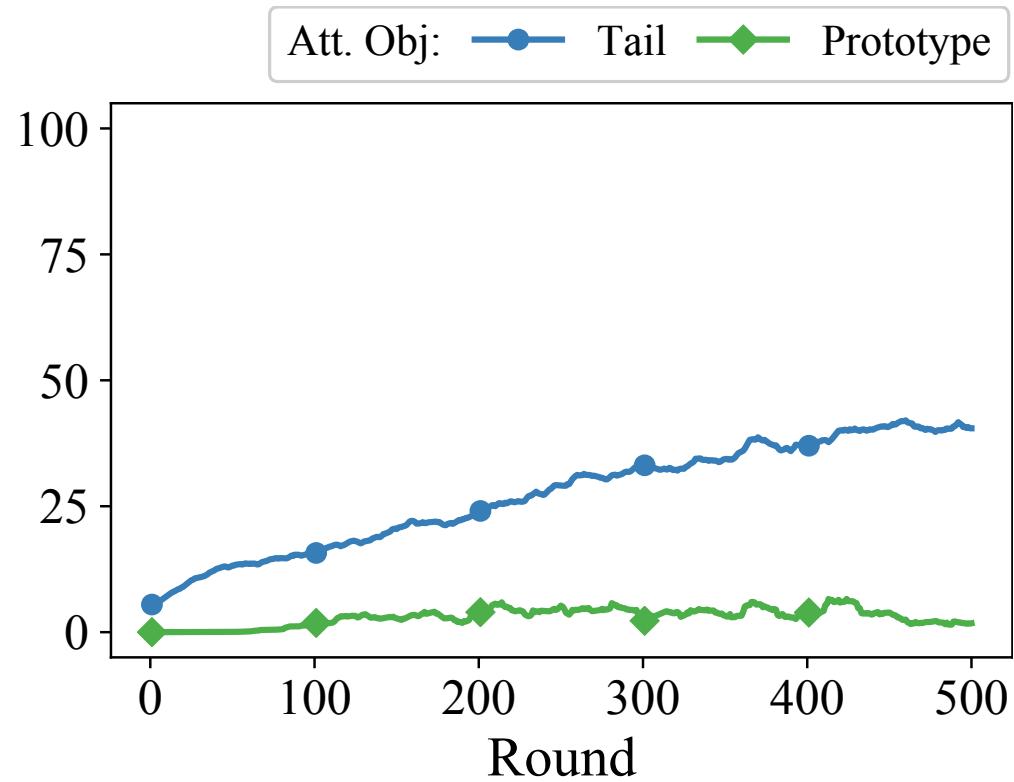
Success of Backdoor Attacks

Prototypical Targets

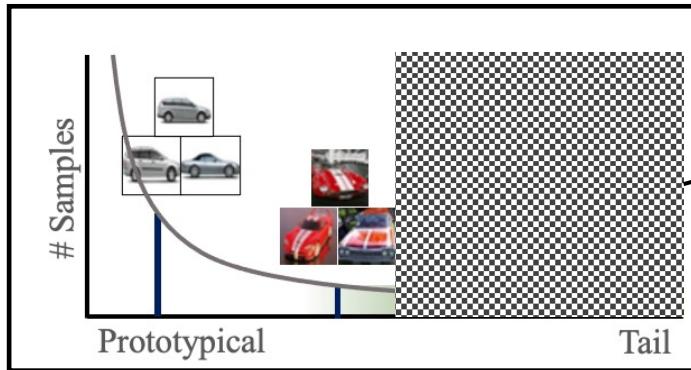


Success of Backdoor Attacks

Tail Targets



Suppressing the Long-Tail



Approaches

- Noise Addition (Differential Privacy)
- Compression

Understand trade-offs between objectives we care about



Robustness



Accuracy



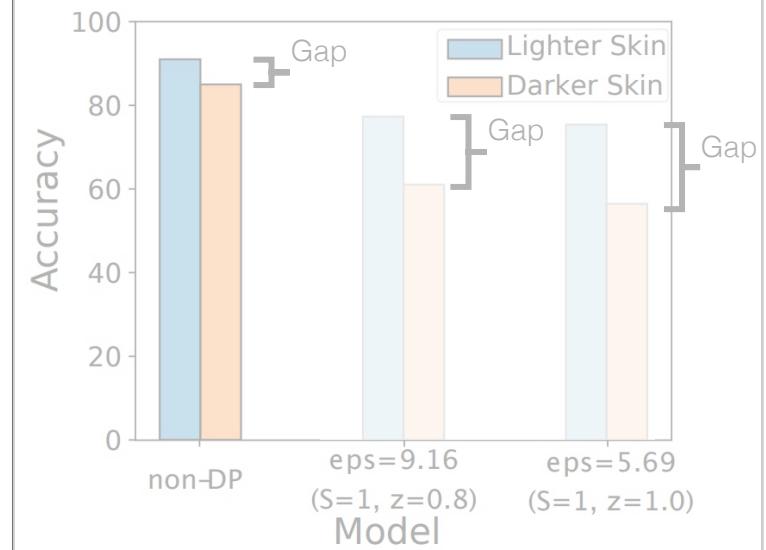
Fairness



Privacy

Leads to Fairness Problems

Differential Privacy disproportionately impacts underrepresented attributes
[Bagdasaryan et al. NeurIPS 2019]



More Resources ...

“

Understanding how capacity impacts (fairness, robustness, privacy) is an increasingly urgent question.

”

-- Sarah Hooker



In the Talk

The myth of interpretable, robust, compact and high performance DNNs

“

Understanding the generalization properties of learning systems (...) is an area of great practical importance.

”

-- Vitaly Feldman



In the Paper

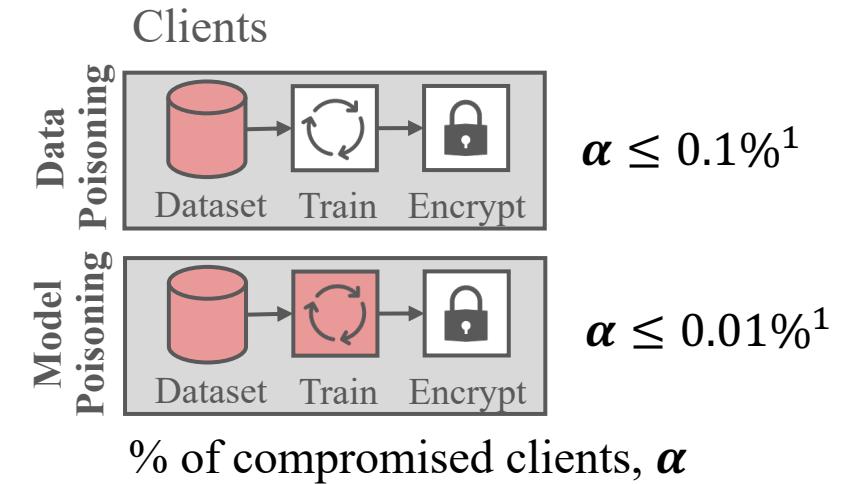
Does Learning Require Memorization? A Short Tale about a Long Tail

Binary View of Robustness

Where can Norm Bound Help?

Attack Type / Attack Target	Prototypical	Tail
Single-shot	✓	✓
Continuous	✓	○

Requires attacker to be consistently selected $\alpha > 2.5\%$



Norm Bound Provides Practical Robustness Guarantees

RoFL: Attestable Robustness for Secure FL

Lukas Burkhalter*, Hidde Lycklama*, Nicolas Küchler, Alexander Viand, Anwar Hithnawi

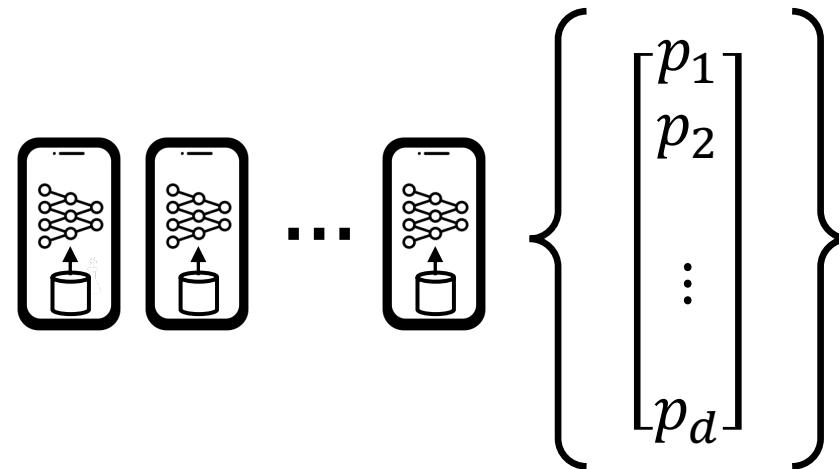
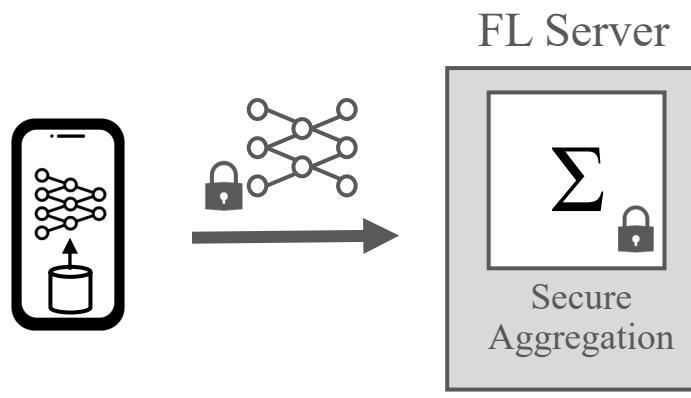
Understand
Vulnerabilities in FL



Cryptographically
Enforce Constraints



Goal: Augment existing secure FL with Zero-Knowledge Proofs to enforce constraints on model updates

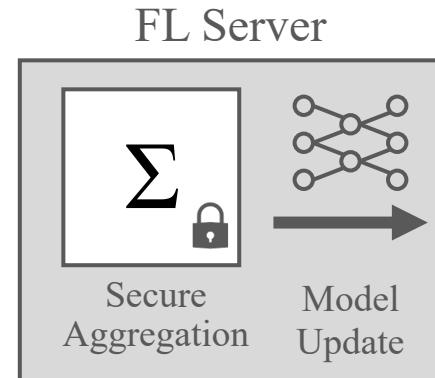
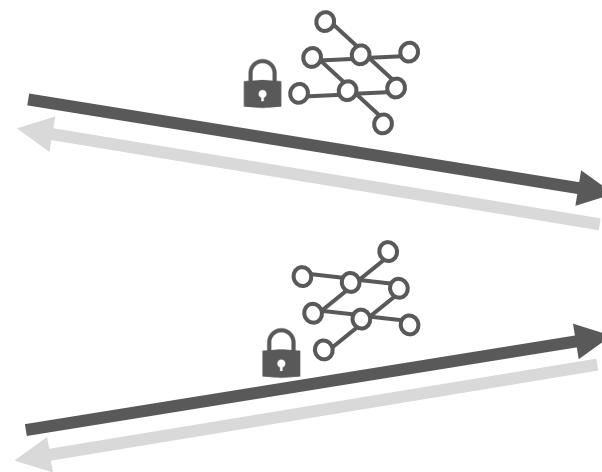
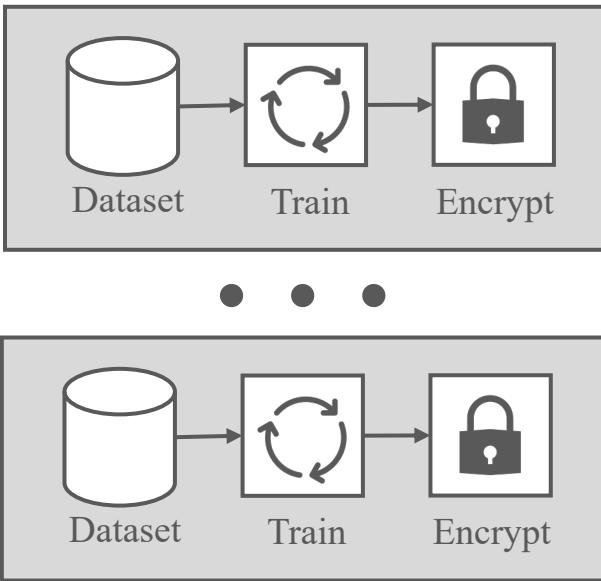


Compatibility

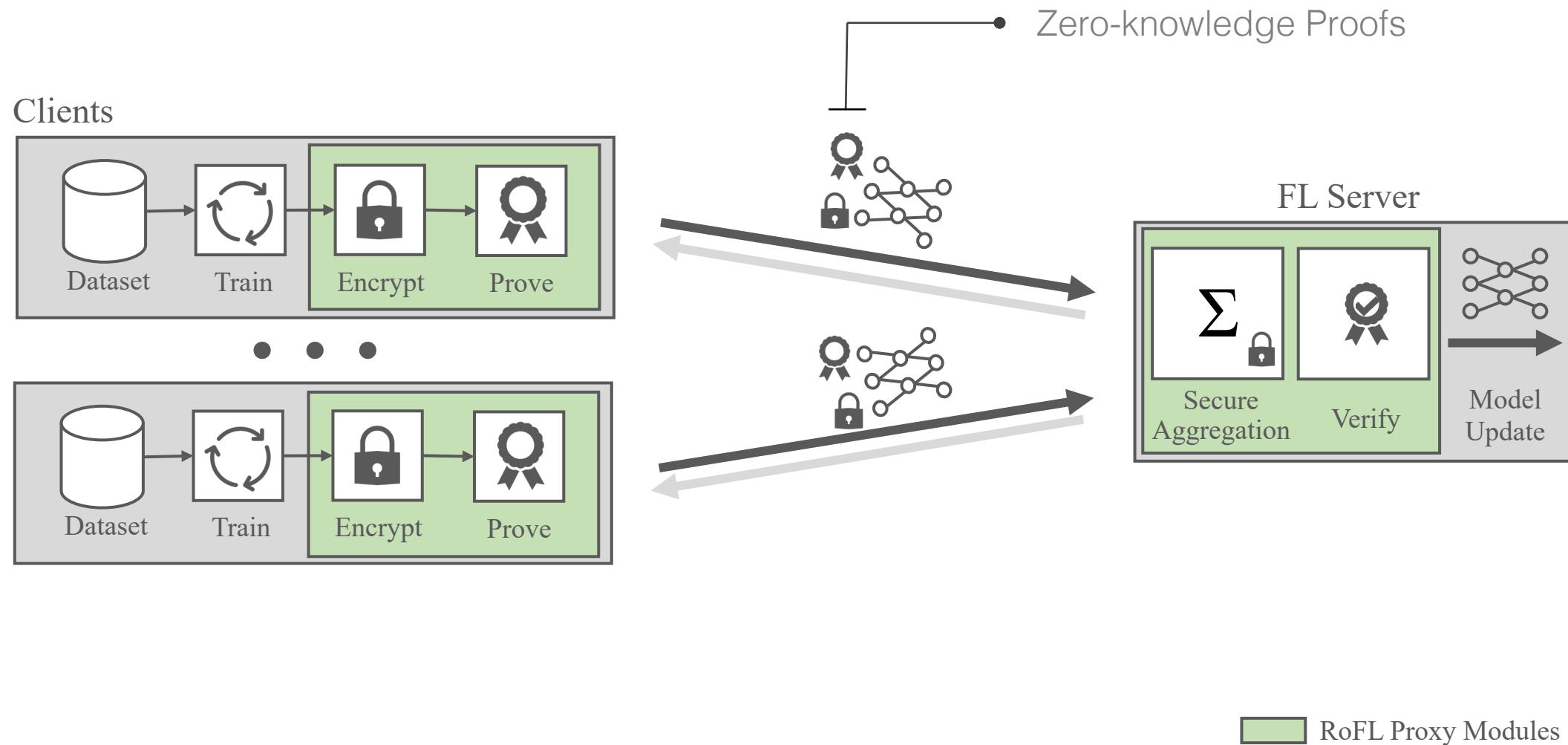
Scalability

Secure Federated Learning

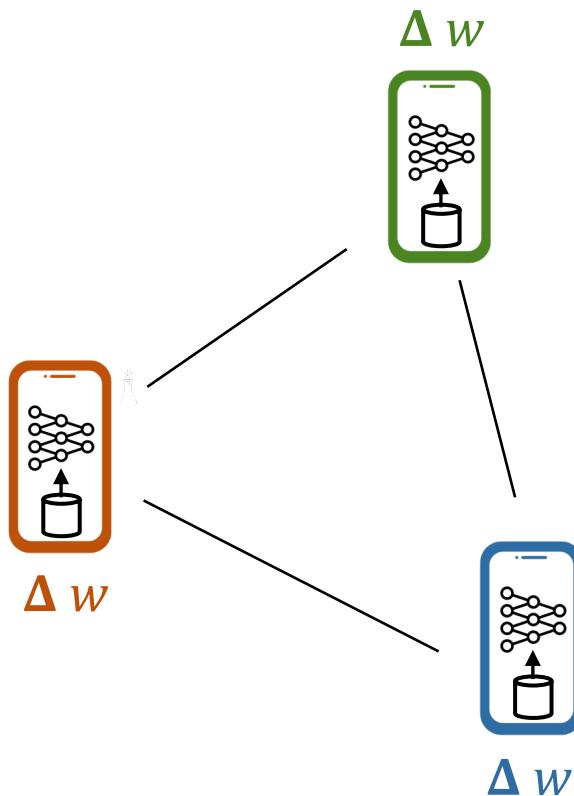
Clients



RoFL Augments Secure FL



Secure Aggregation



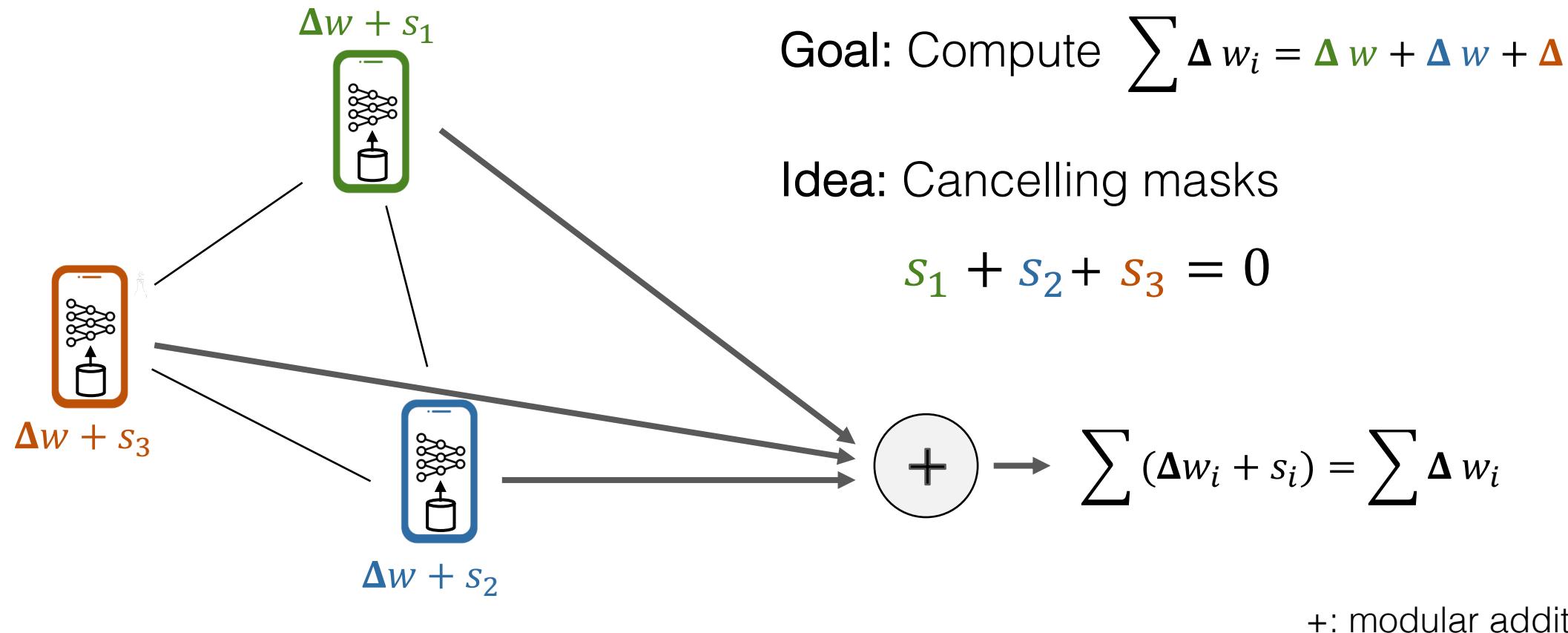
Goal: Compute $\sum \Delta w_i = \Delta w + \Delta w + \Delta w$

Idea: Cancelling masks

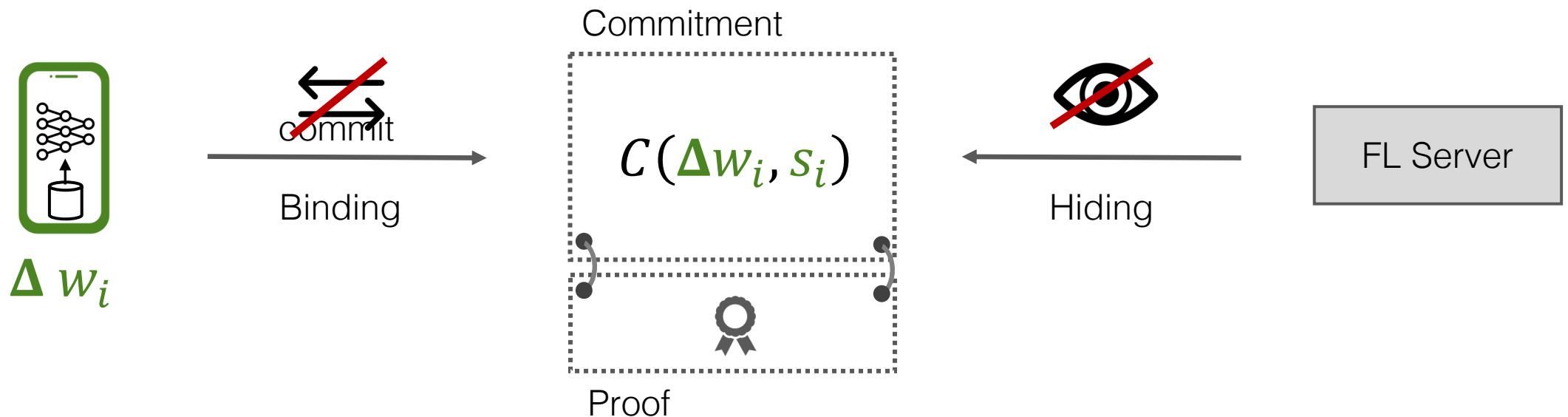
$$s_1 + s_2 + s_3 = 0$$

+: modular addition

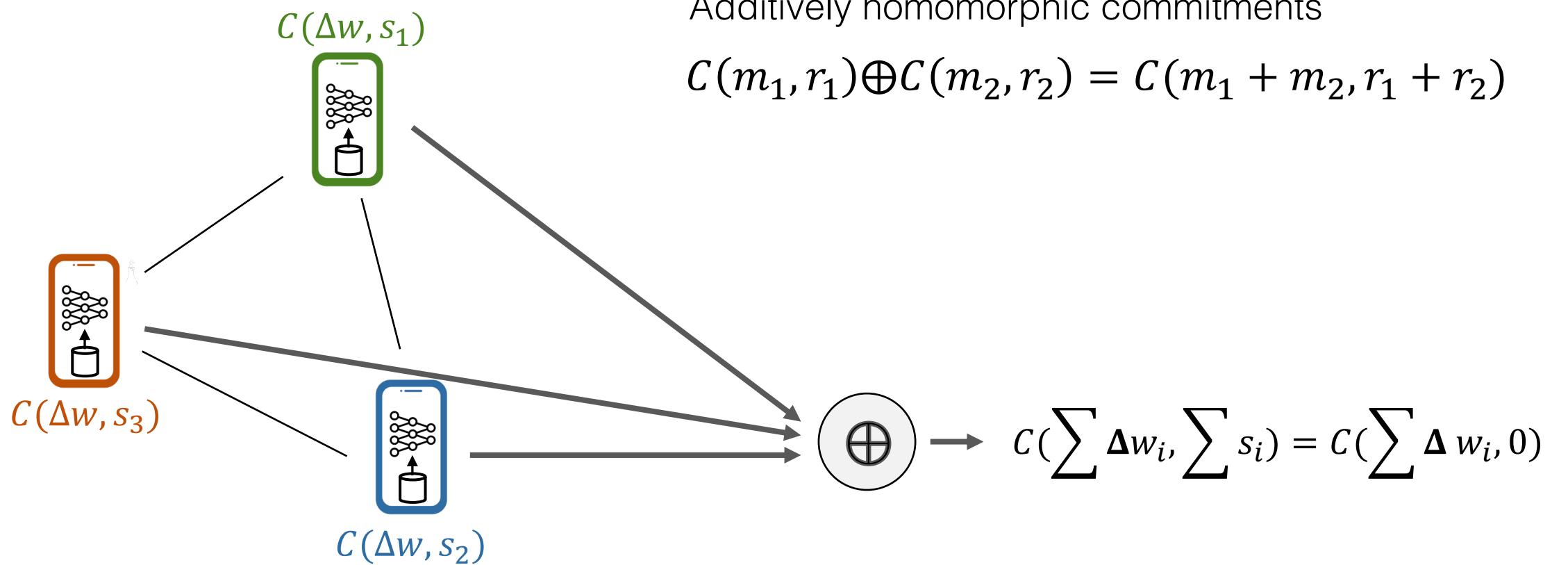
Secure Aggregation



Cryptographic Commitments



Switching to Homomorphic Commitments



Zero-knowledge Proofs for Norm Constraints

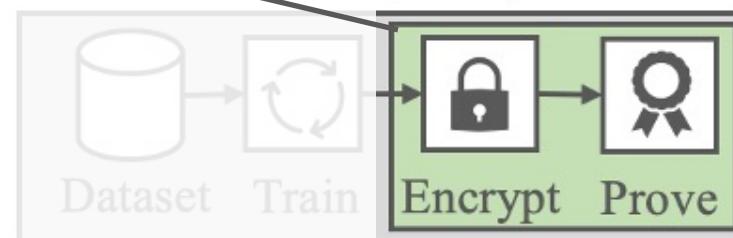
$$\begin{bmatrix} C(p_1, r_1) \\ C(p_2, r_2) \\ \vdots \\ C(p_d, r_d) \end{bmatrix}$$

ElGamal commitments



L_∞ -, L_2 -norm

Non-Interactive Zero-Knowledge
Proofs



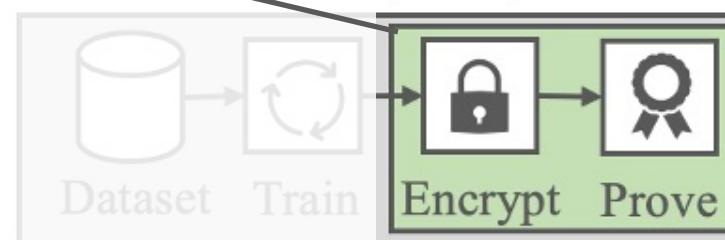
Enforcing L_∞ -norm

$$\begin{bmatrix} C(p_1, r_1) \\ C(p_2, r_2) \\ \vdots \\ C(p_d, r_d) \end{bmatrix}$$

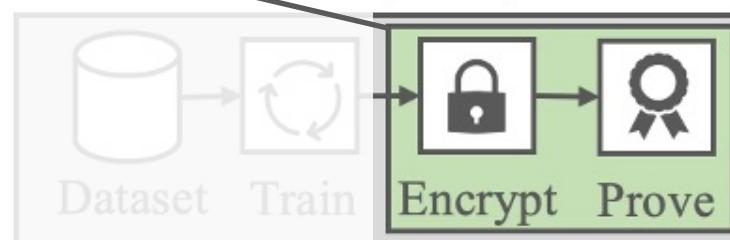
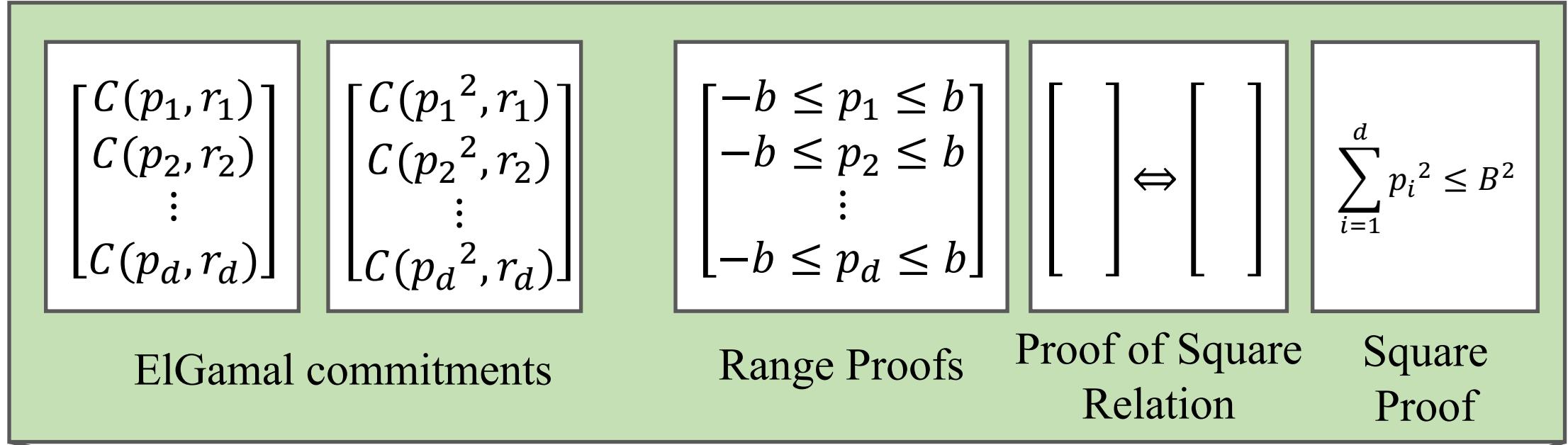
ElGamal commitments

$$\begin{bmatrix} -b \leq p_1 \leq b \\ -b \leq p_2 \leq b \\ \vdots \\ -b \leq p_d \leq b \end{bmatrix}$$

Bulletproof Range Proofs



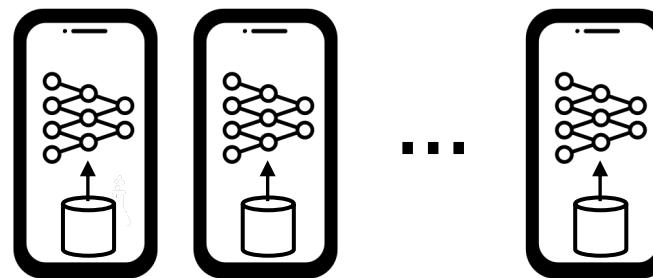
Enforcing L_2 -norm



Problem: Scalability

$$\begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_d \end{bmatrix} \quad \left. \right\} > 100k$$

High-dimensional updates

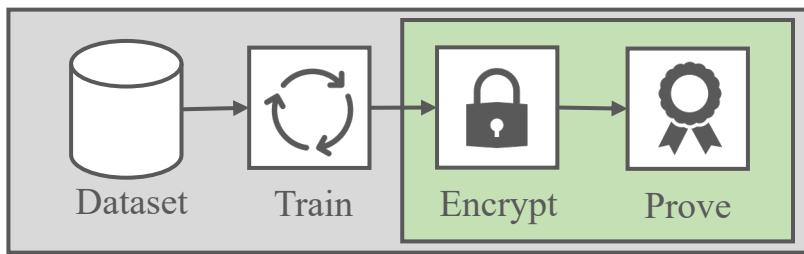


Number of clients

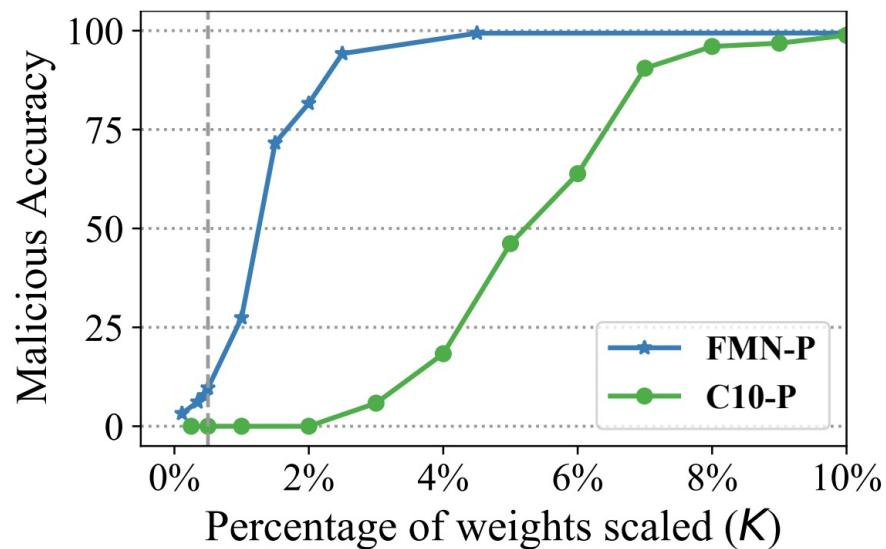
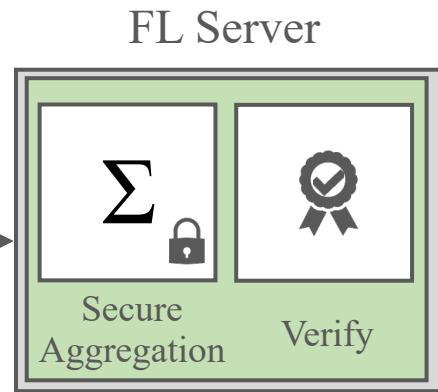
Can we reduce the number of proofs while maintaining
the same level of security?

Optimizing L_∞

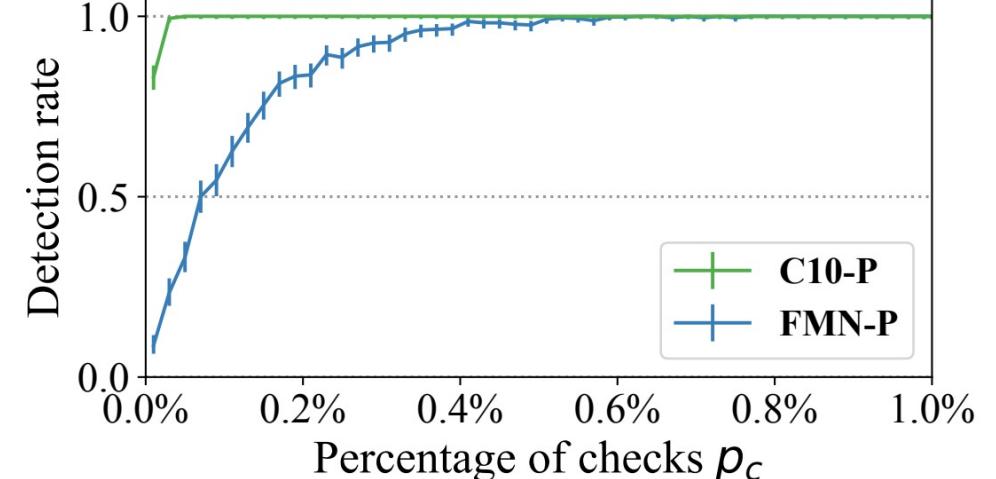
Clients



Commitments Δw , range proofs 🏅

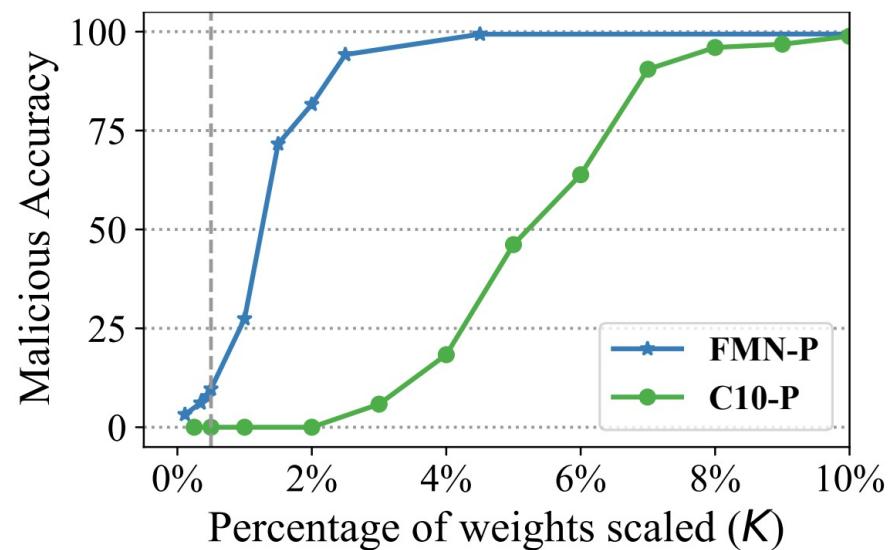
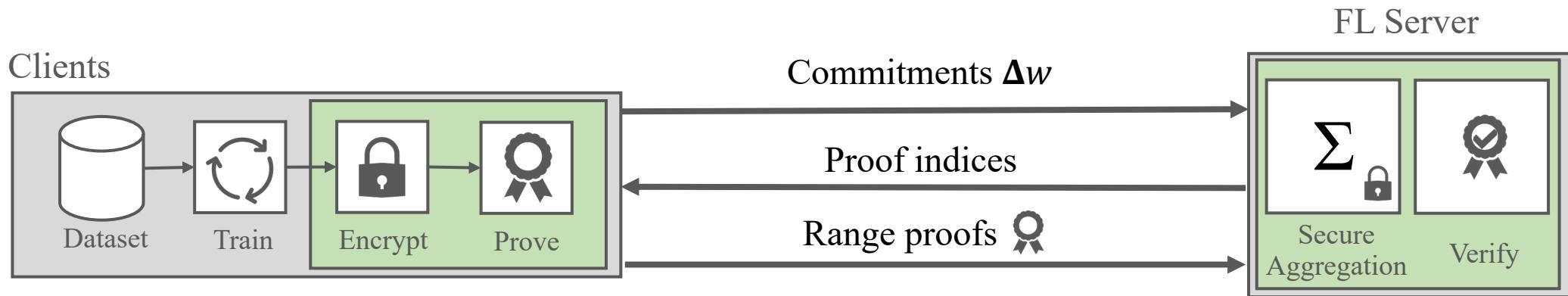


Number of scaled weights required

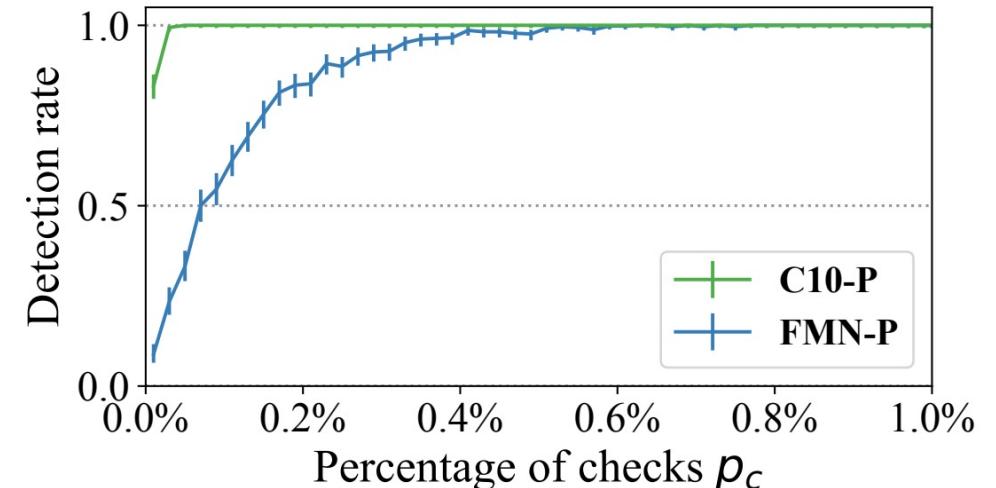


Number of required checks

Optimizing L_∞



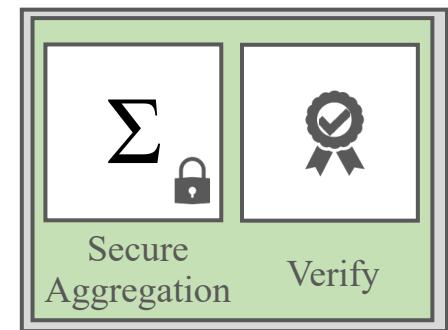
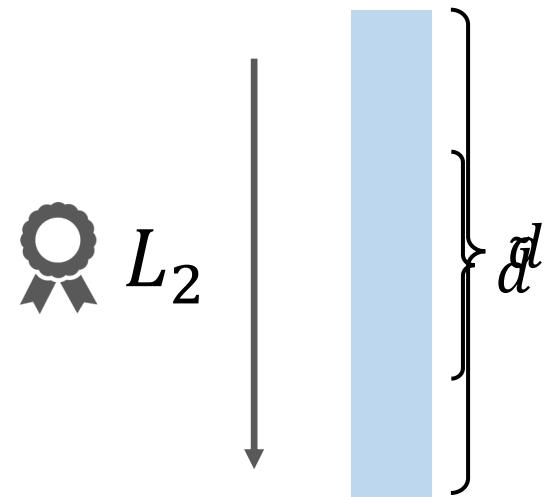
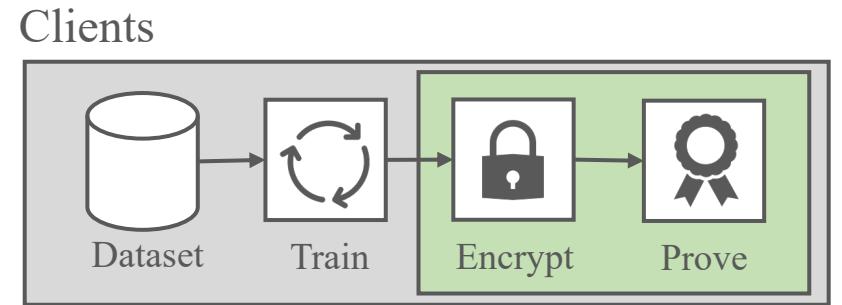
Number of scaled weights required



Number of required checks

Optimizing L_2

$$\left. \begin{array}{c} \text{Model} \\ \text{parameters} \end{array} \right\} d = \left. \begin{array}{c} P \\ \times \end{array} \right\} \left. \begin{array}{c} \text{Update} \\ \text{parameters} \end{array} \right\} \tilde{d}$$



FL Server

65

RoFL: End-To-End Performance

CIFAR-10 Model 270k Parameters

Setup: 48 Clients, 160 rounds

Plaintext

Accuracy: 0.86

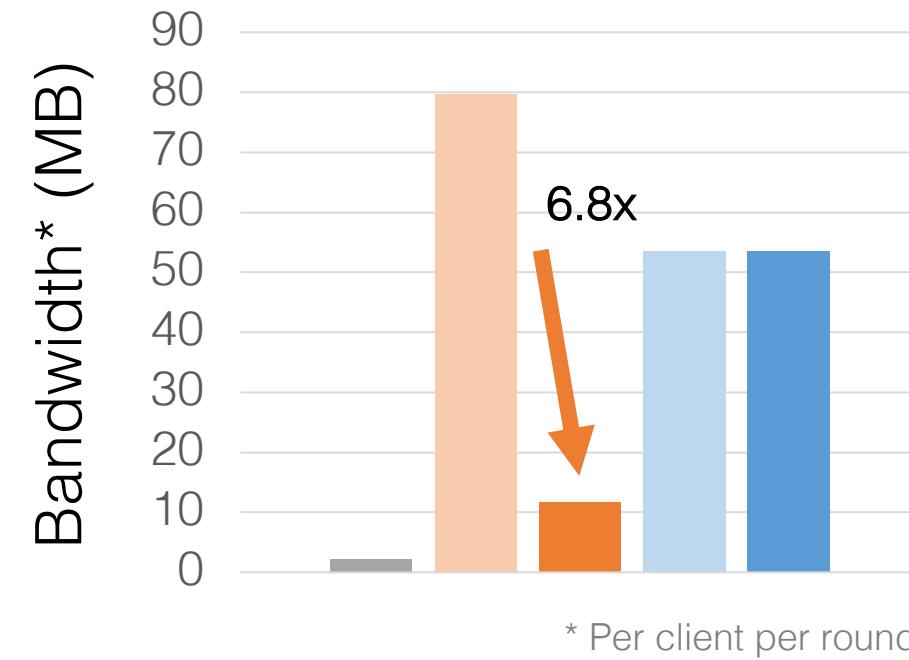
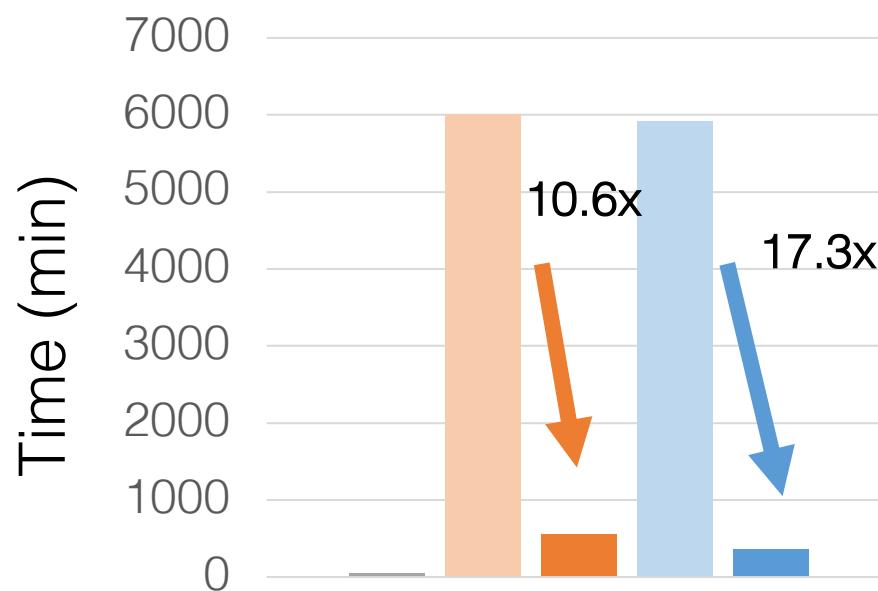
L_2

0.85
0.82

L_2 Optimized

L_∞

L_∞ Optimized



Evaluation: End-To-End

Shakespeare Model 818k Parameters

Setup: 48 Clients, 20 rounds

Plaintext

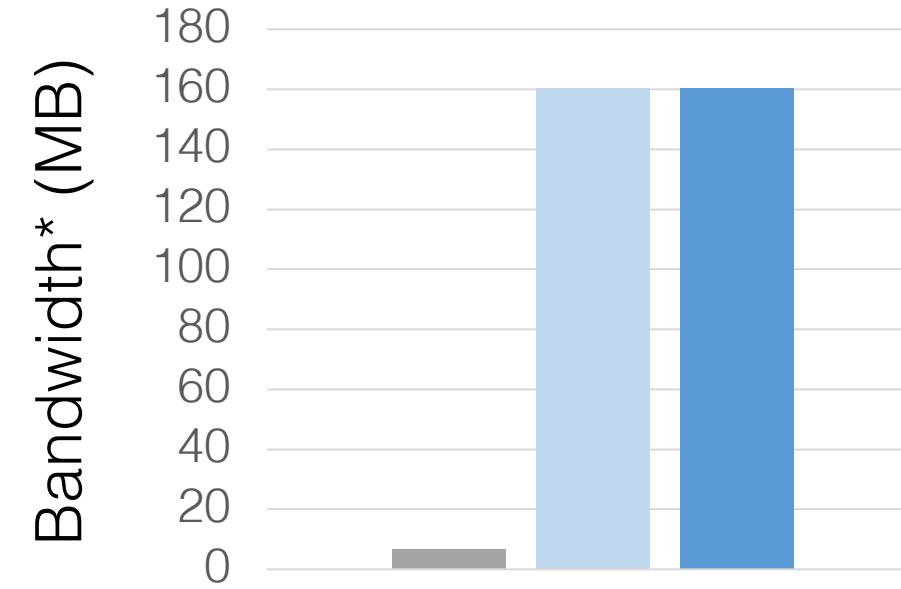
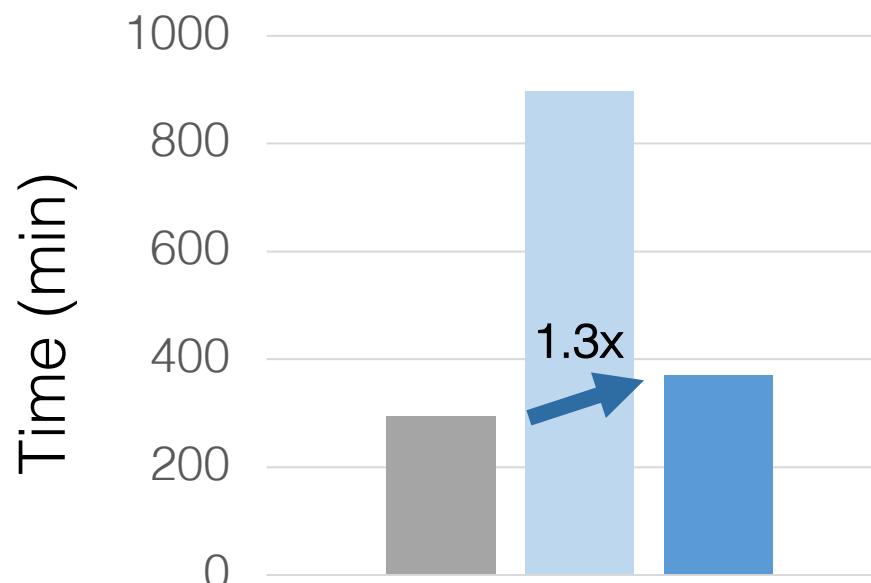
Accuracy: 0.57

L_{∞}

0.57

L_{∞} Optimized

0.57



* Per client per round



<https://arxiv.org/pdf/2107.03311.pdf> (Preprint)



Analysis Code: <https://github.com/pps-lab/fl-analysis>

RoFL Code: <https://github.com/pps-lab/rofl-project-code>