

RoFL: Robustness of Secure Federated Learning



Hidde Lycklama*



Lukas Burkhalter*



Alexander Viand

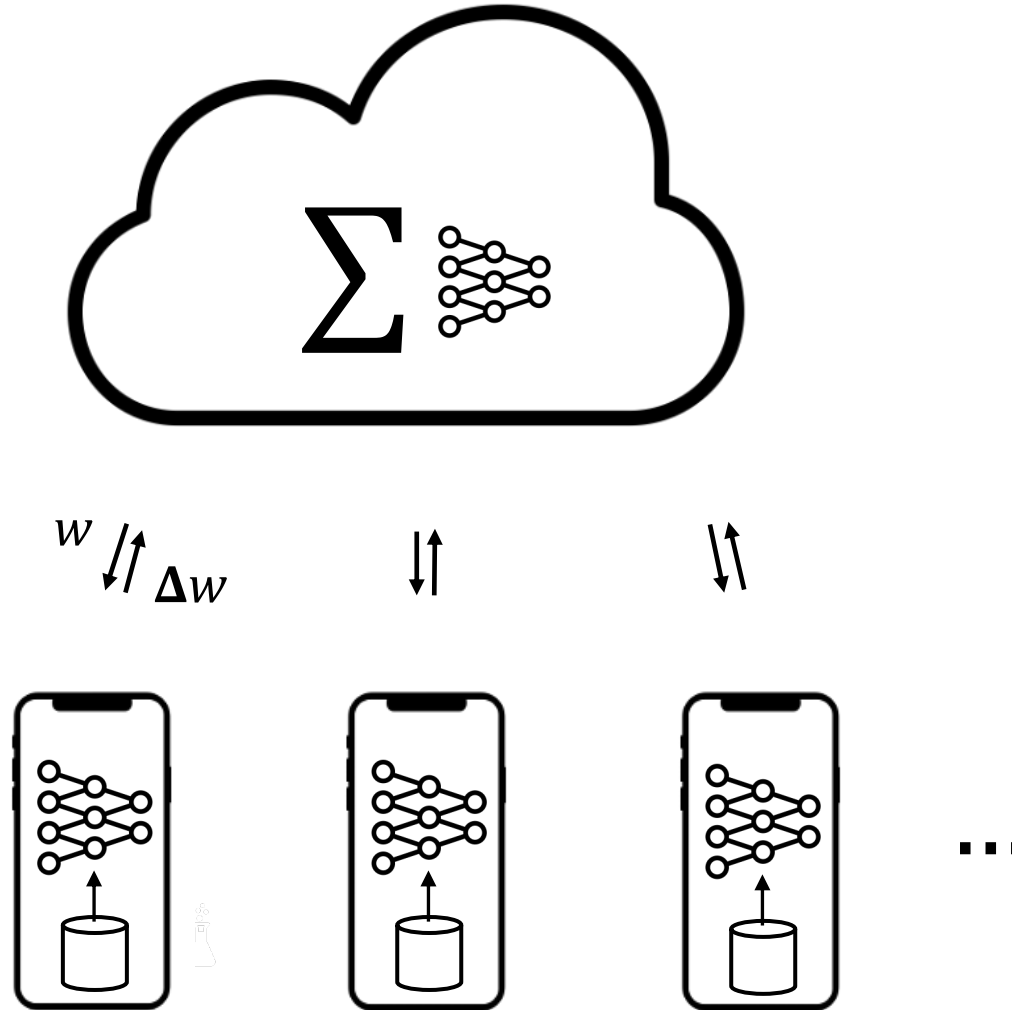


Nicolas Küchler



Anwar Hithnawi

Federated Learning



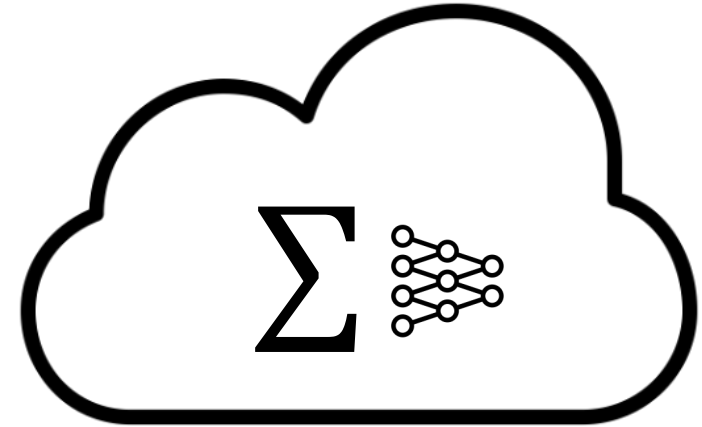
Federated Learning



Purpose Limitation



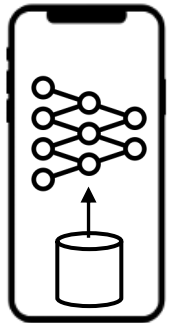
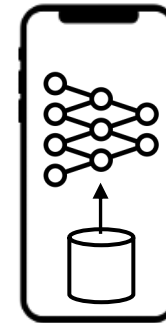
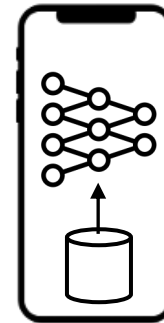
Data Minimization



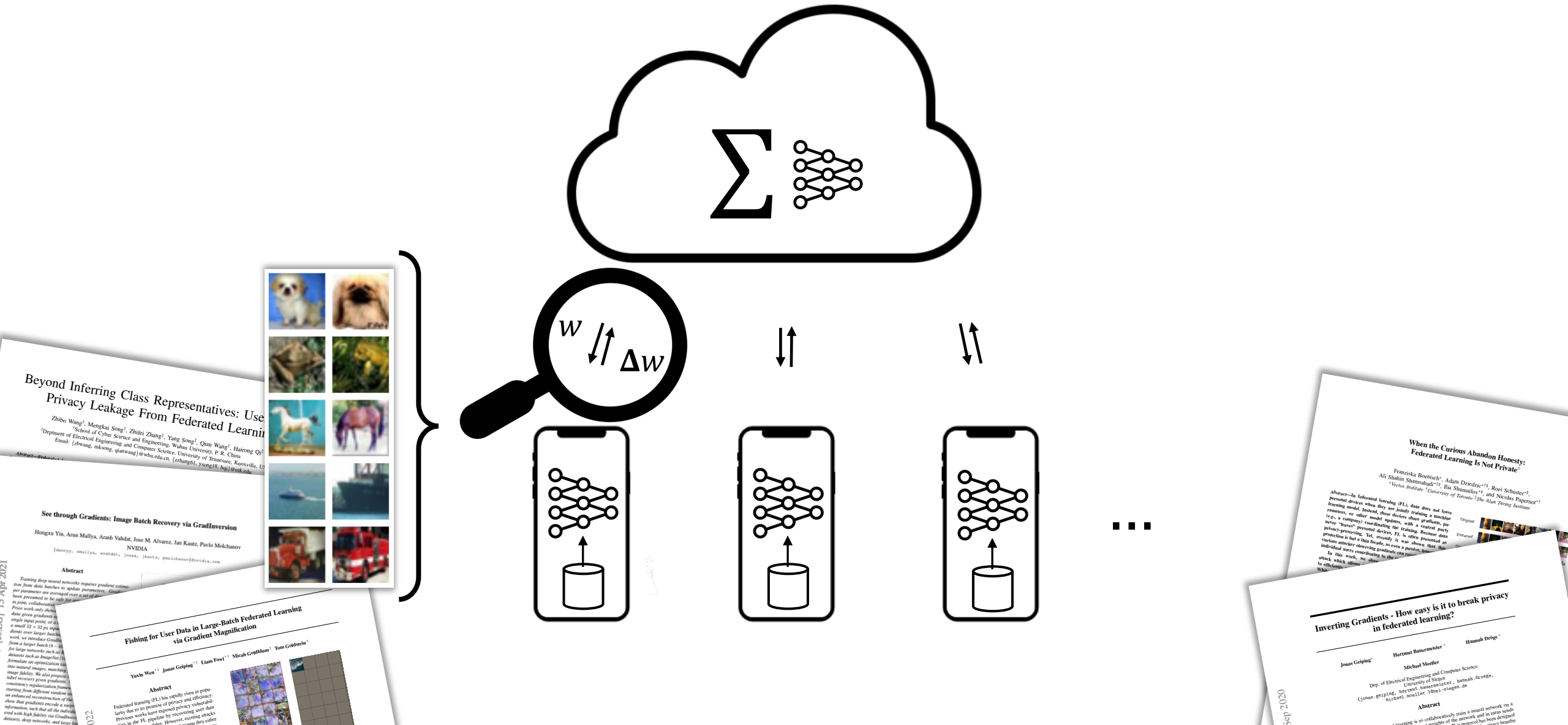
$w \rightleftarrows \Delta w$

\updownarrow

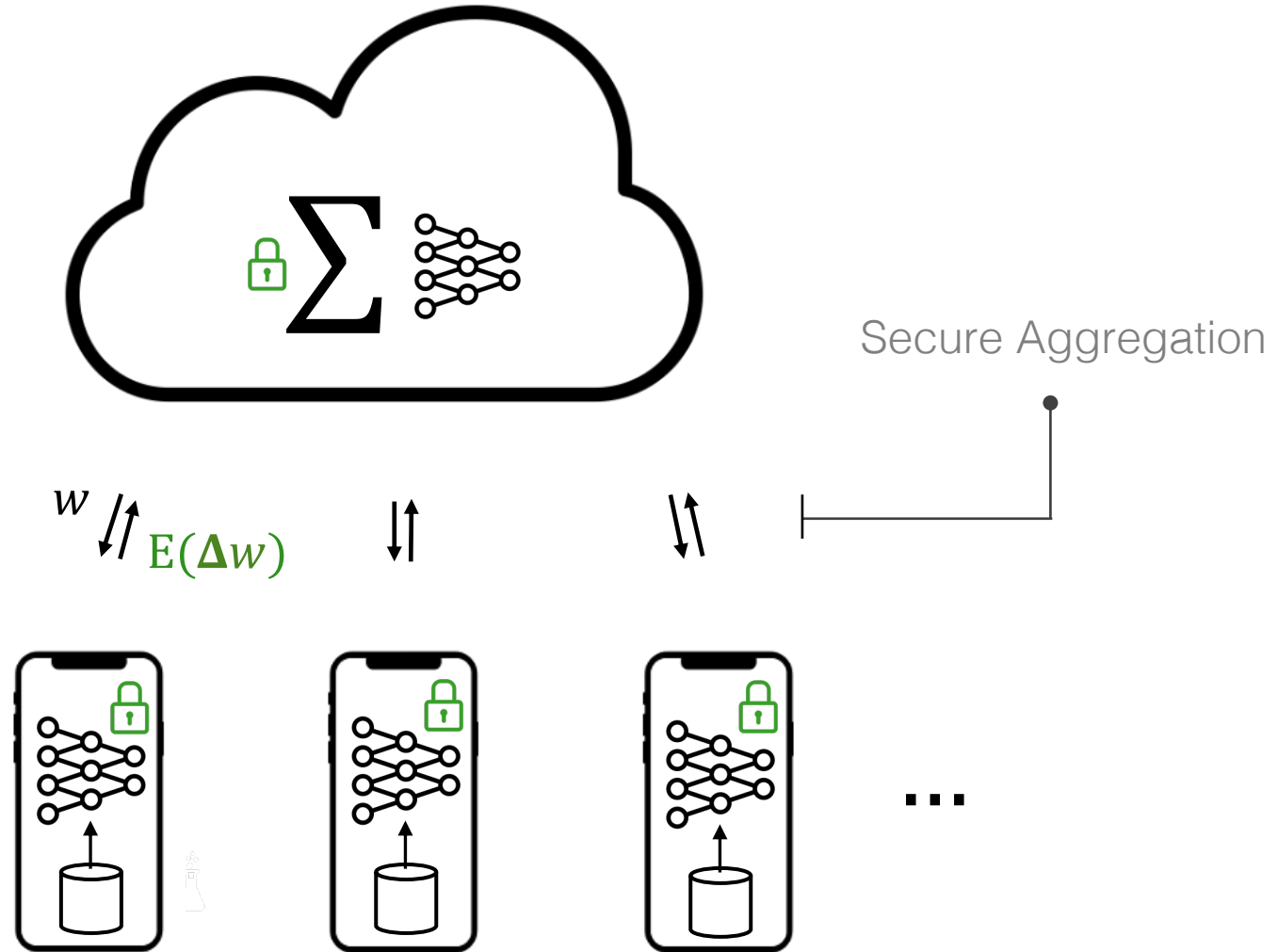
\updownarrow



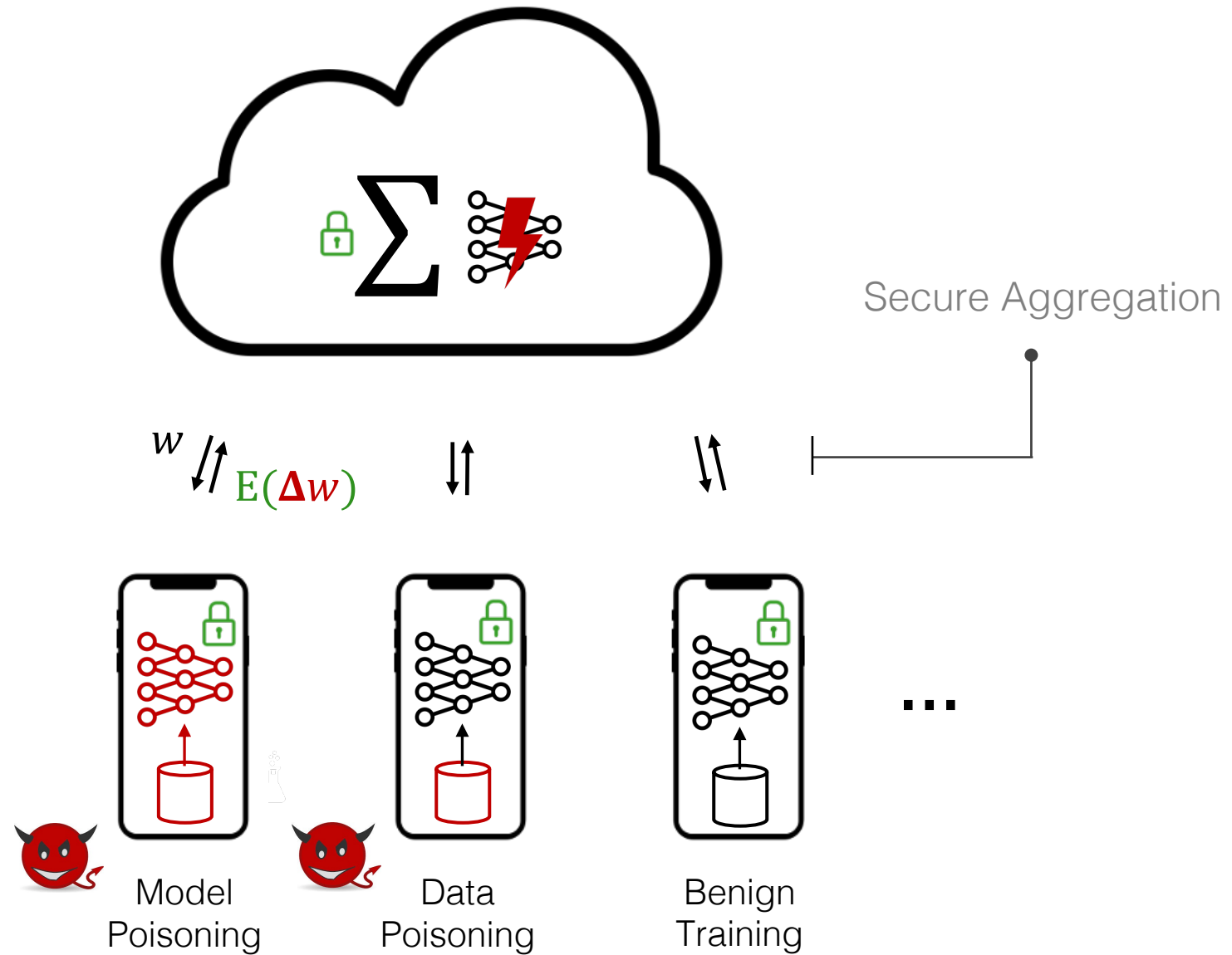
Federated Learning: Input Privacy



Secure Federated Learning



Malicious Clients



RoFL: Robustness of Secure Federated Learning

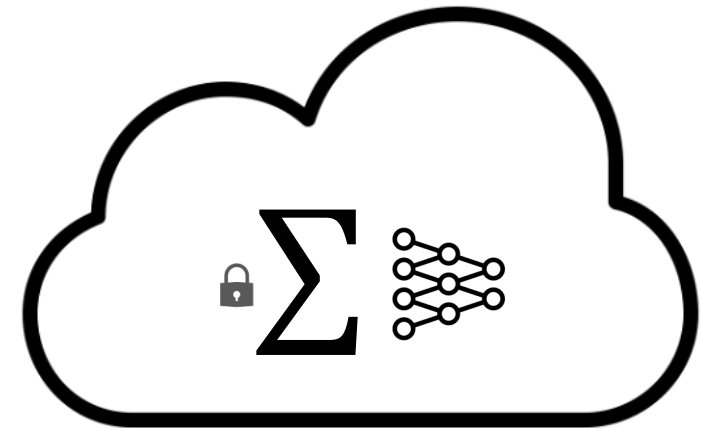
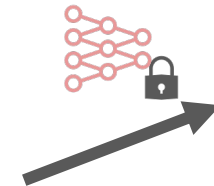
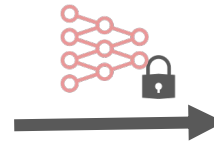
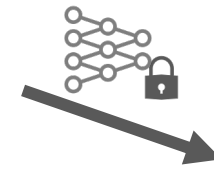
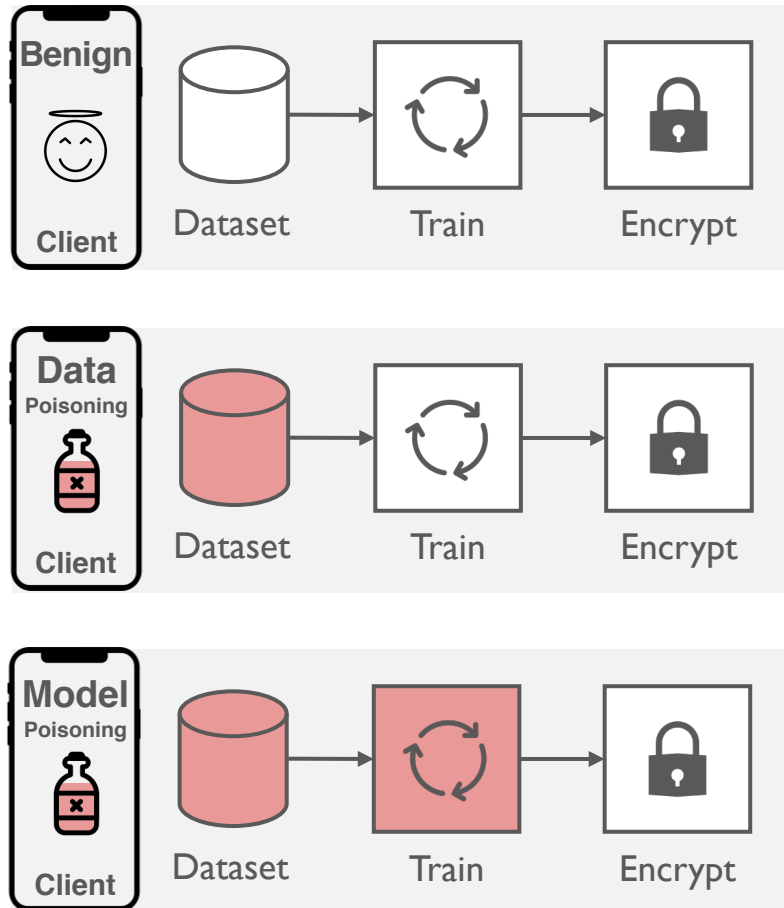
Understand
Vulnerabilities in FL



Cryptographically
Enforce Constraints



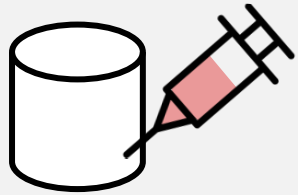
Adversarial Clients



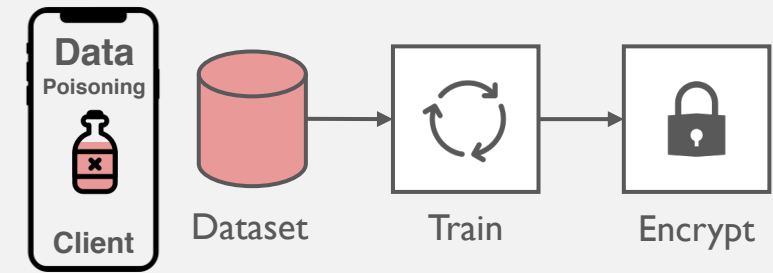
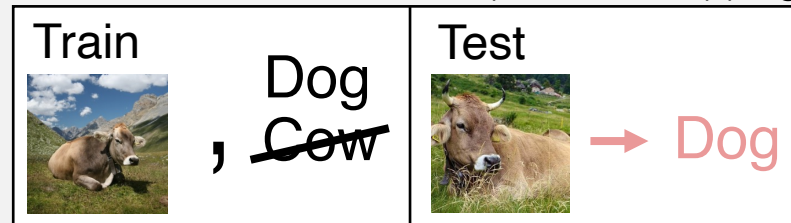
Adversarial Clients

Data Poisoning

adversary controls training data

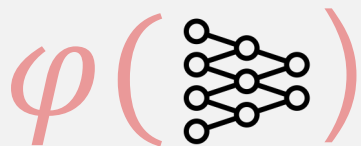


Example: Label Flipping

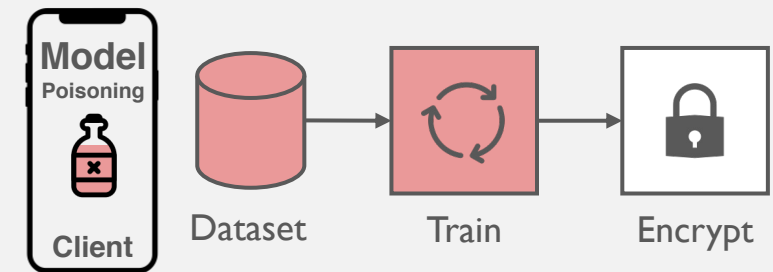
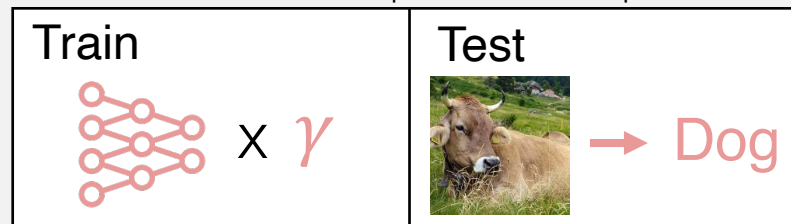


Model Poisoning

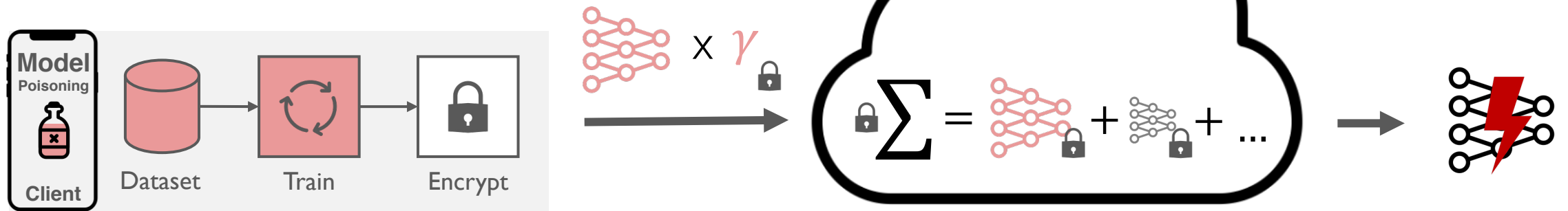
adversary controls model updates



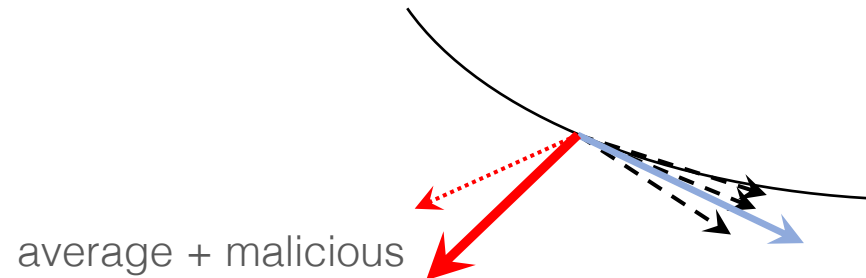
Example: Model Replacement



Adversarial Clients



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]

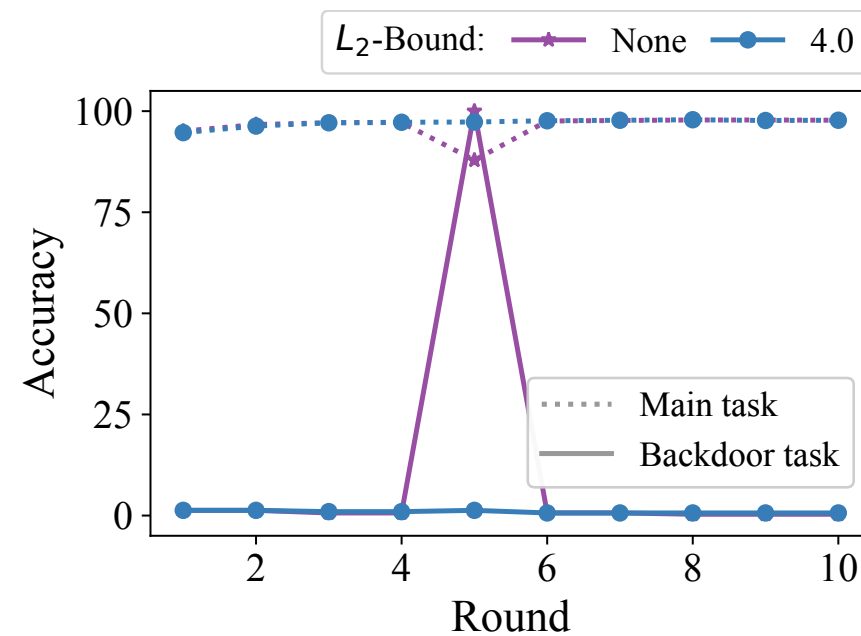
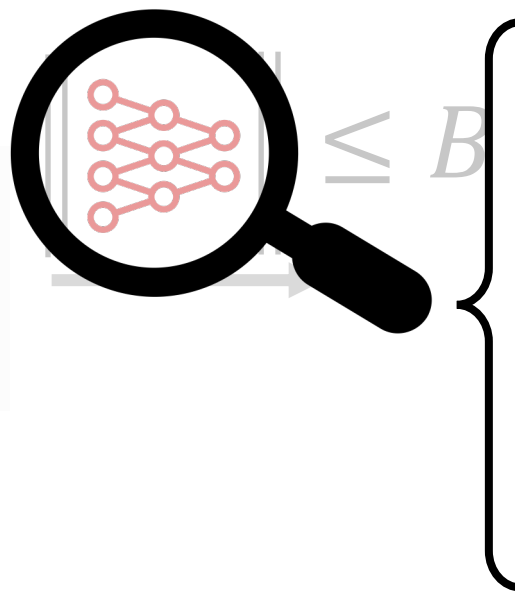
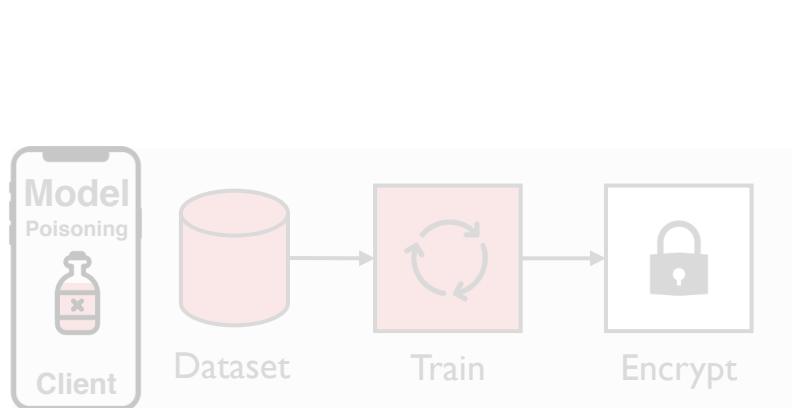
Security: Private Data-Collection Systems

- 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

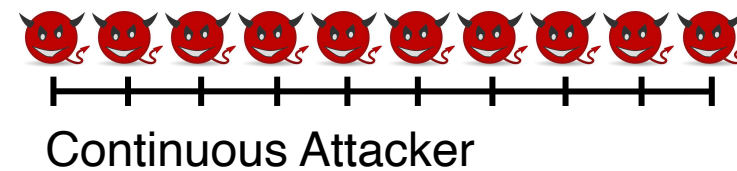
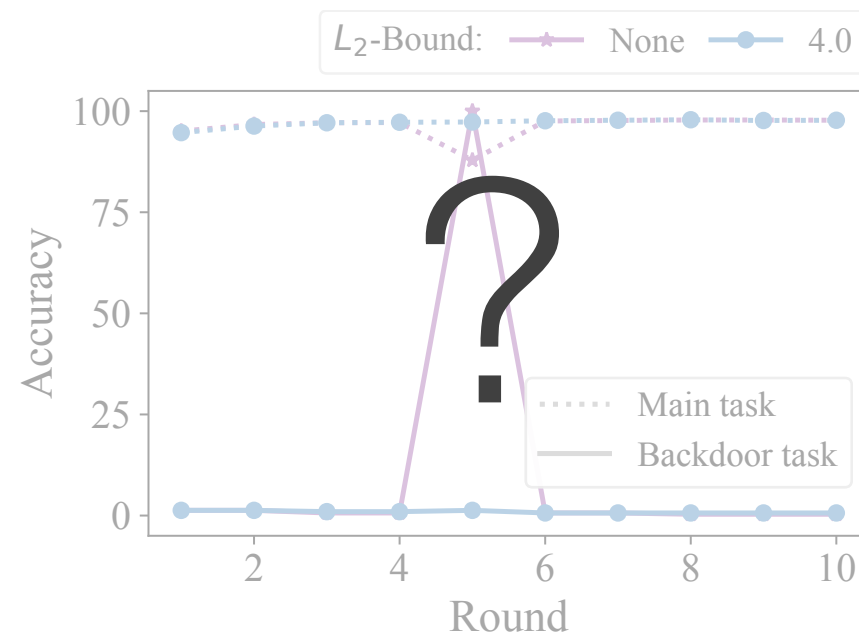
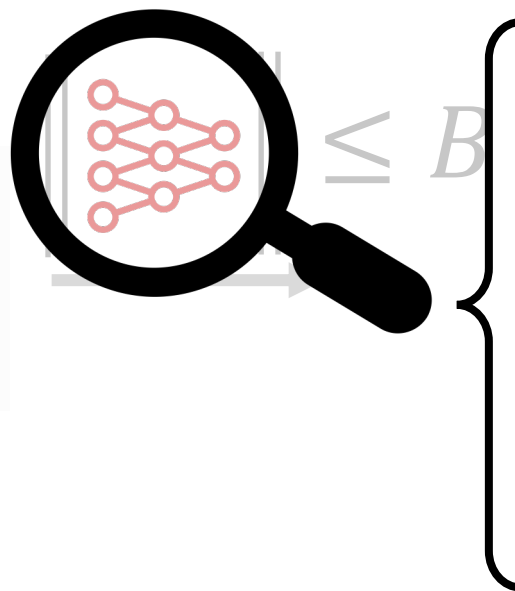
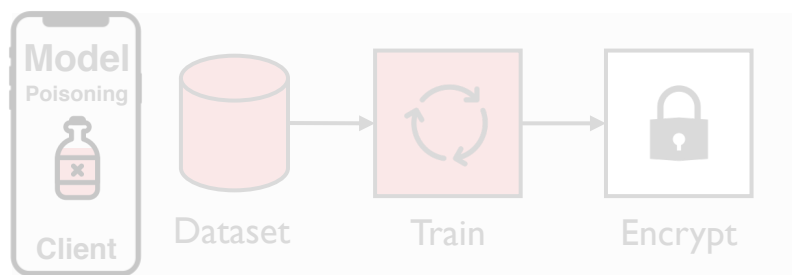
A Well-Formed Client Submission in Federated Learning

Norm bound



Single-shot Attacker (round 5)

Norm bound



Is the norm bound actually effective?

How To Backdoor Federated Learning

Can You Really Backdoor Federated Learning?

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

Long Tail ...

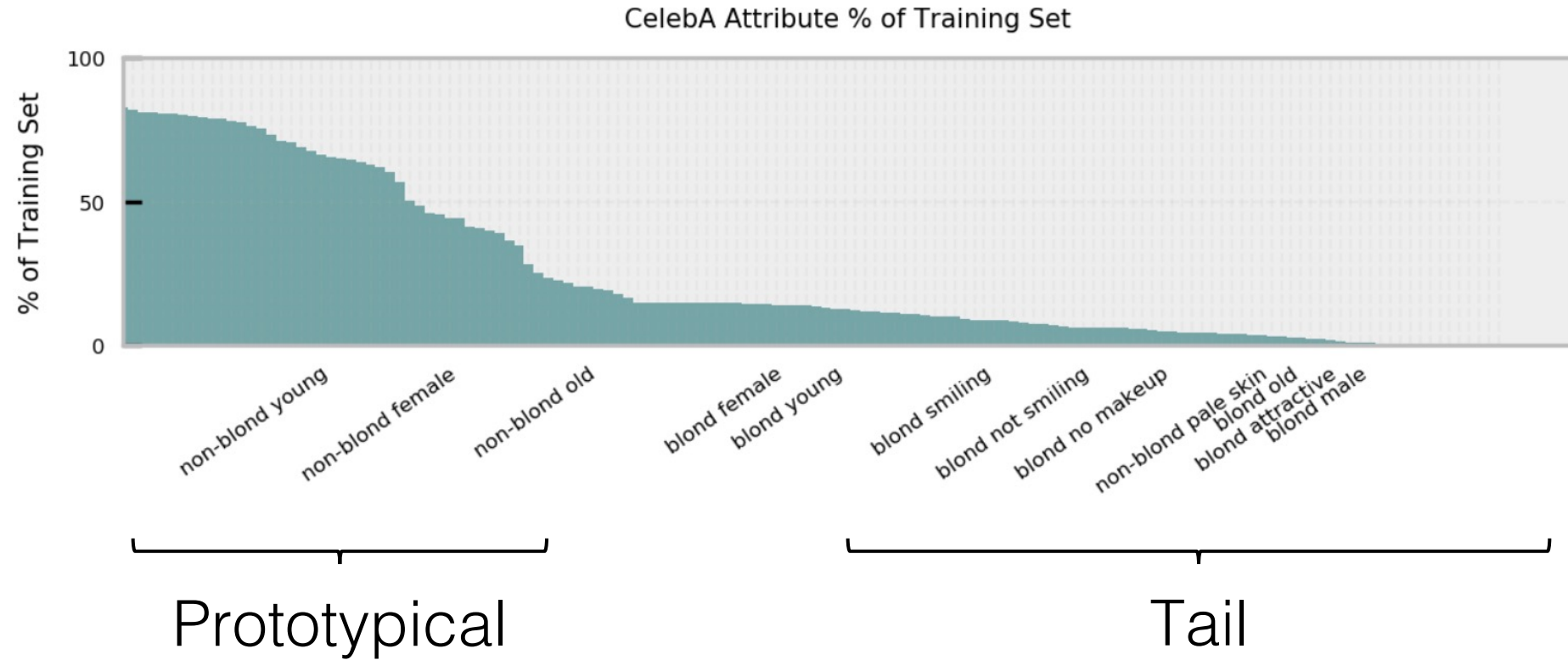
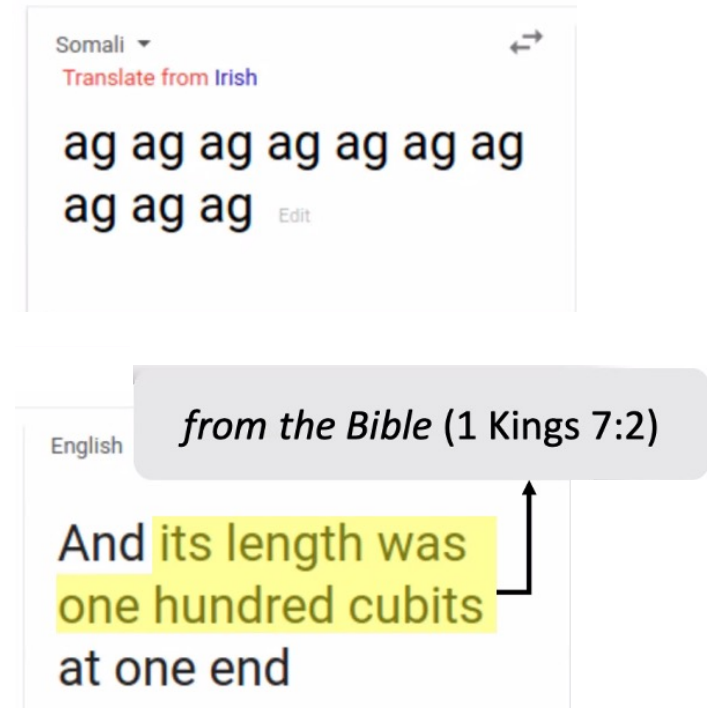
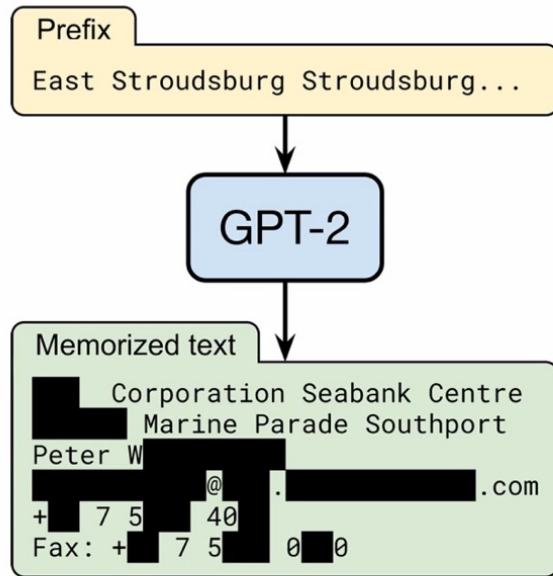


Fig: Hooker, Moorosi et al., 2020.

Model Capacity Implications on Privacy ...



Memorization leads to leakage of private text

Analysis: Understanding FL Robustness



Adaptive attacks

MP-PD: Projected Gradient Descent [Sun et al., FLDPC@NeurIPS'19]

MP-NT: Neurotoxin [Zhang et al., ICML'22]

MP-AT: Anticipate [Wen et al., AdvML@ICML'22]

Considered:

Attack
Objective

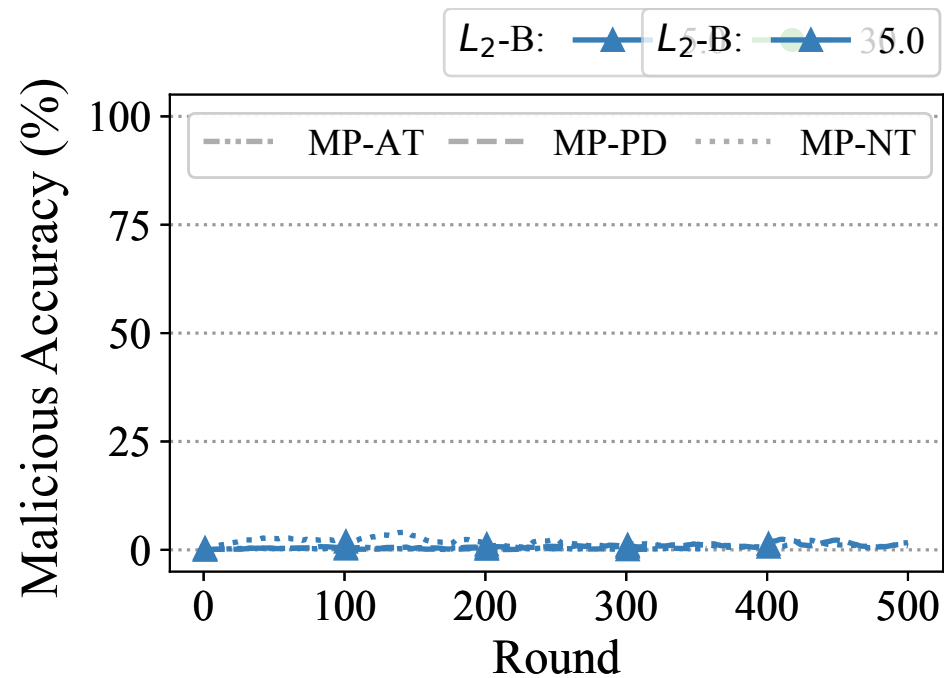
Number of
Attackers

Bound
Selection

Pixel-Pattern
Backdoors

Untargeted
Attacks

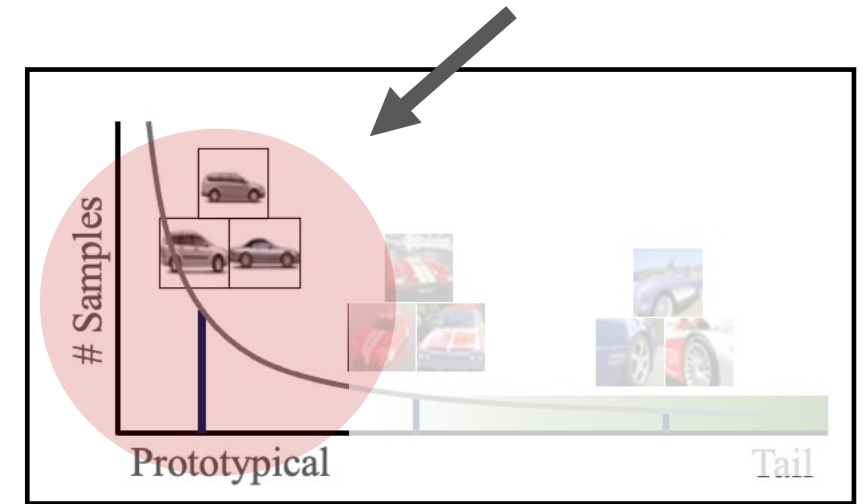
Impact of Attack Objective on Backdoor Attacks



CIFAR-10

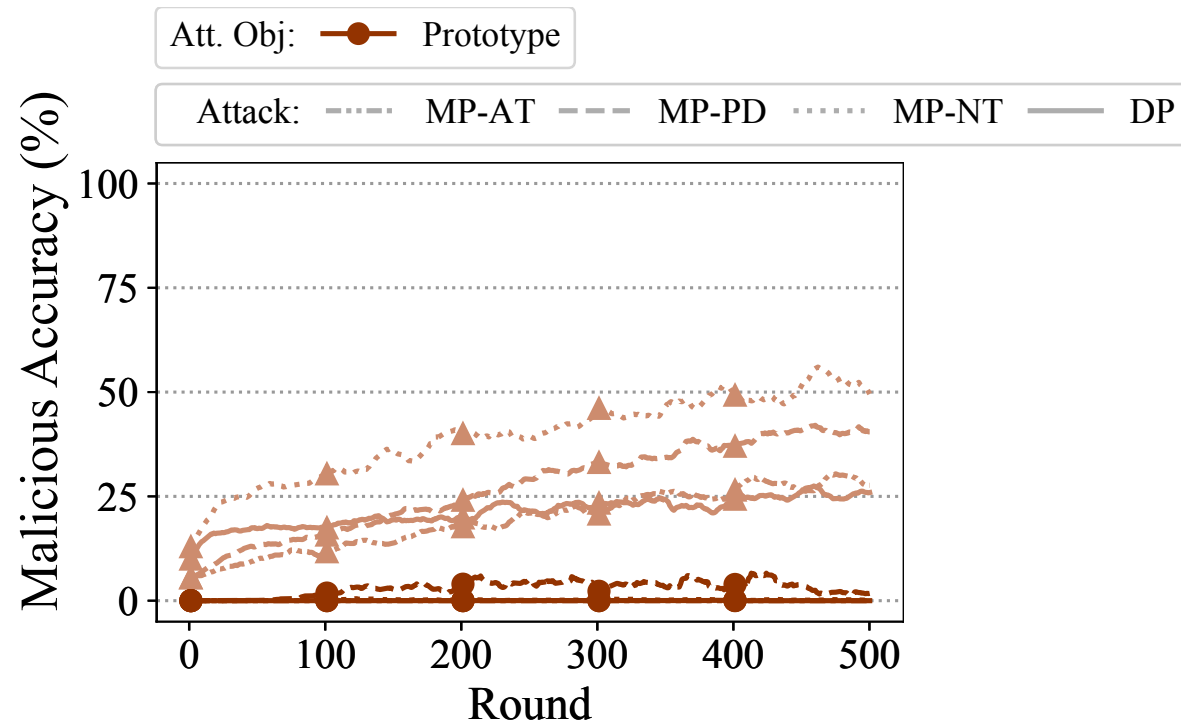


Continuous Attacker



Prototypical Targets

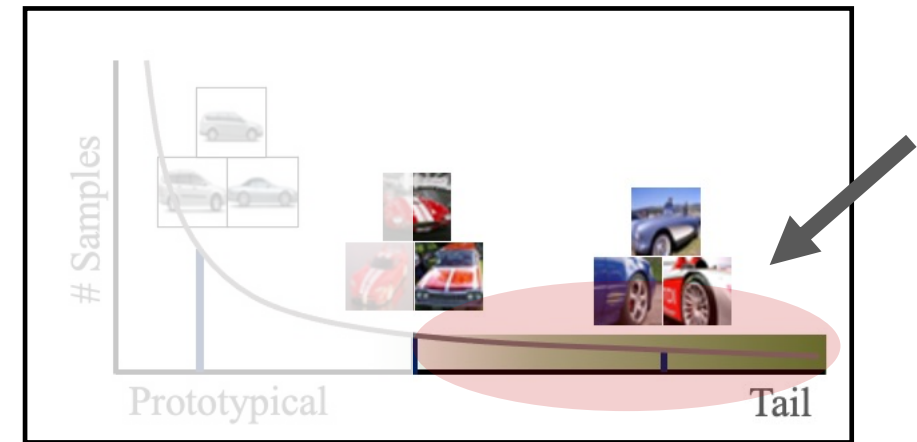
Impact of Attack Objective on Backdoor Attacks



CIFAR-10

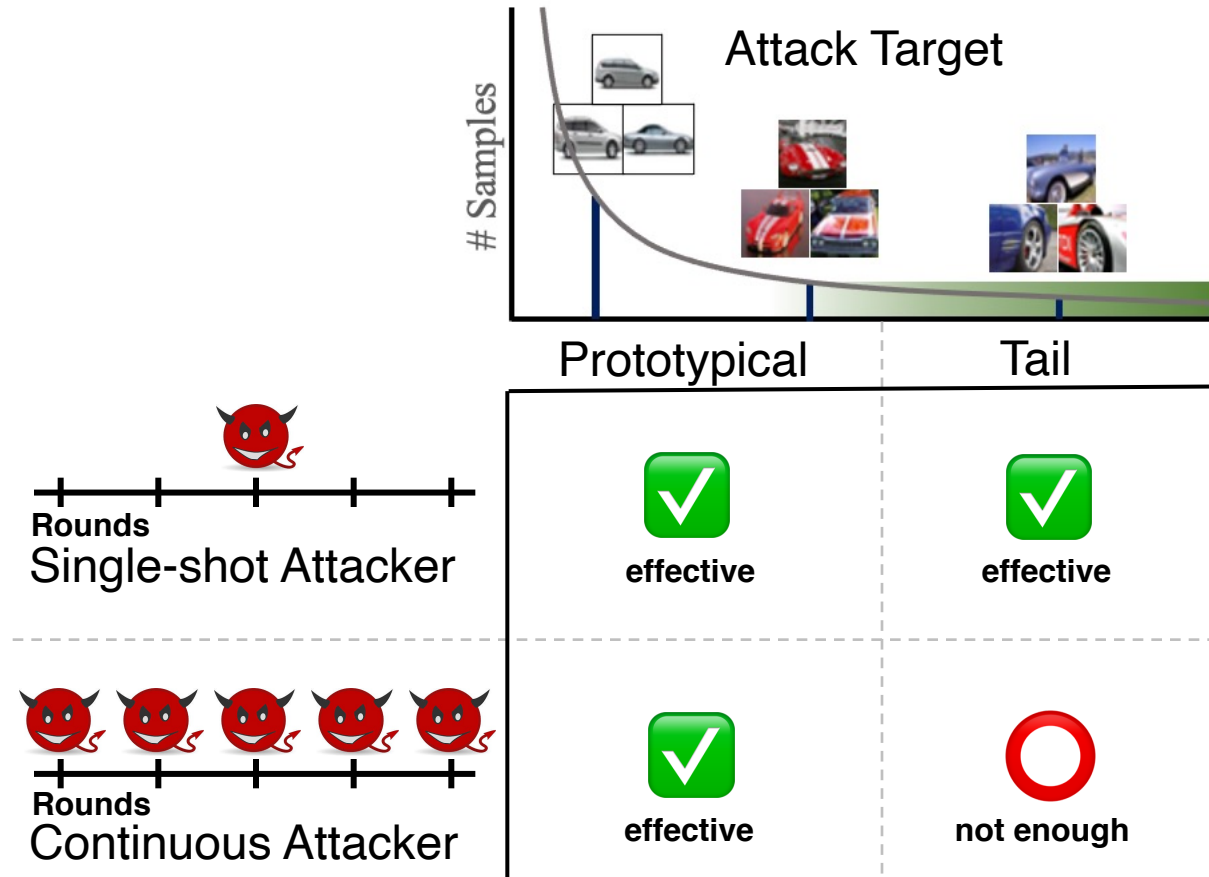


Continuous Attacker



Tail Targets

Norm Bound Provides Practical Robustness Guarantees



... requires a strong attacker that is consistently selected and targets a tail sample

RoFL: Robustness of Secure Federated Learning

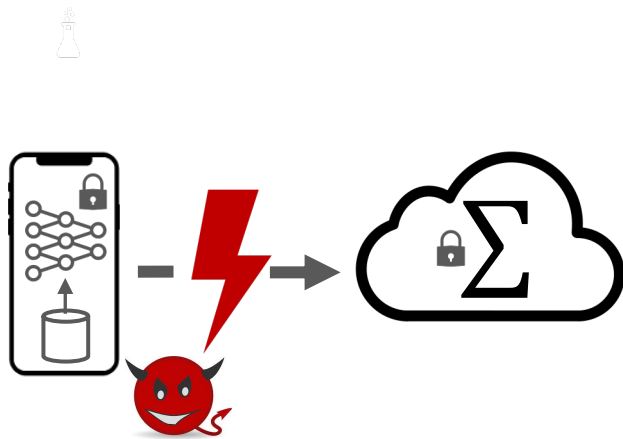
Understand
Vulnerabilities in FL



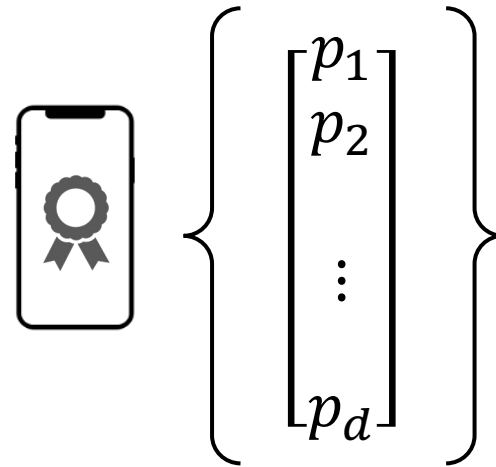
Cryptographically
Enforce Constraints



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



Correctness



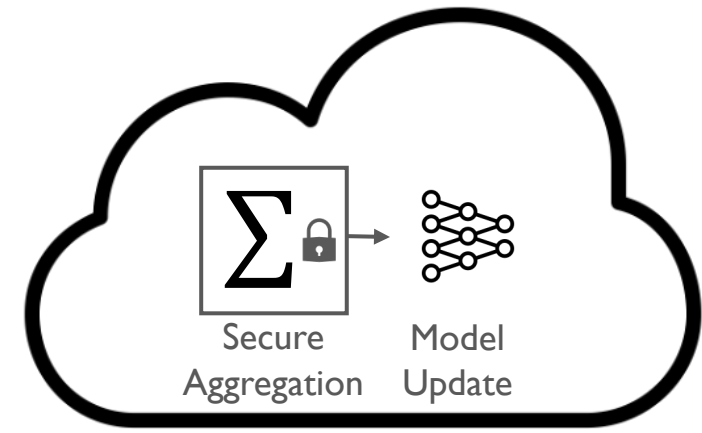
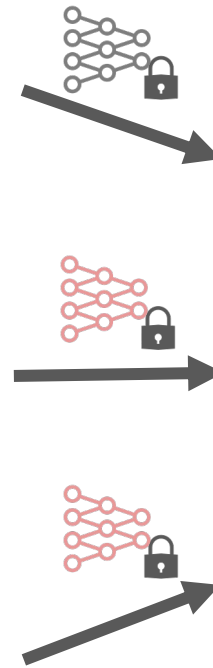
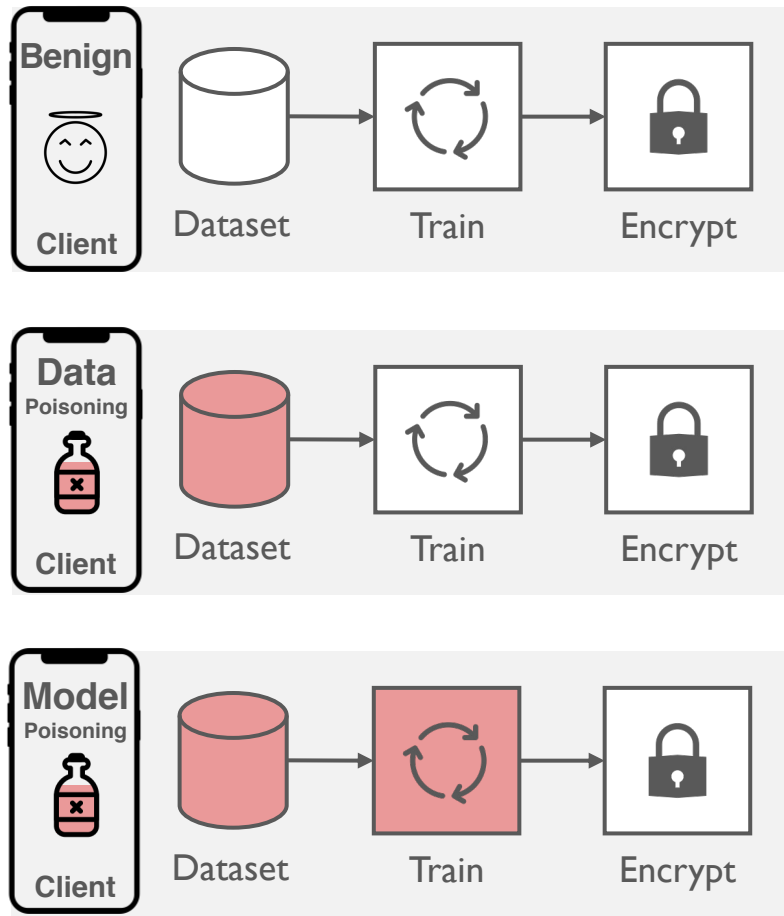
Private Input Validation



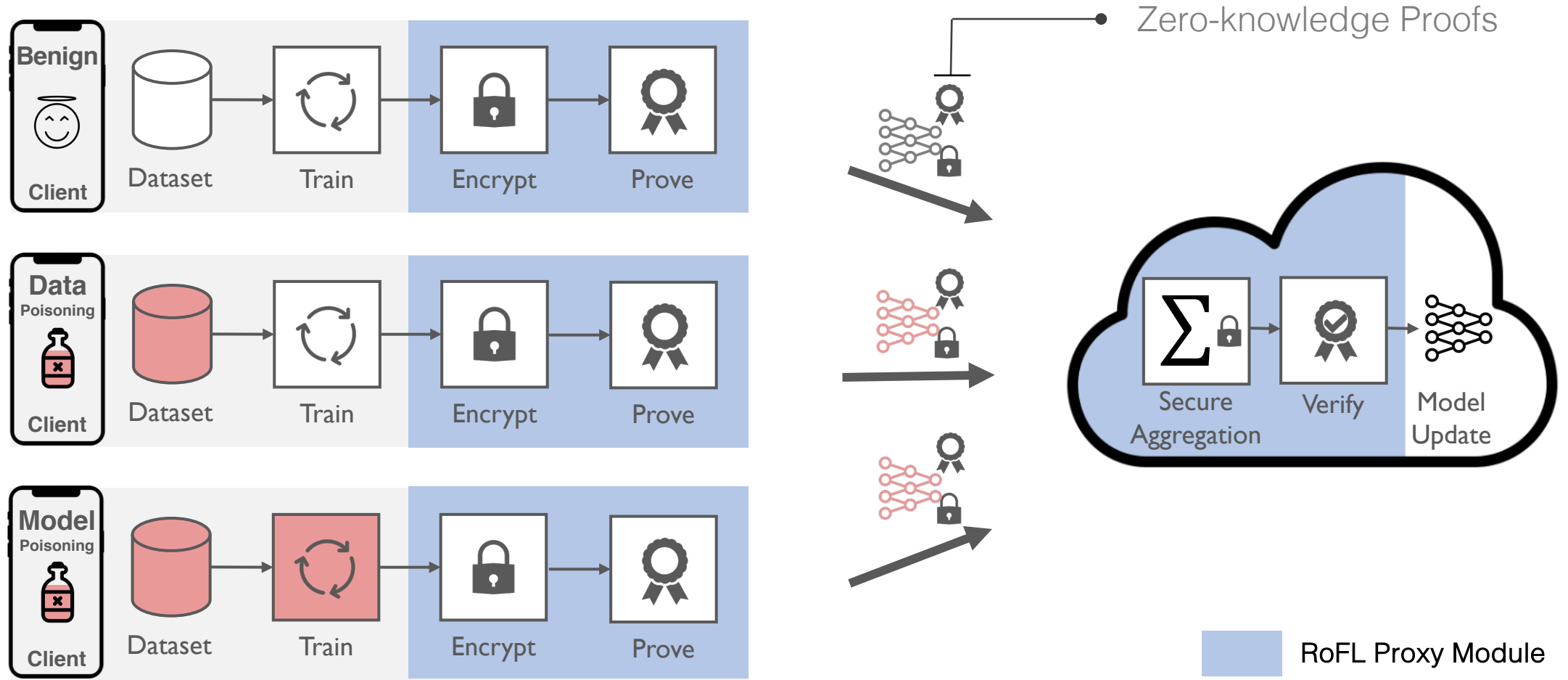
- Compressed Sigma protocols
- Optimistic continuation
- Probabilistic checking
- Subspace learning

Optimizations

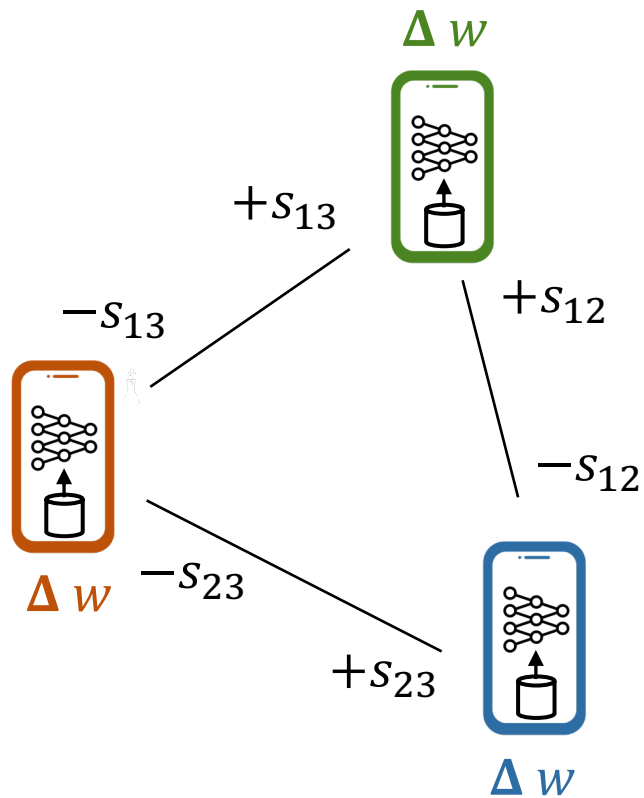
Secure Federated Learning



RoFL Augments Secure Federated Learning



Secure Aggregation



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

Idea: Additive masks based on pairwise secrets s_{ij}

$$r_1 + r_2 + r_3 = 0$$

where

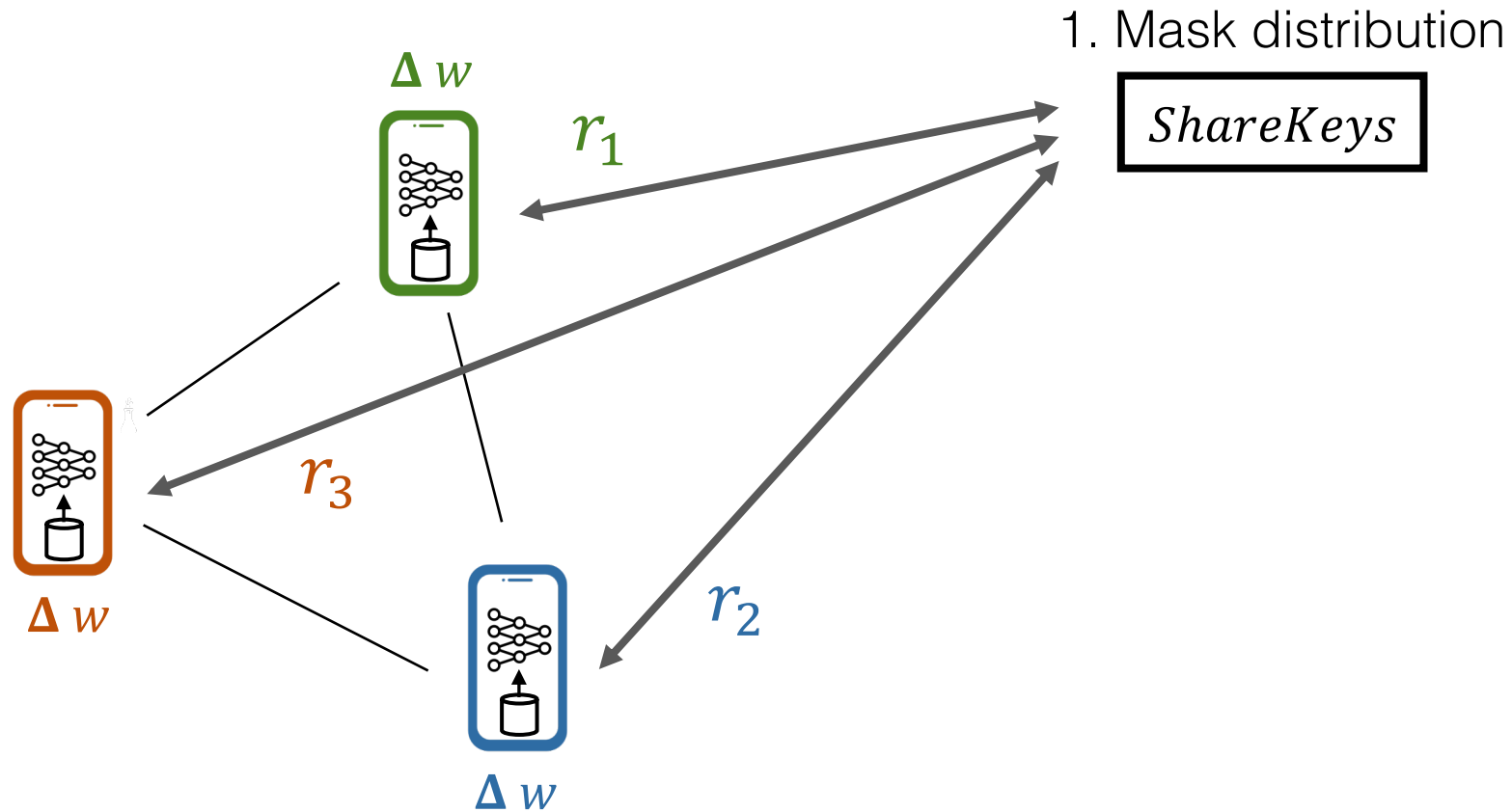
$$r_1 = s_{12} + s_{13}$$

$$r_2 = -s_{12} + s_{23}$$

$$r_3 = -s_{13} - s_{23}$$

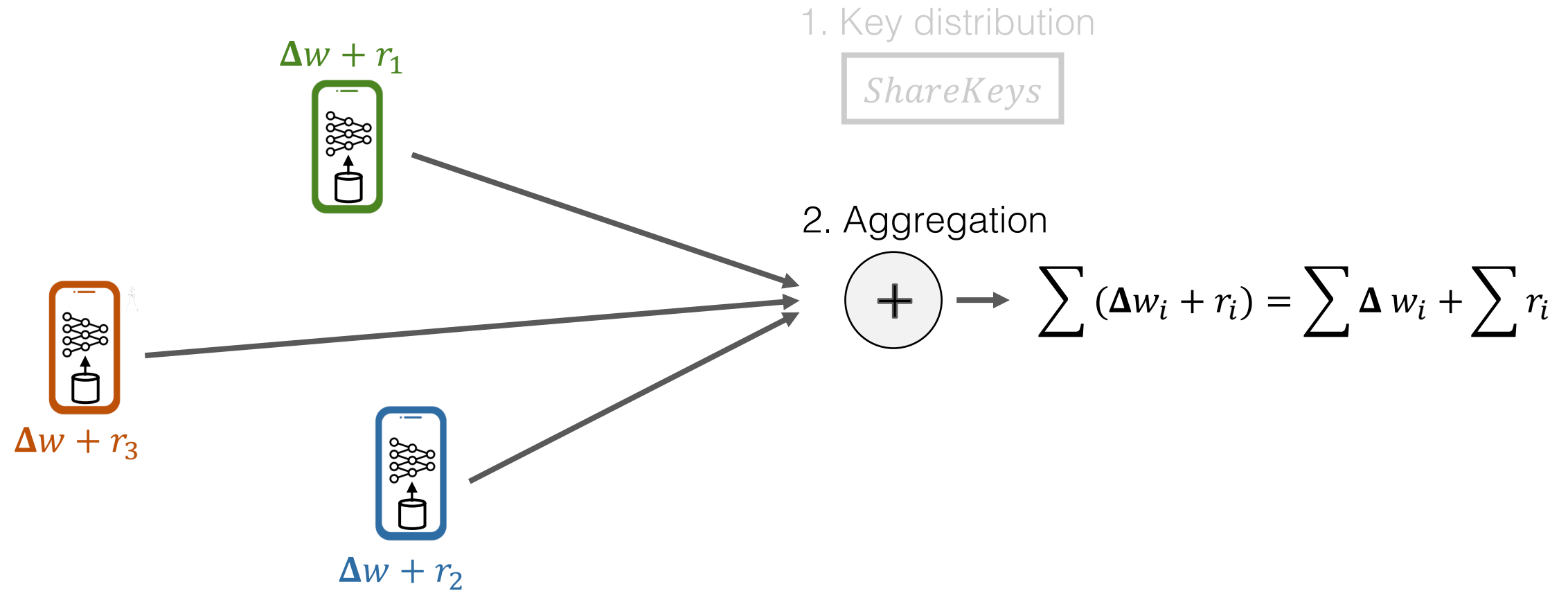
+: modular addition

Secure Aggregation



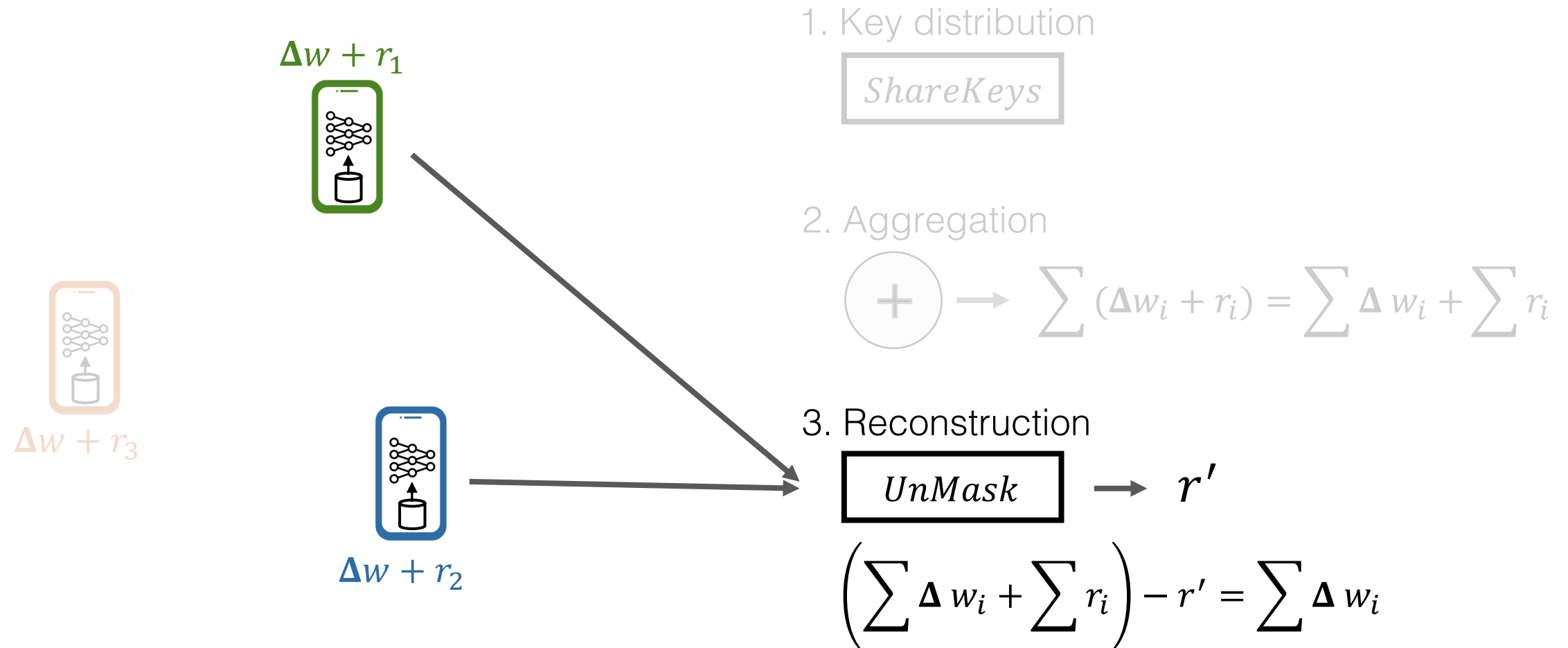
+: modular addition 27

Secure Aggregation



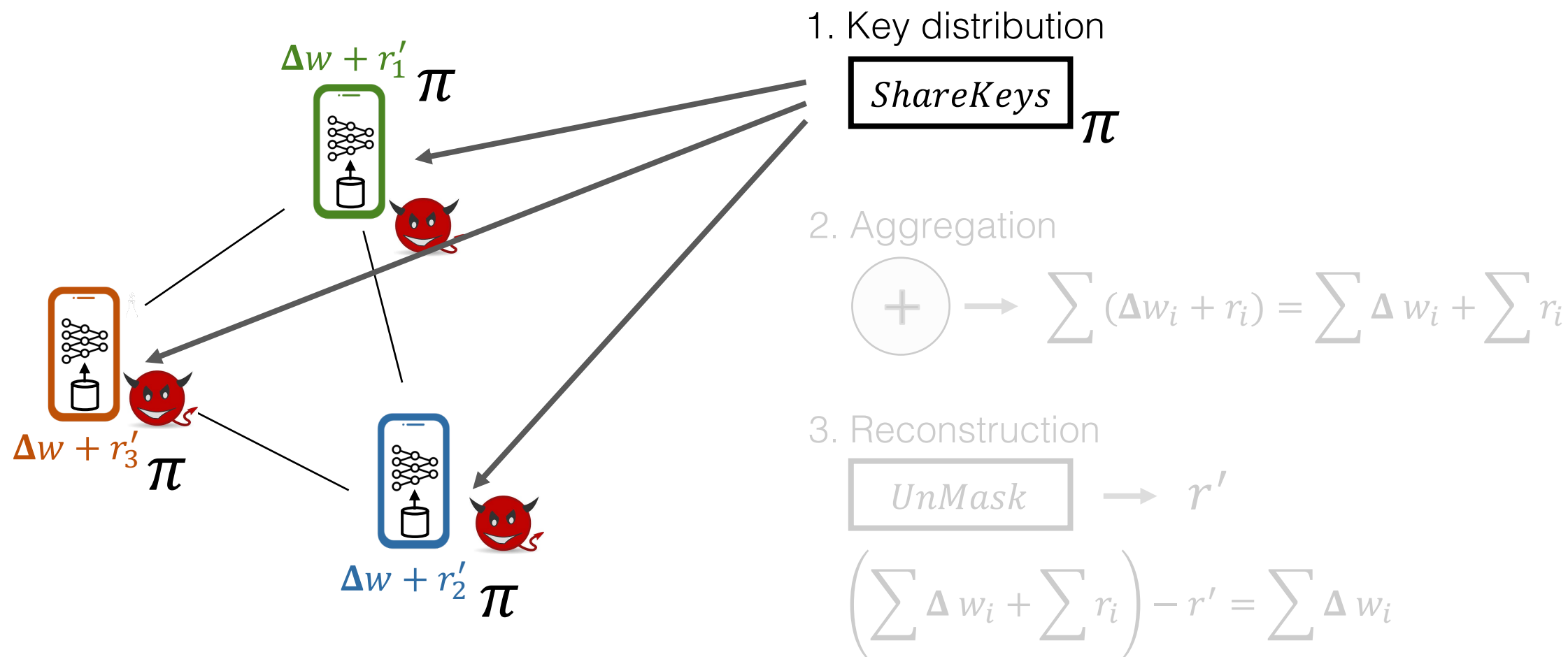
+: modular addition 28

Secure Aggregation

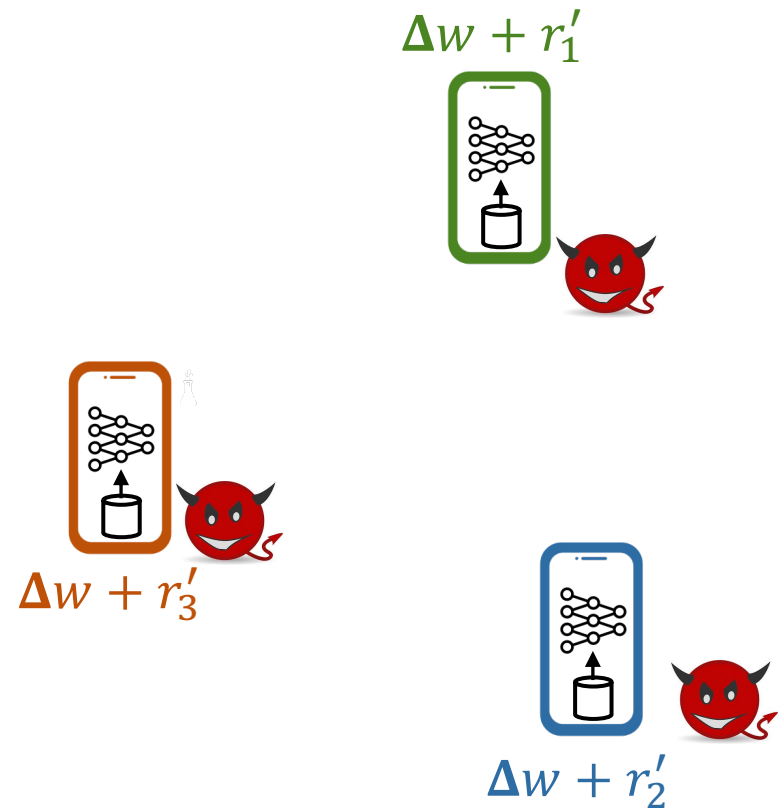


\oplus : modular addition 29

Limitation: Correctness with malicious clients



Insight: Checking $\sum r_i = r'$ sufficient for correctness



1. Key distribution

ShareKeys

2. Aggregation

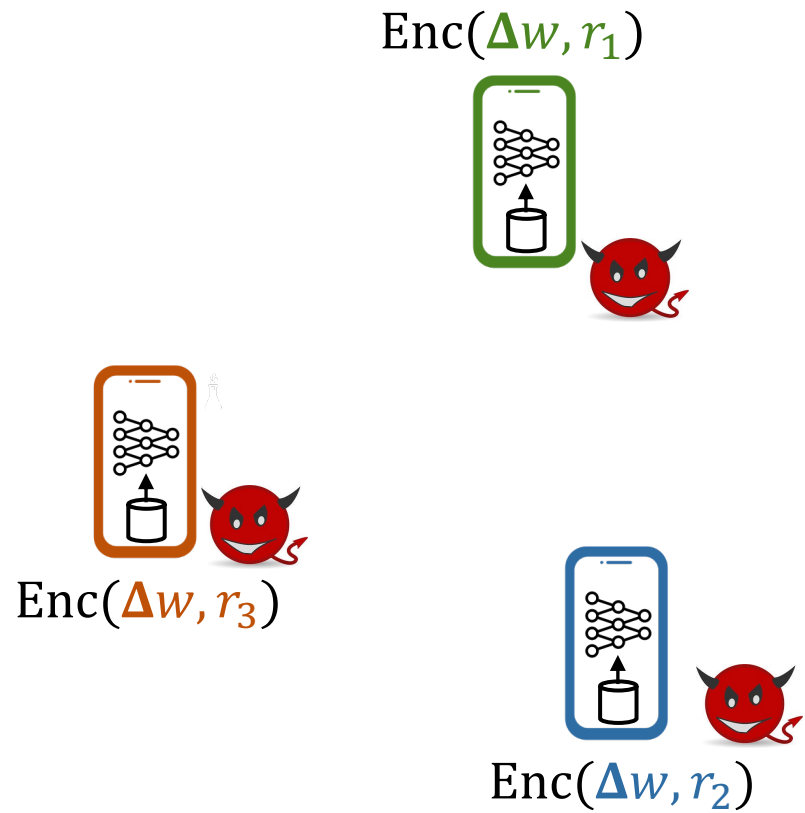
$$\bigoplus \rightarrow \sum (\Delta w_i + r_i) = \sum \Delta w_i + \sum r_i$$

3. Reconstruction

UnMask $\rightarrow r'$

$$\left(\sum \Delta w_i + \sum r_i \right) - r' = \sum \Delta w_i$$

Insight: Checking $\sum r_i = r'$ sufficient for correctness



1. Key distribution

ShareKeys

2. Aggregation

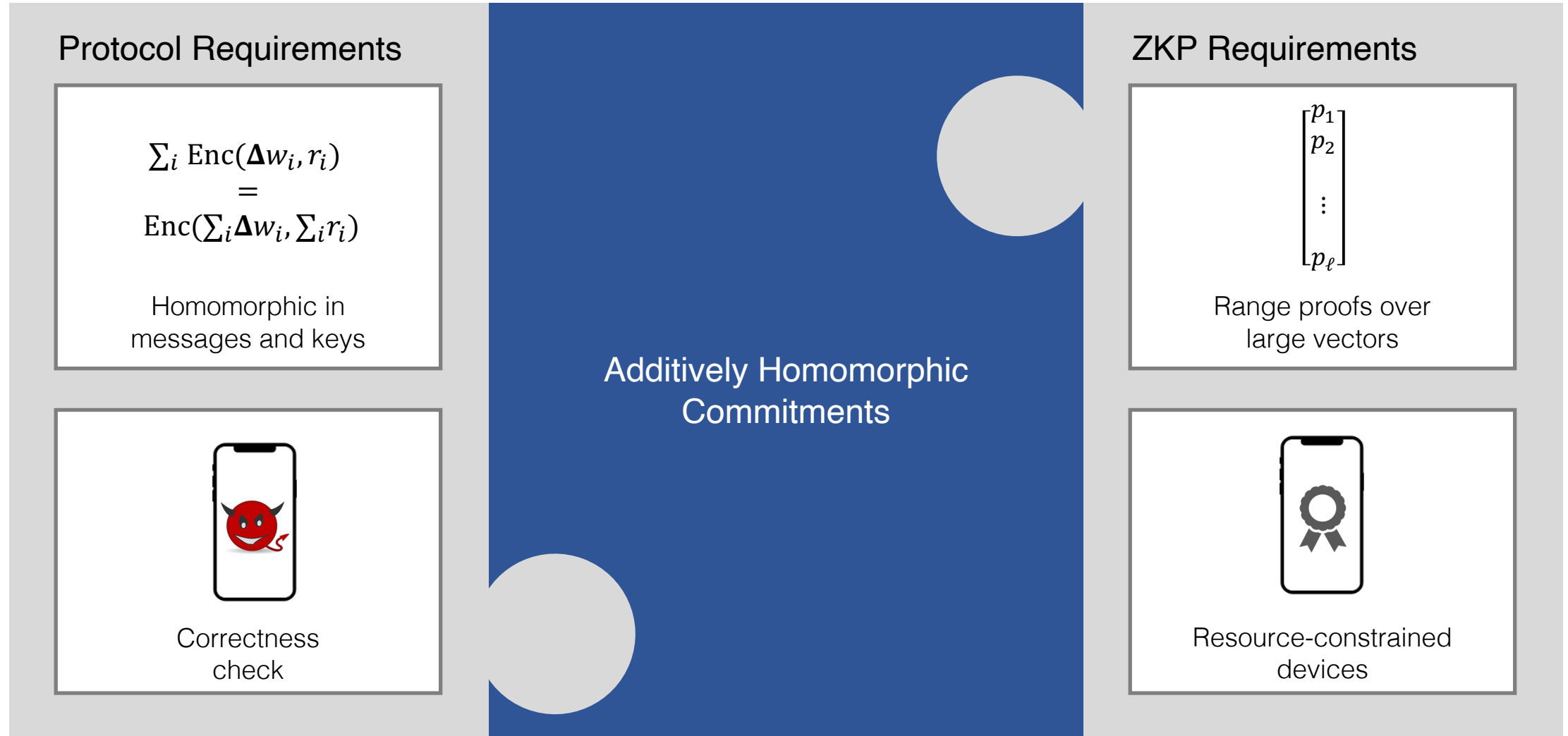
$$\bigoplus \rightarrow \sum \text{Enc}(\Delta w_i, r_i) = \text{Enc}(\sum \Delta w_i, \sum r_i)$$

3. Reconstruction

UnMask $\rightarrow r'$

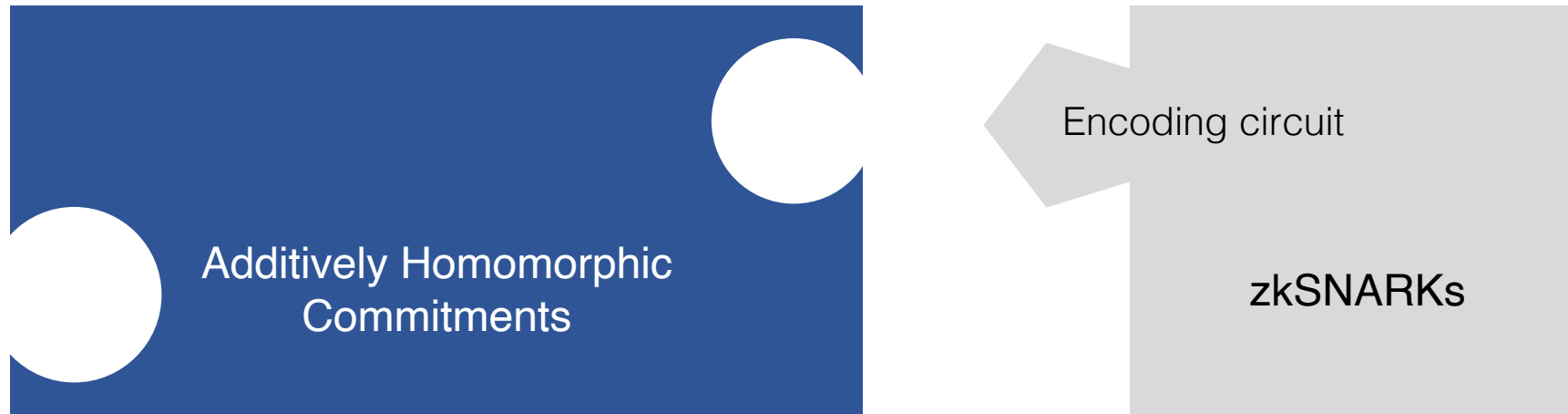
$$\text{Dec}(\text{Enc}(\sum \Delta w_i, \sum r_i), r') = \sum \Delta w_i$$

Efficiency hinges on compatibility with zero-knowledge proofs



Compatibility with Commitments

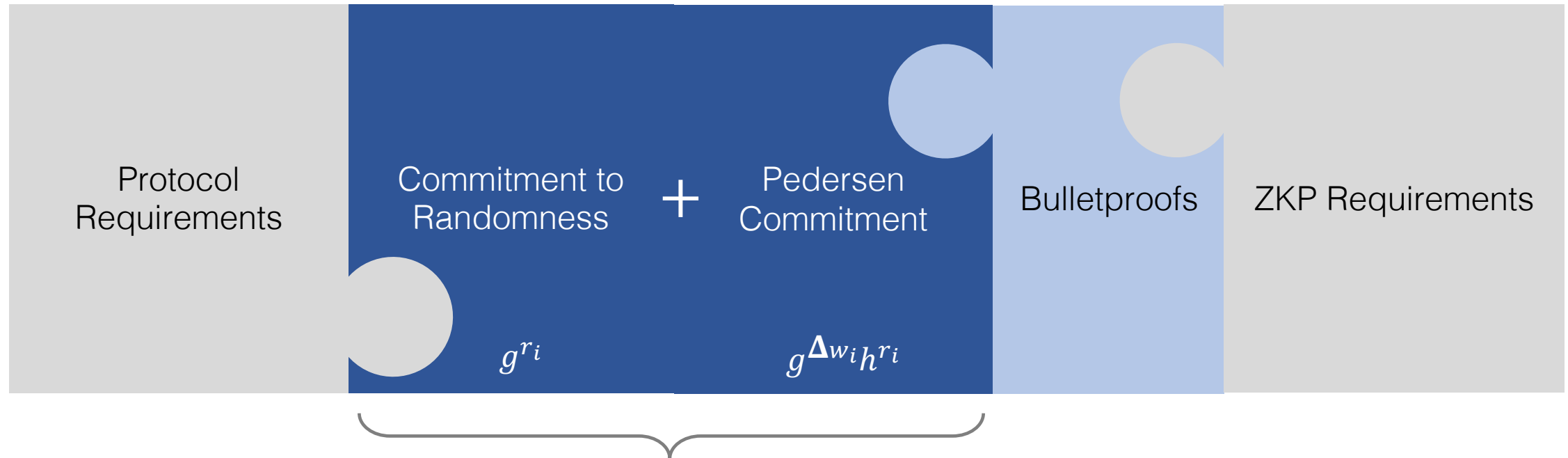
	GGPR-style zkSNARKs
Proof size	$O(1)$
Prover time	$O(\ell \log(\ell))$
Verification time	$O(1)$



Compatibility with Commitments

	GGPR-style zkSNARKs	Bulletproofs
Proof size	$O(1)$	$O(\log(\ell))$
Prover time	$O(\ell \log(\ell))$	$O(\ell)$
Verification time	$O(1)$	$O(\ell)$
Operates directly on additively homomorphic commitments	✗	✓
Specialized range proof construction	✗	✓
No trusted setup	✗	✓

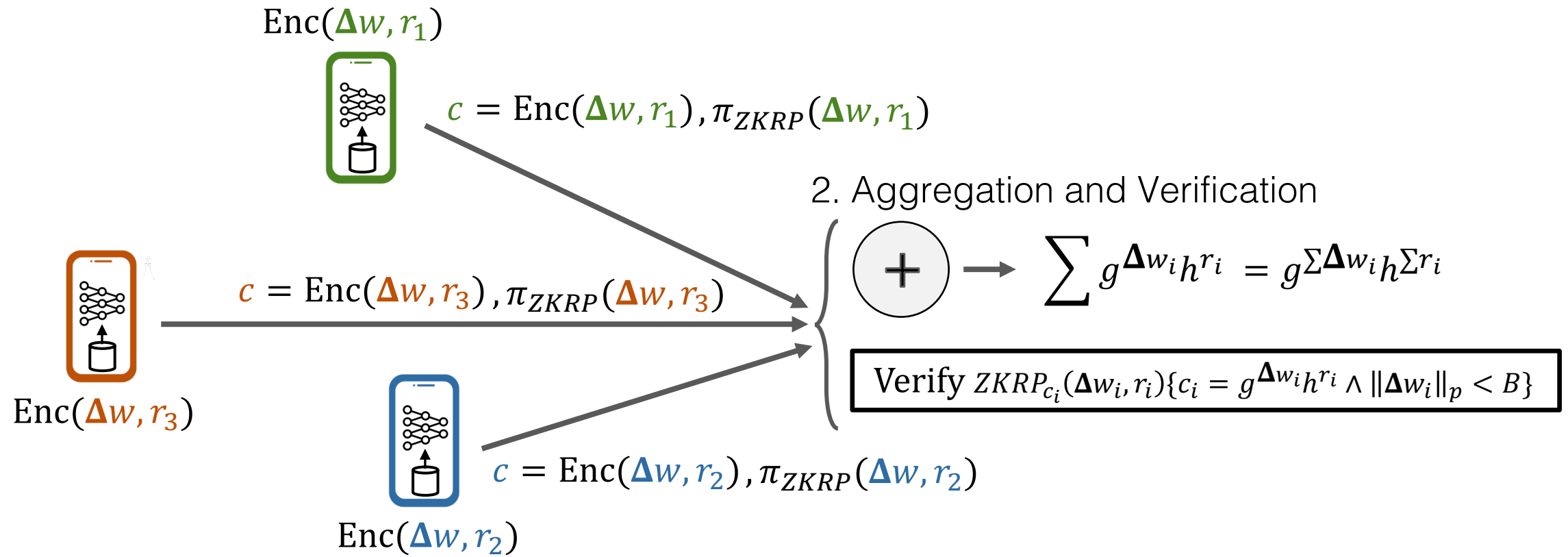
Extending Pedersen commitments for correctness



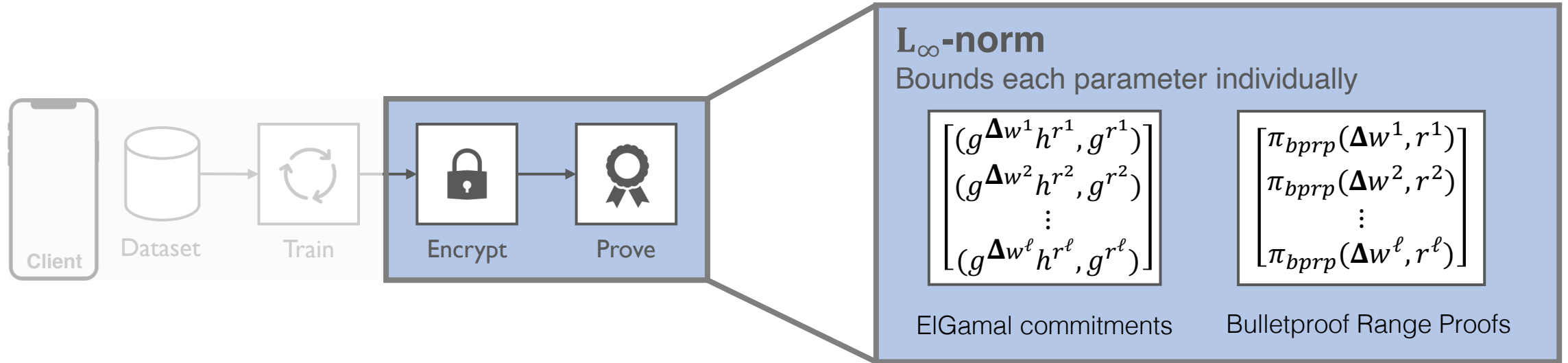
ElGamal commitment

- Server can compare $\sum g^{r_i} \leftrightarrow g^{r'}$
- Clients generate non-interactive proof-of-knowledge to proof well-formedness, i.e., r_i is the same in $(g^{\Delta w_i h^{r_i}}, g^{r_i})$

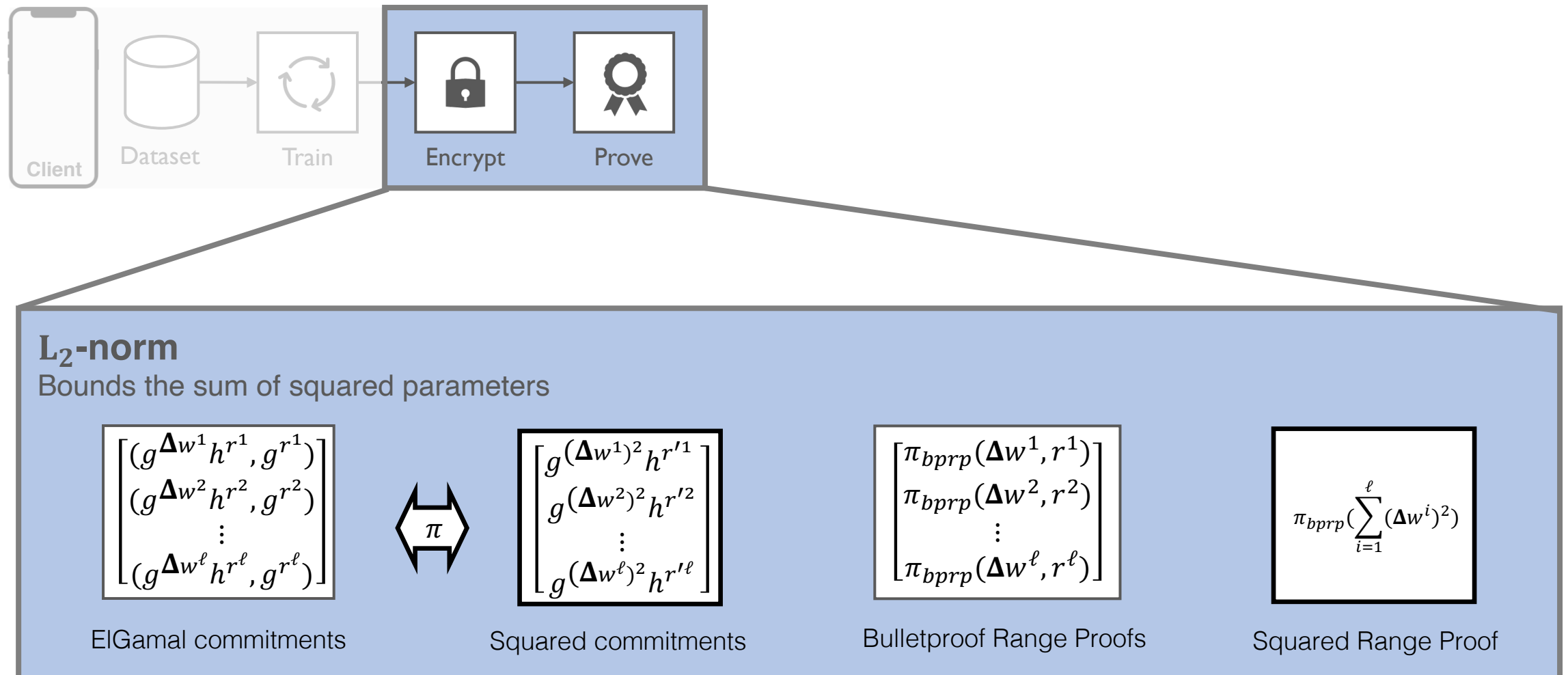
Secure Aggregation with Input Constraints



Enforcing Norm Bounds



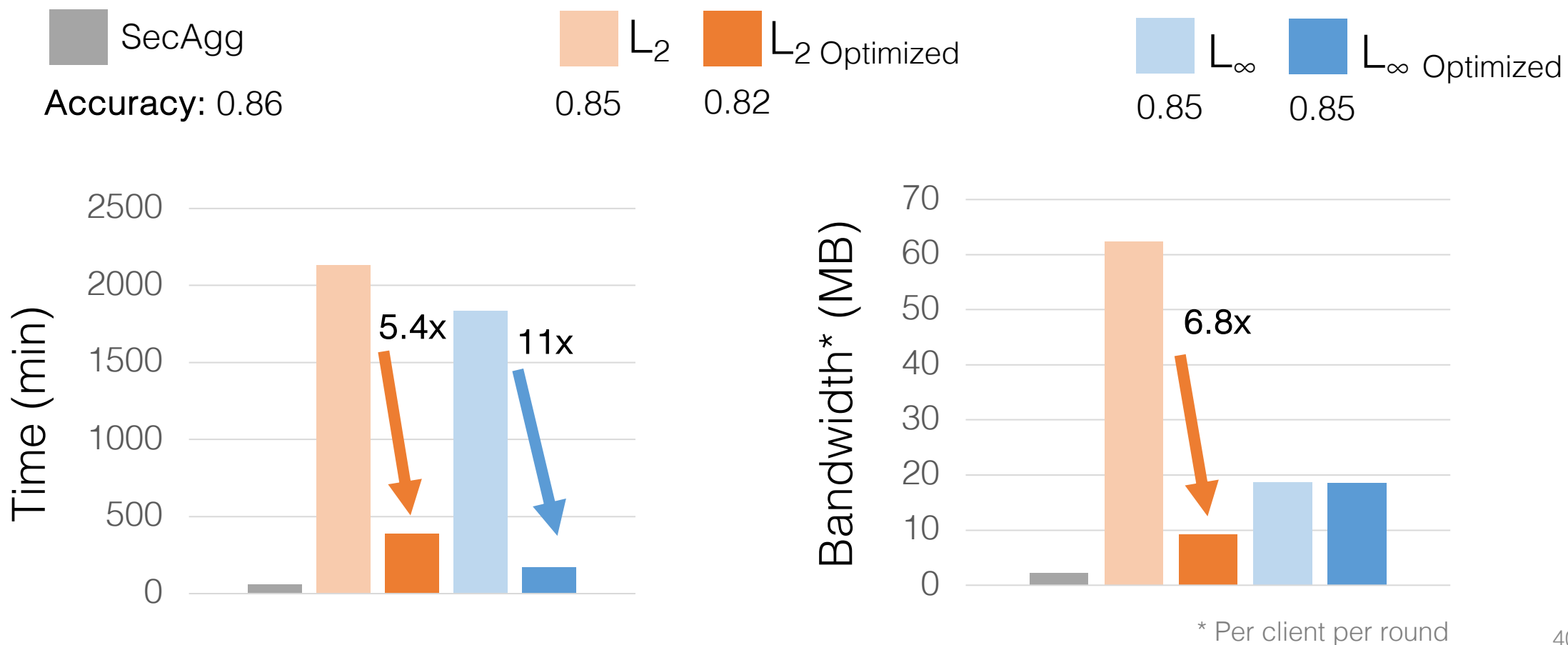
Enforcing Norm Bounds



RoFL: End-To-End Performance

CIFAR-10 Model 273k Parameters

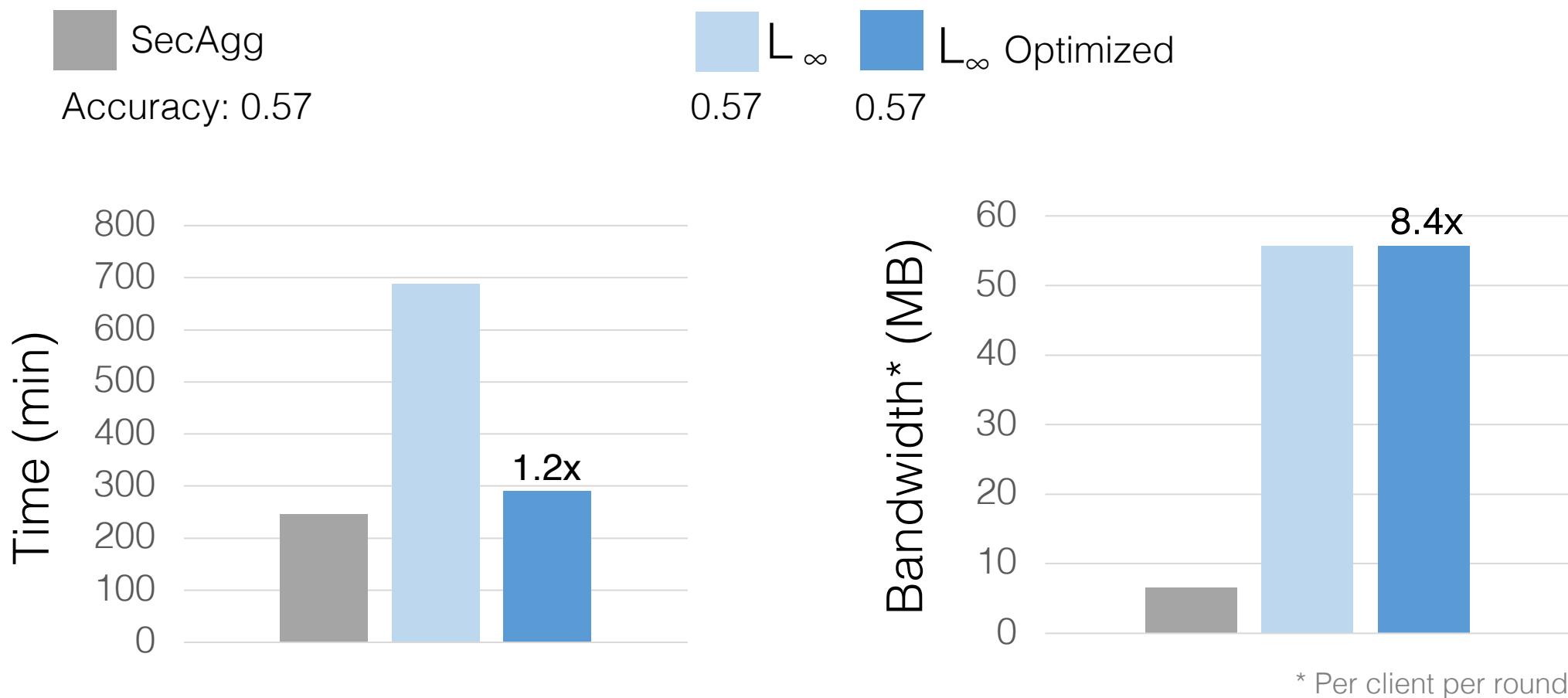
Setup: 48 Clients, 160 rounds



RoFL: End-To-End Performance

Shakespeare Model 818k Parameters

Setup: 48 Clients, 20 rounds



This work:

- Understanding FL Robustness
- RoFL: Secure Aggregation with Private Input Validation

Future work:

- Exploring additional client constraints for robustness
- Protocols with better bandwidth overhead
- Efficient ZKPs for resource-constrained provers



pps-lab/fl-analysis



pps-lab/rofl-project-code



pps-lab.com/research/ml-sec



