

Holding Secrets Accountable: Auditing Private ML Algorithms



Hidde Lycklama
ETHzürich



Alexander Viand
intel



Nicolas Küchler
ETHzürich

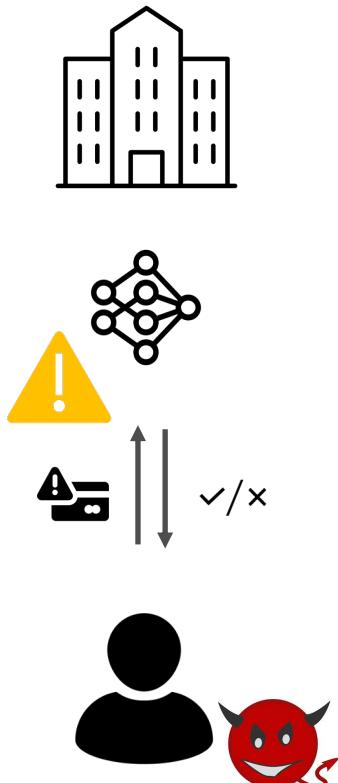


Christian Knabenhans
EPFL

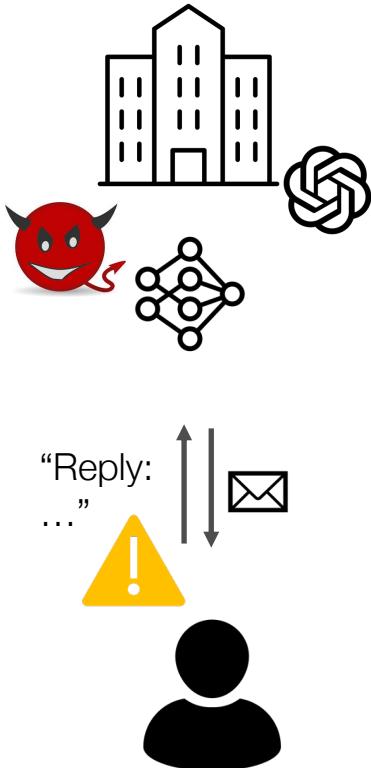


Anwar Hithnawi
ETHzürich

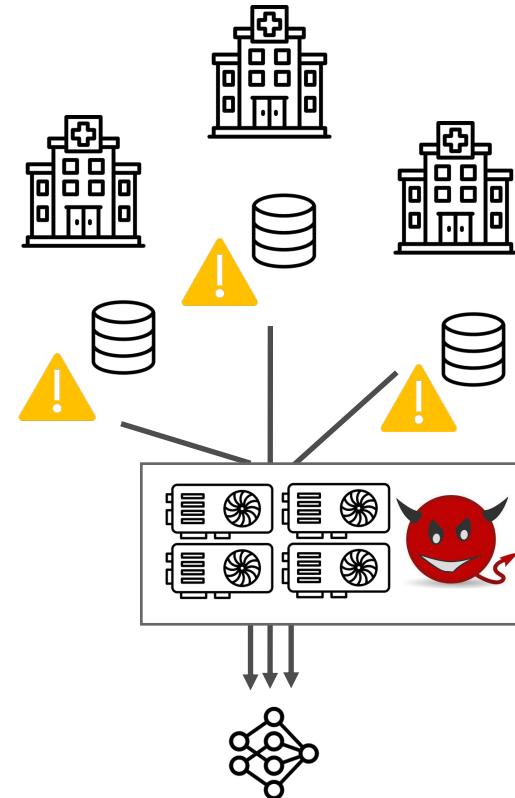
Privacy-Preserving Machine Learning (PPML)



Model Privacy

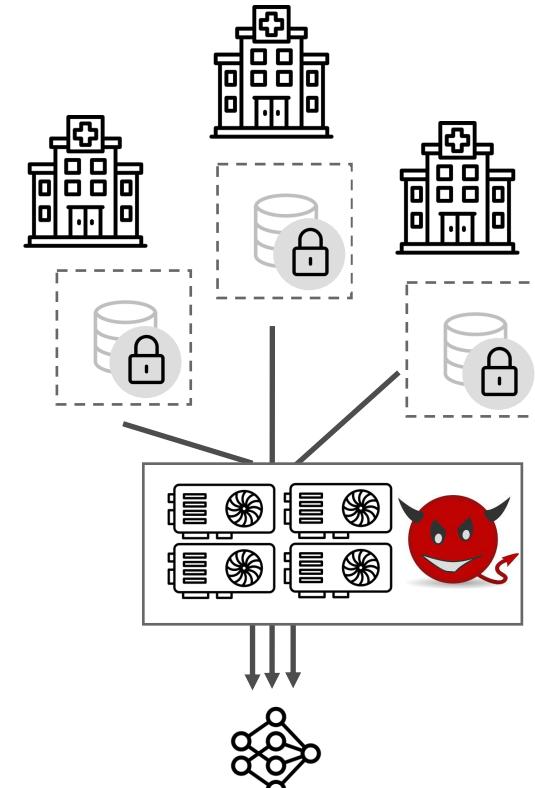
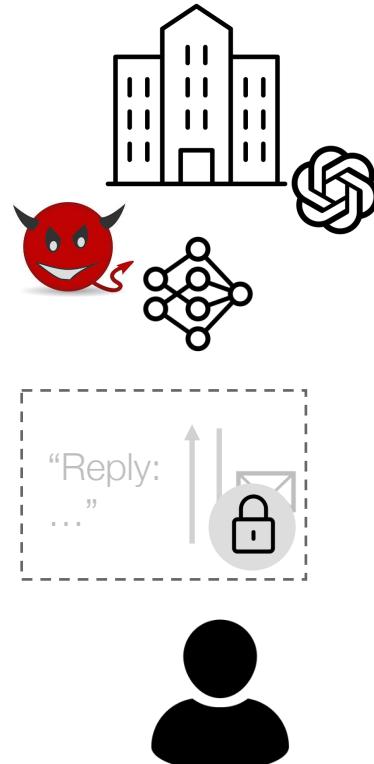
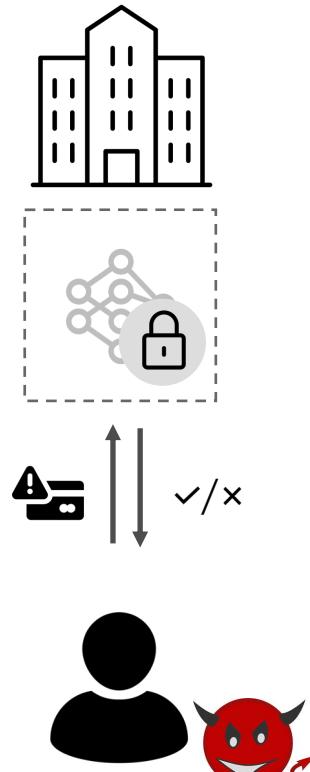


Query Privacy



Data Privacy

Privacy-Preserving Machine Learning (PPML)



Model Privacy

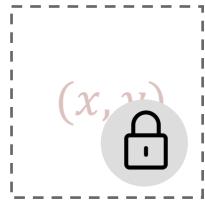
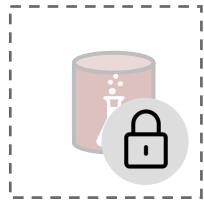
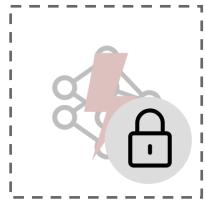
Query Privacy

Data Privacy

Privacy-Transparency Dichotomy



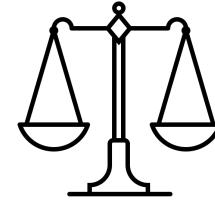
Privacy



Transparency



Accountability

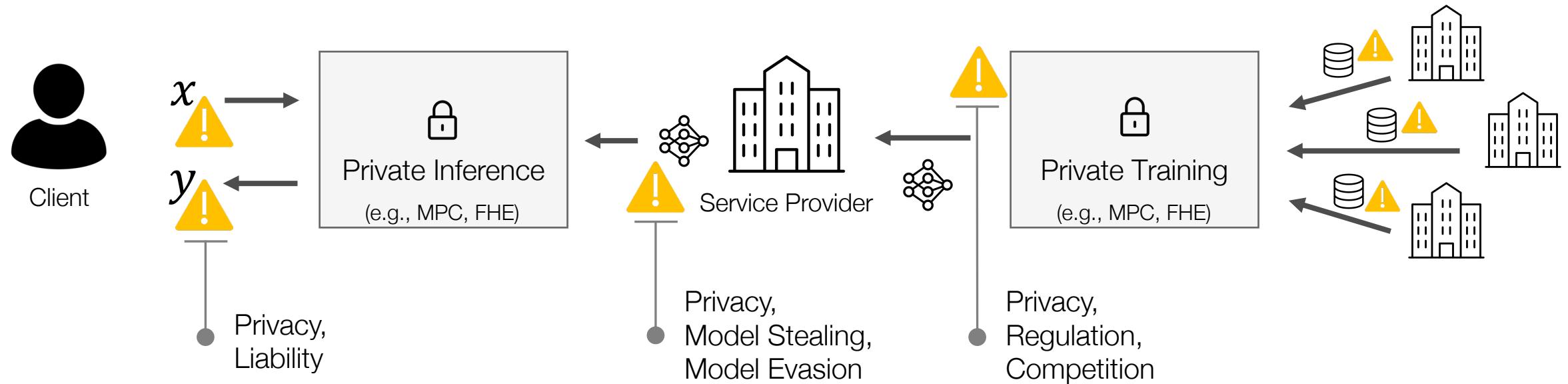


Fairness

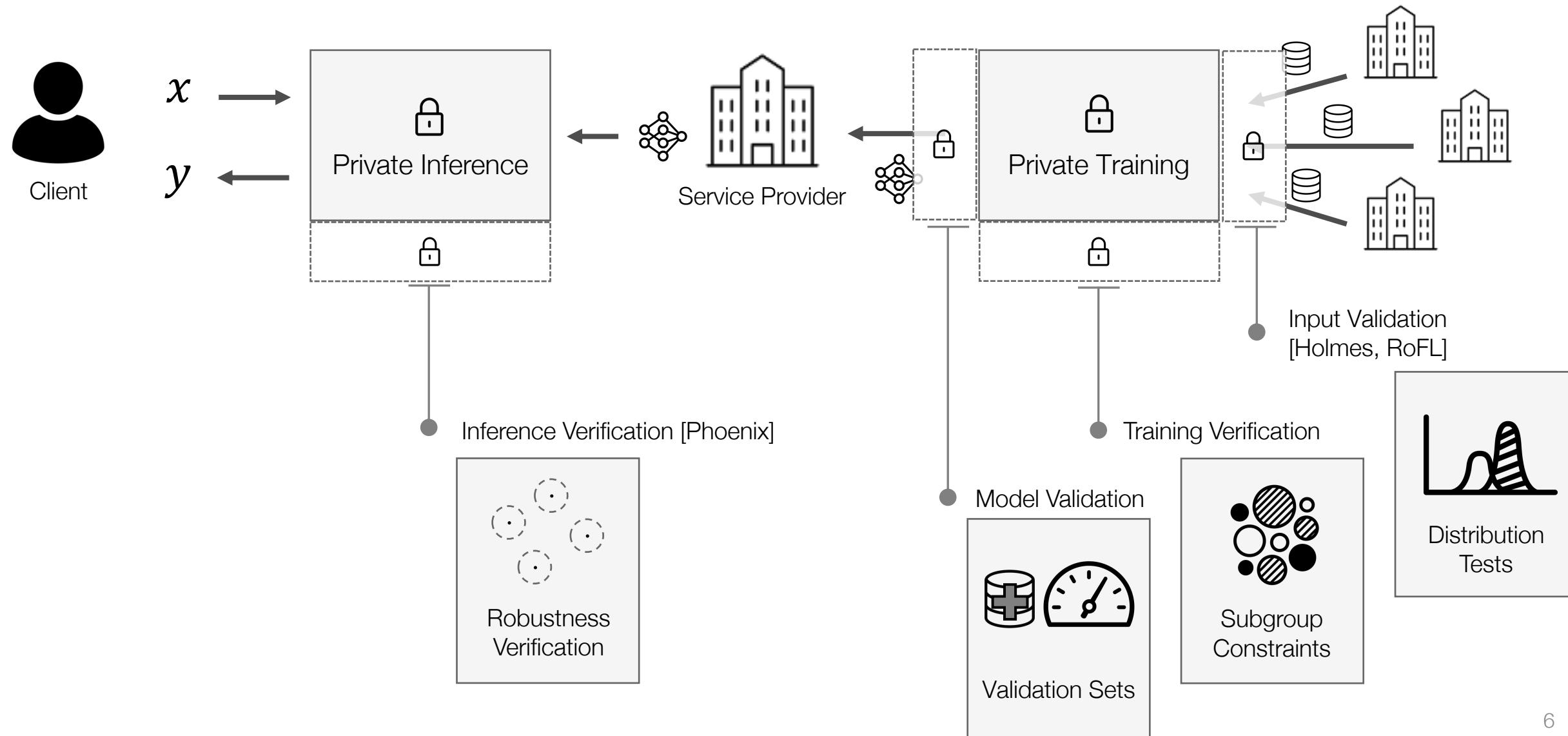


Safety

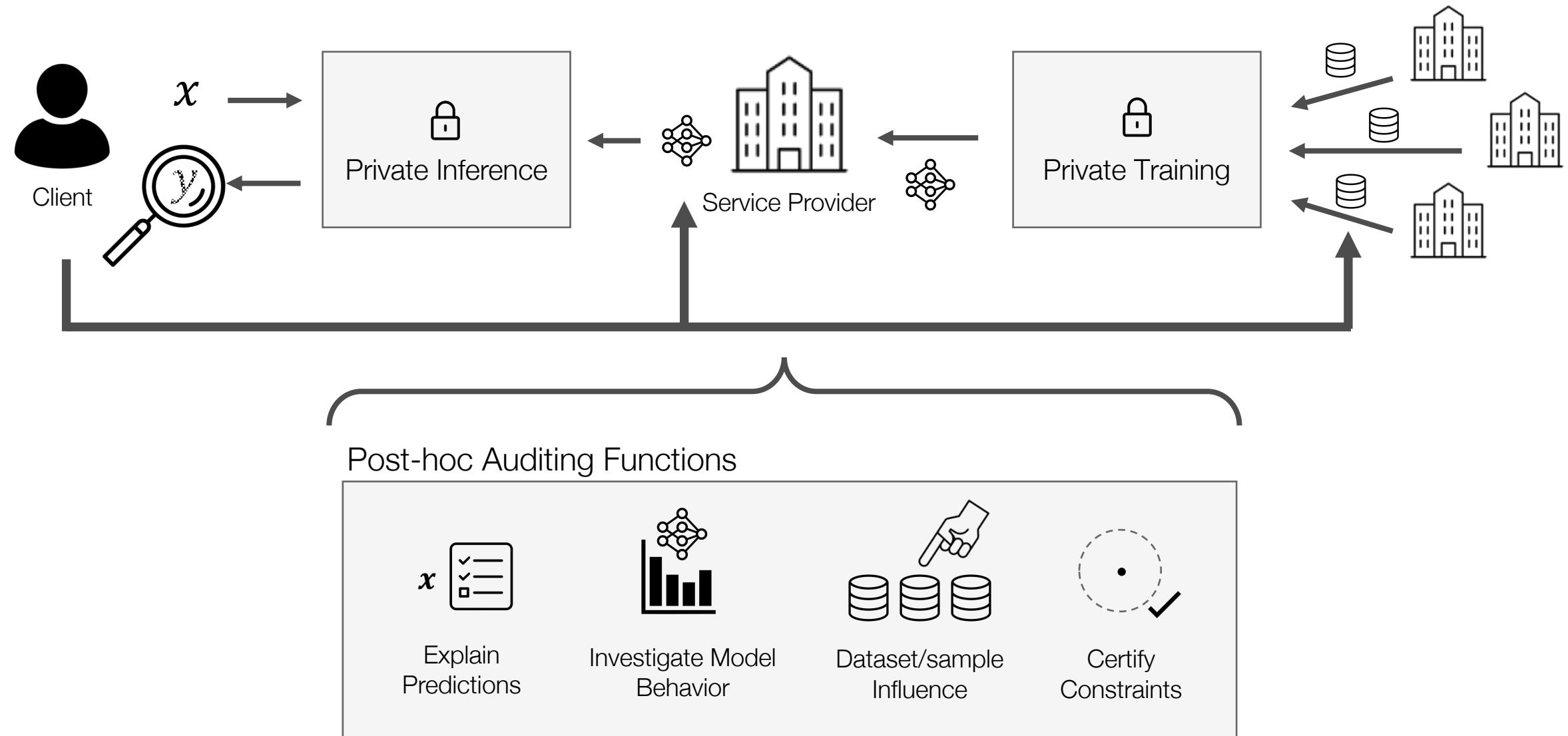
Privacy-Preserving Machine Learning (PPML)



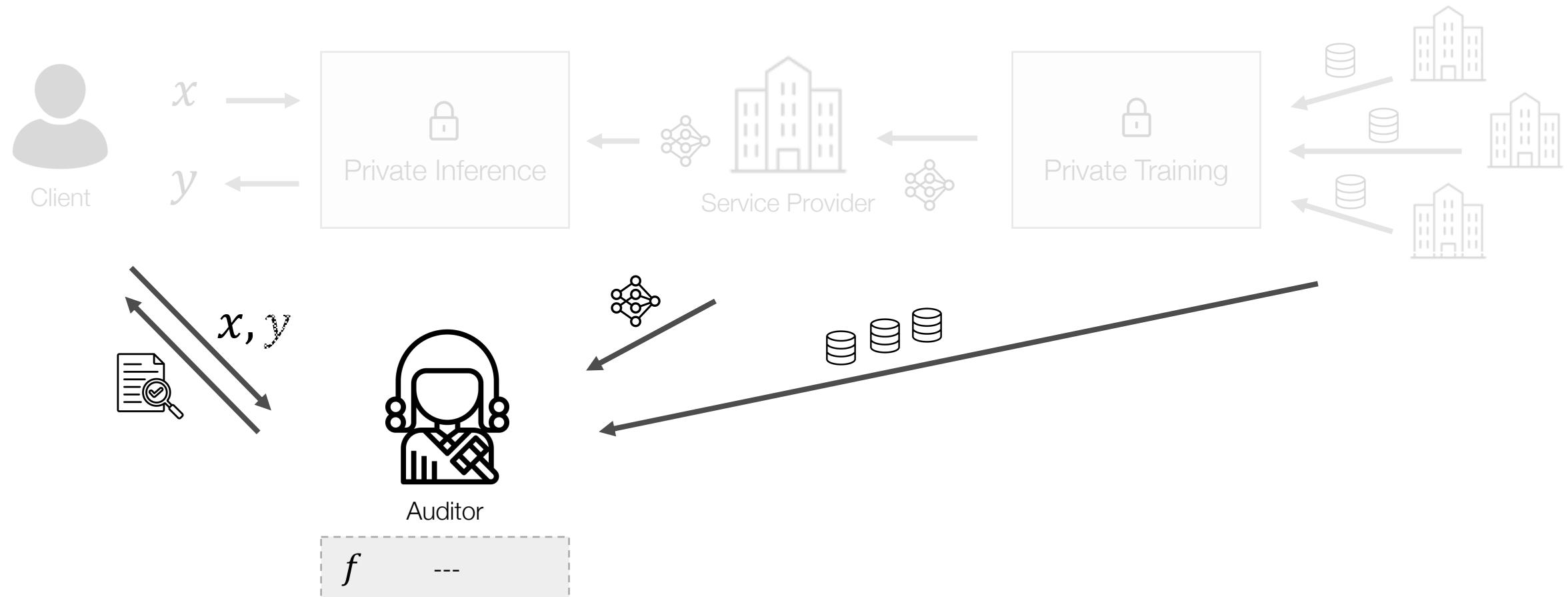
Auditing Privacy-Preserving Machine Learning



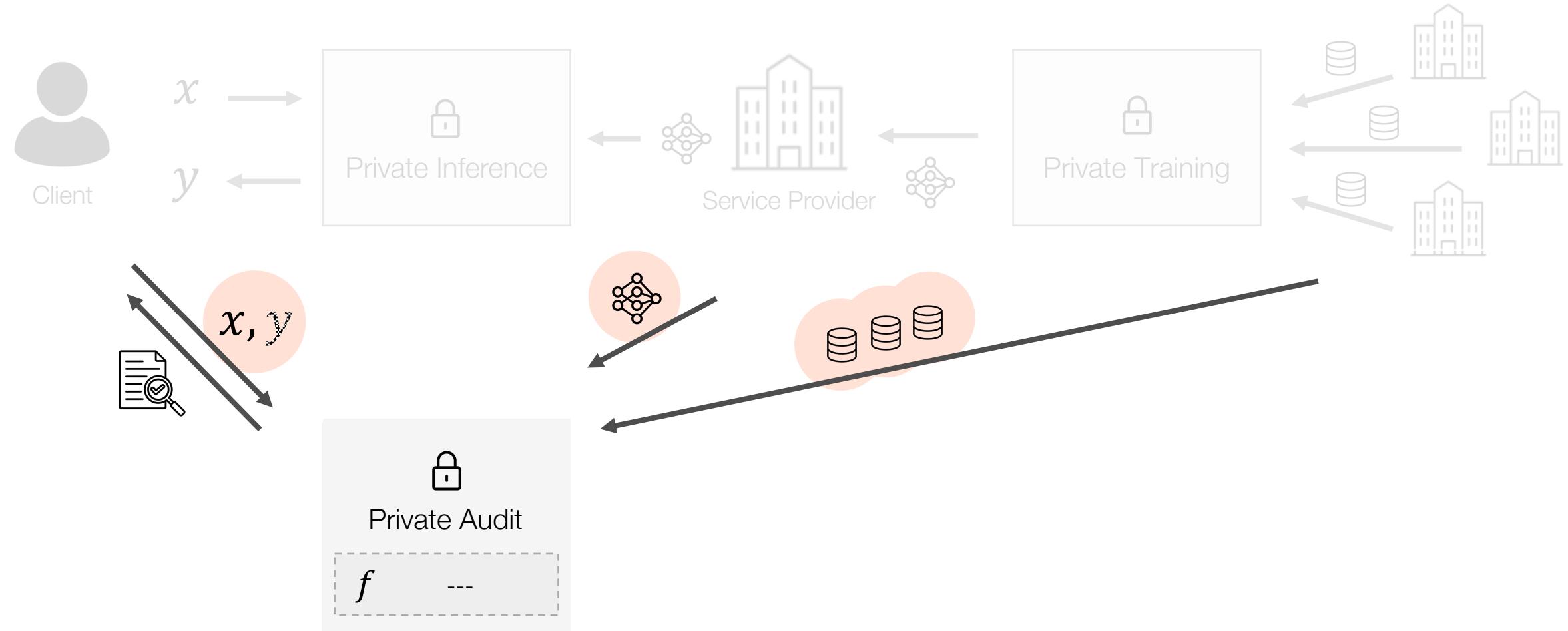
Post-Hoc Auditing of PPML



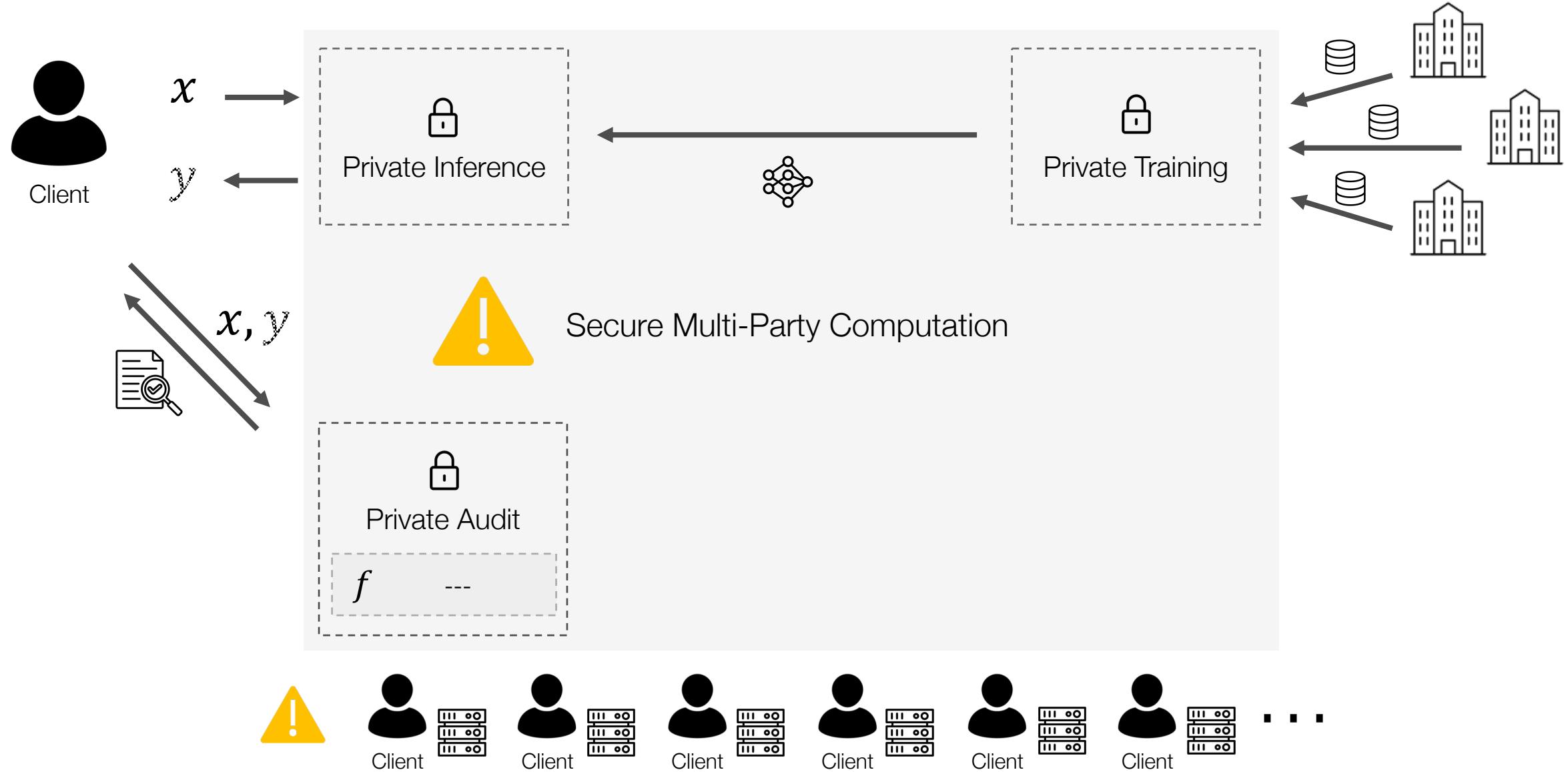
Post-Hoc Auditing of PPML



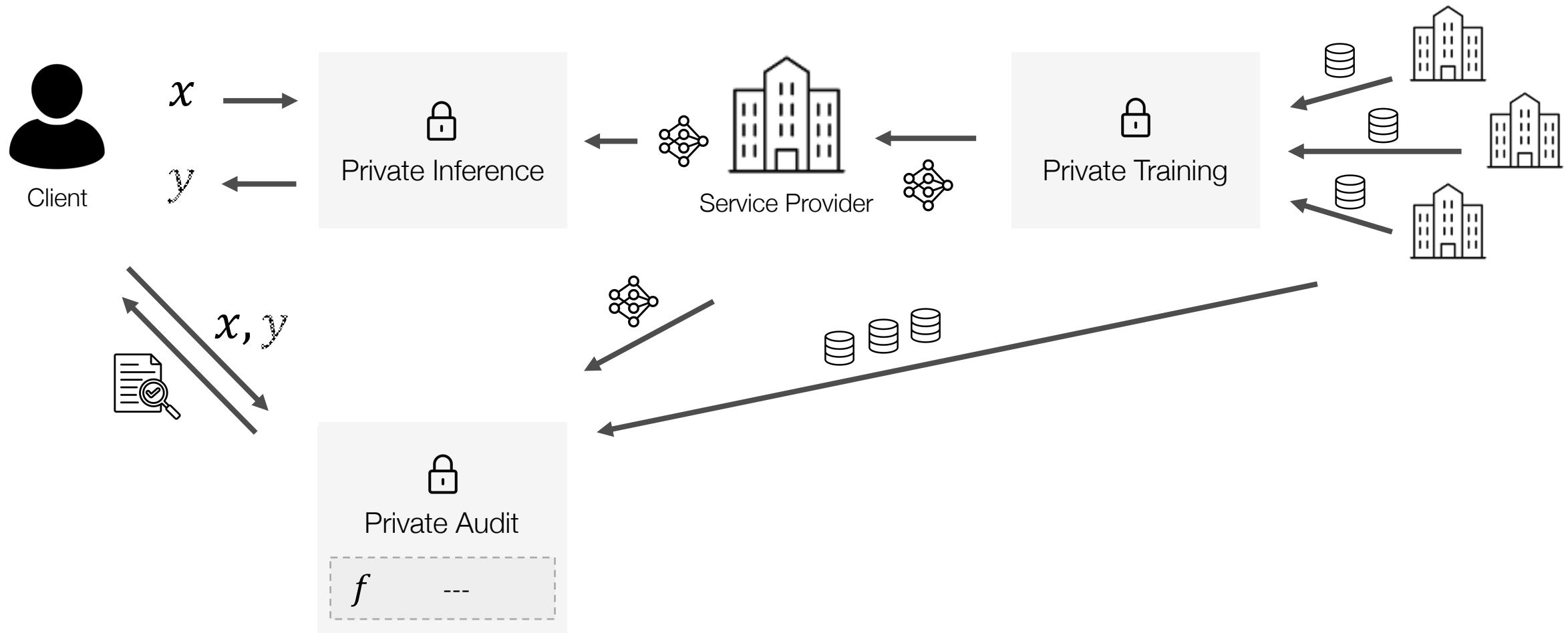
Post-Hoc Auditing of PPML



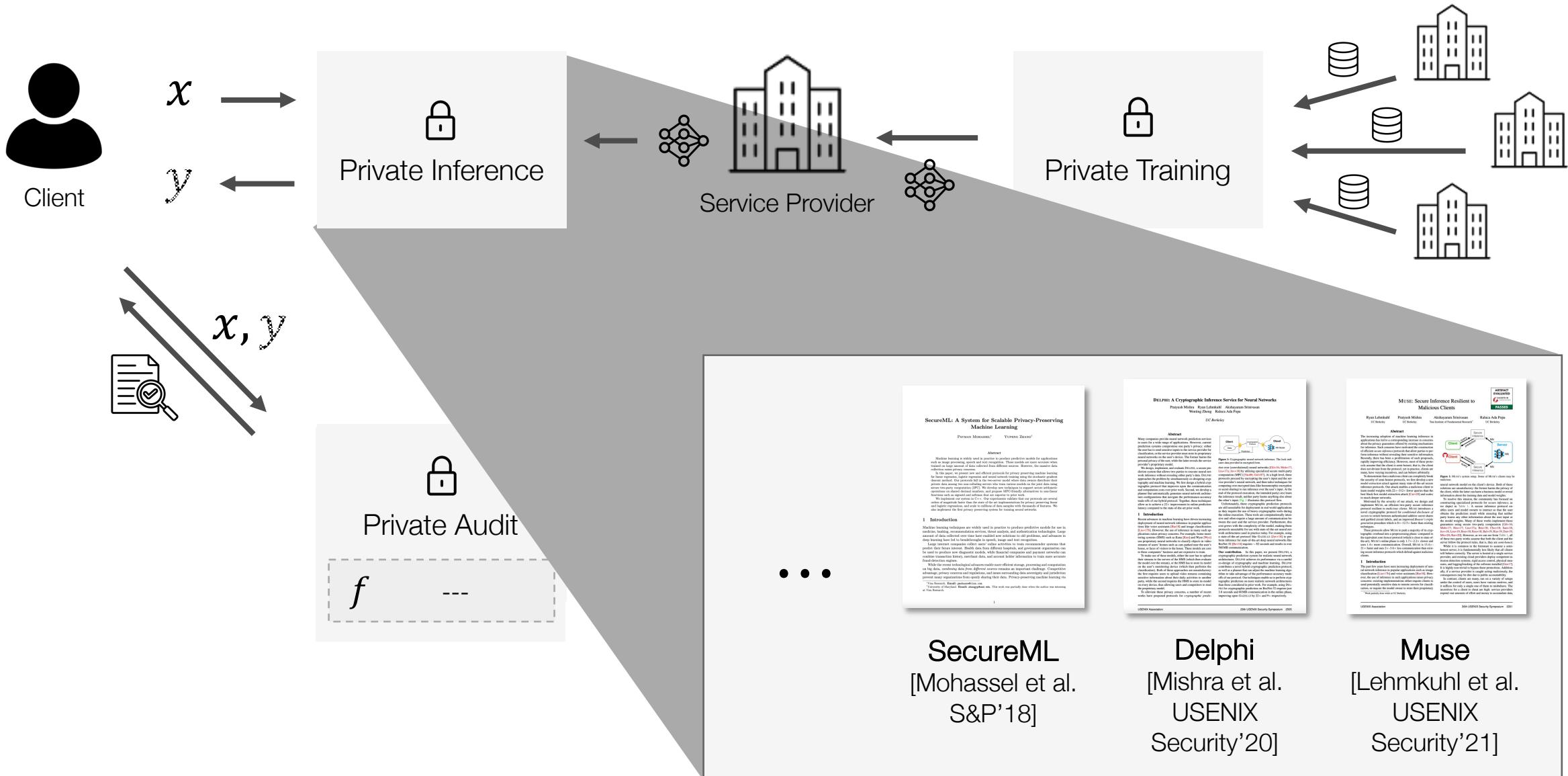
Post-Hoc Auditing of PPML



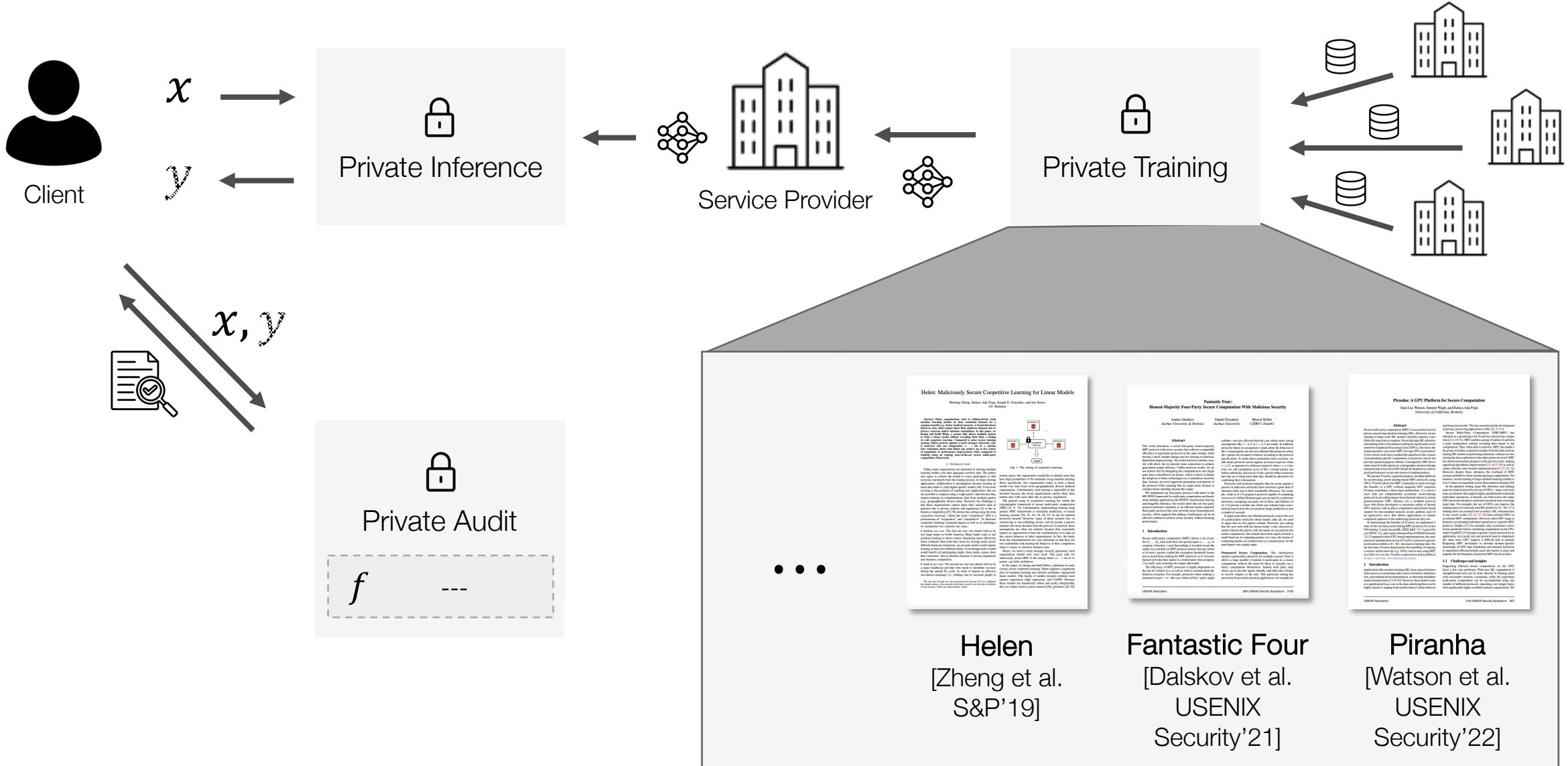
Post-Hoc Auditing of PPML: MPC Phases



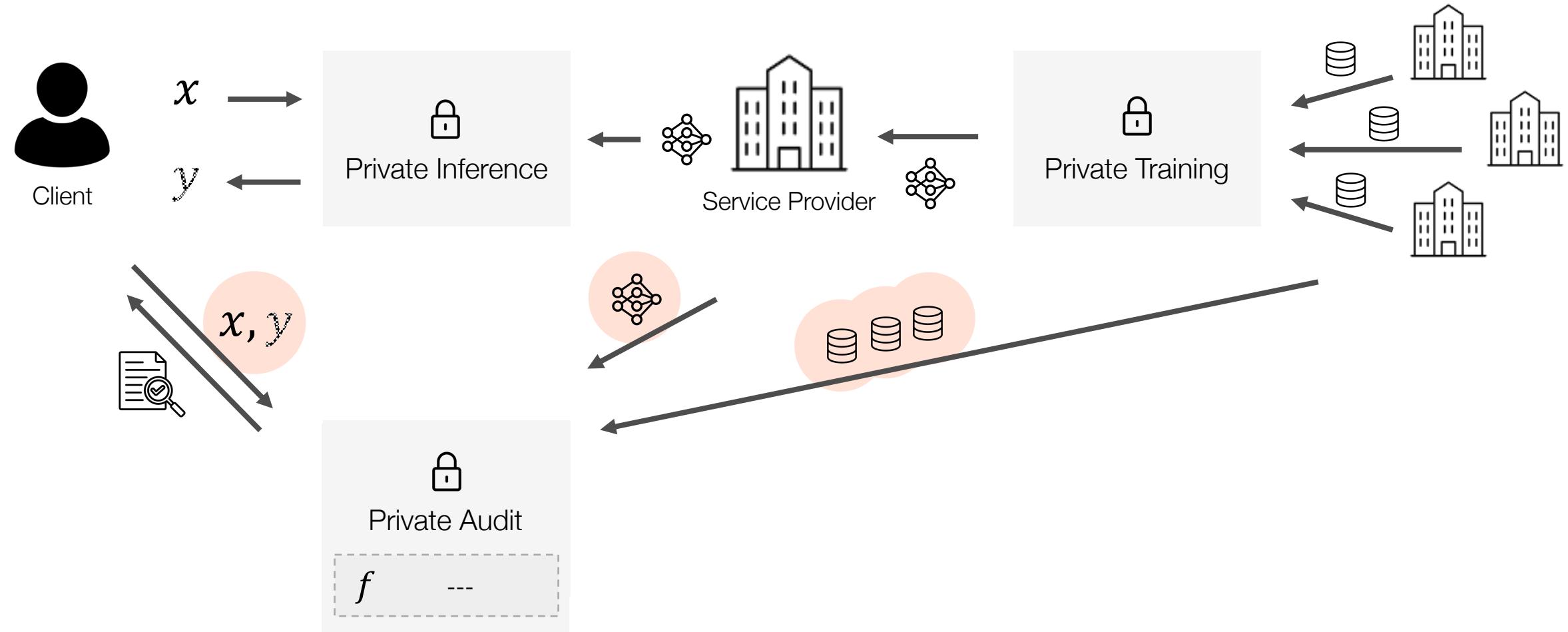
Post-Hoc Auditing of PPML: MPC Phases



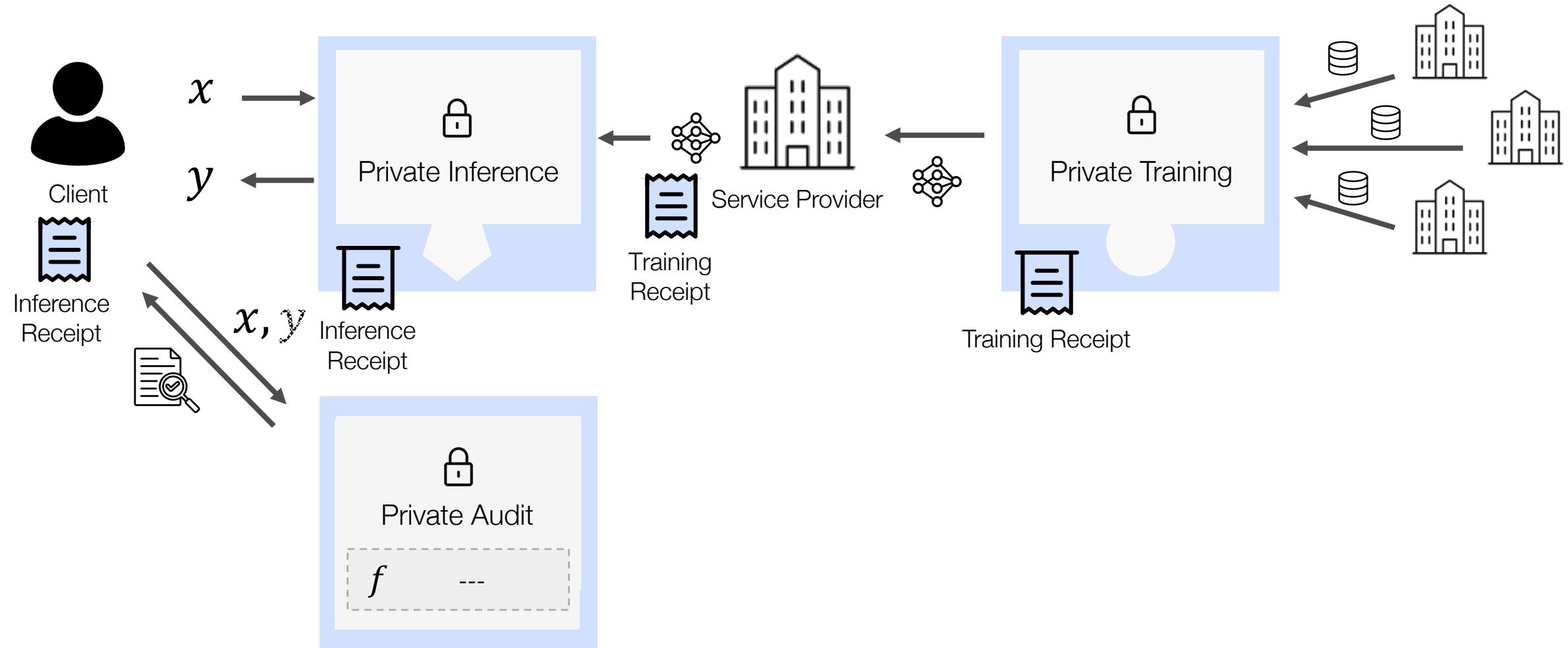
Post-Hoc Auditing of PPML: MPC Phases



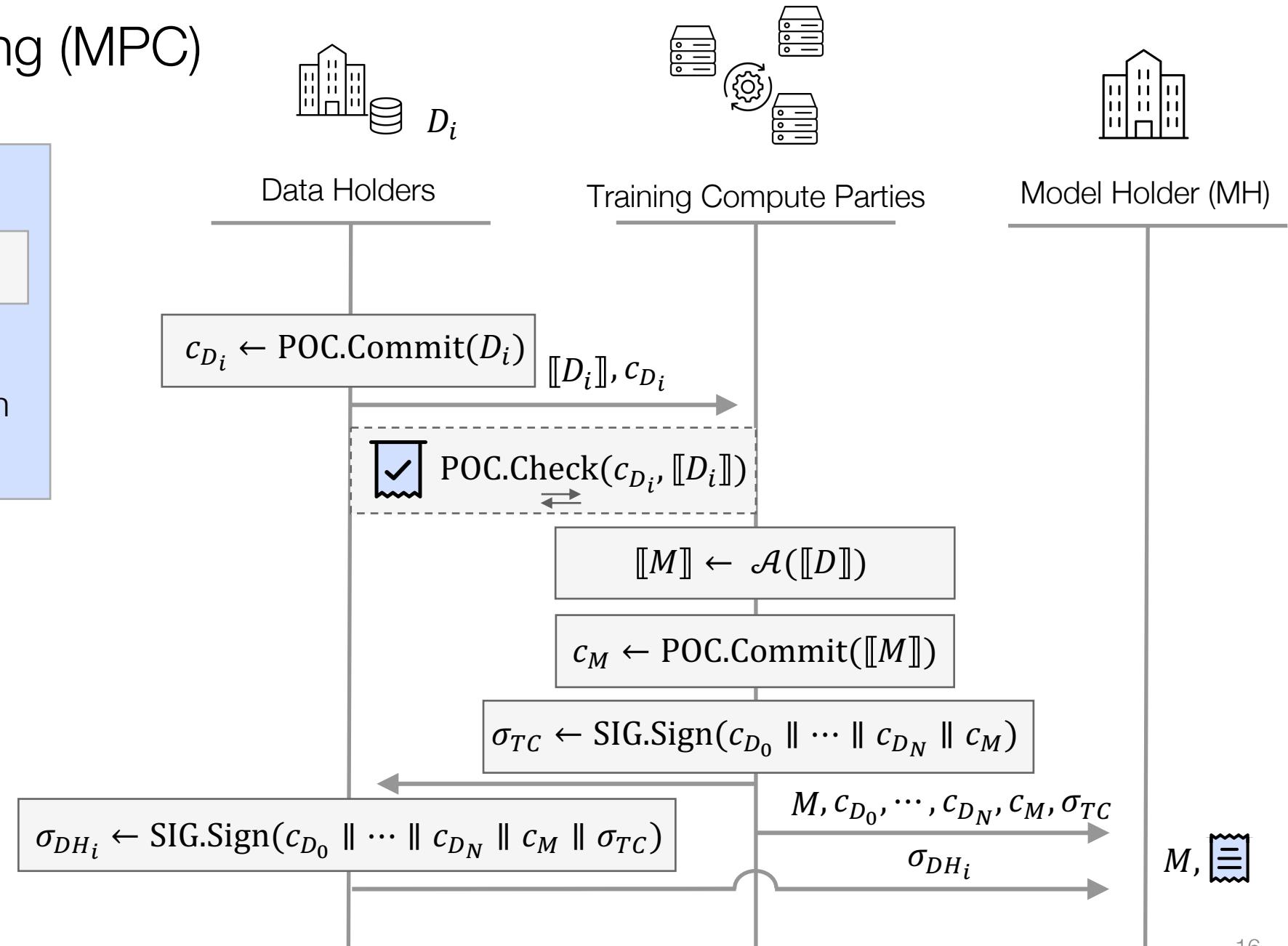
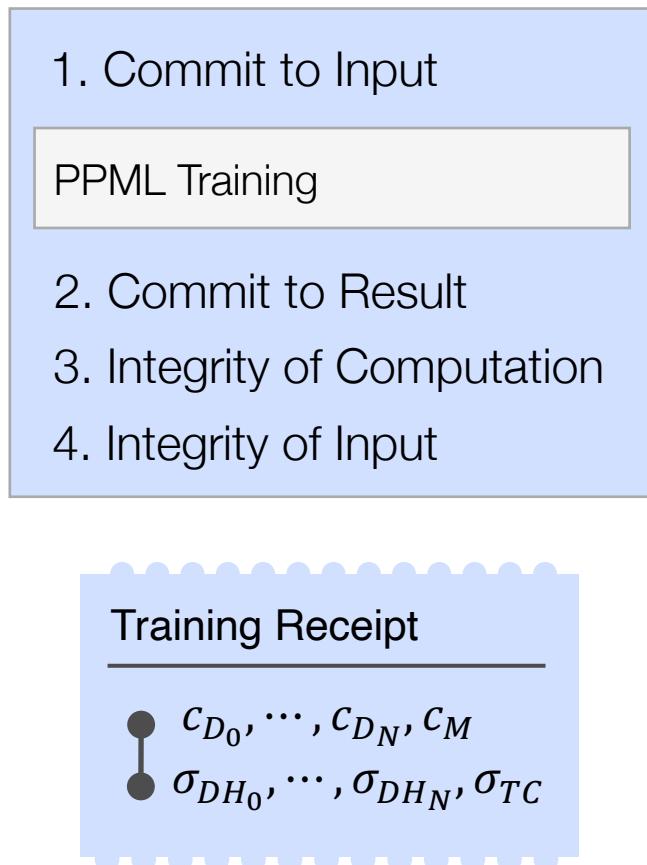
Post-Hoc Auditing of PPML



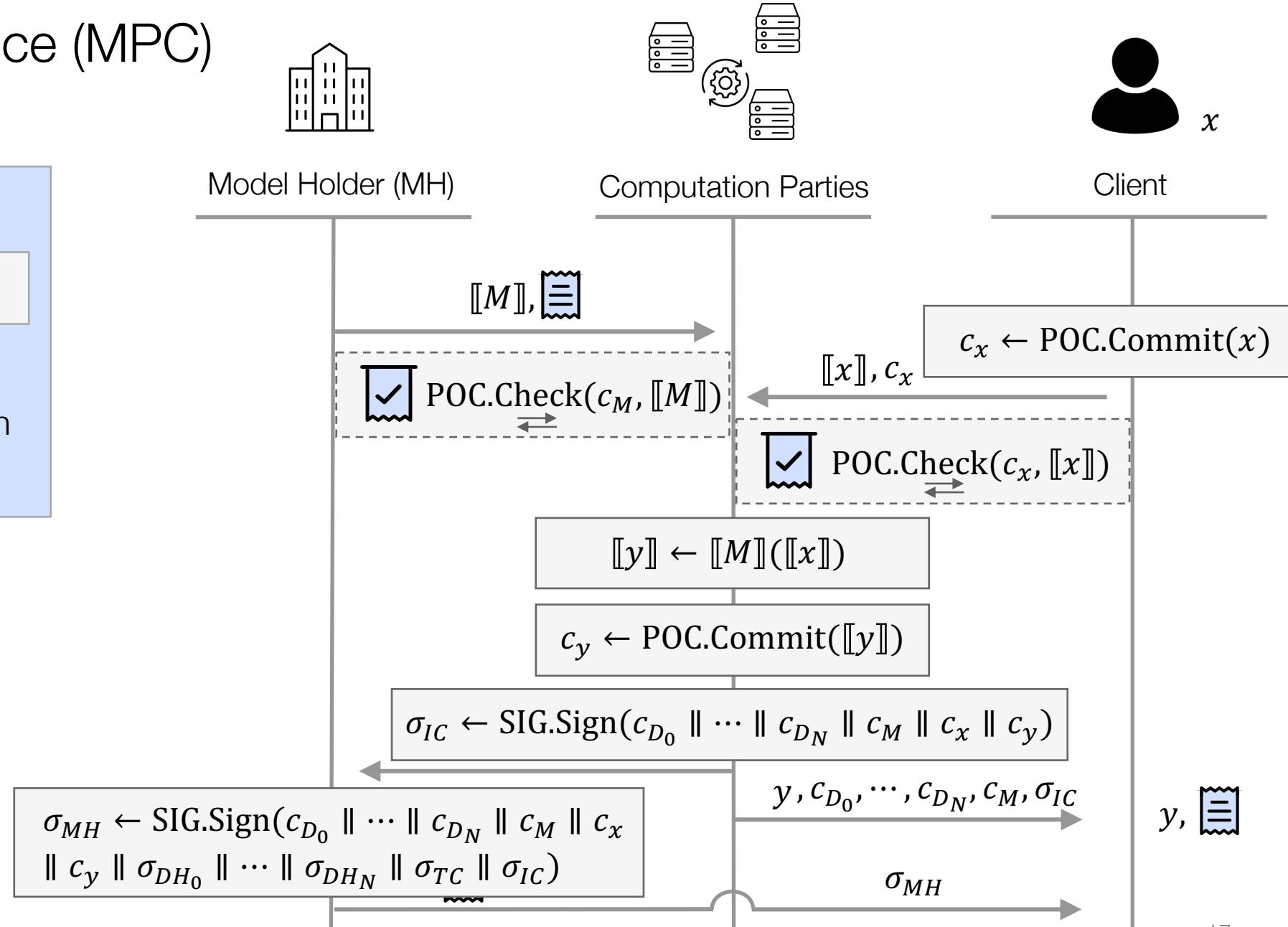
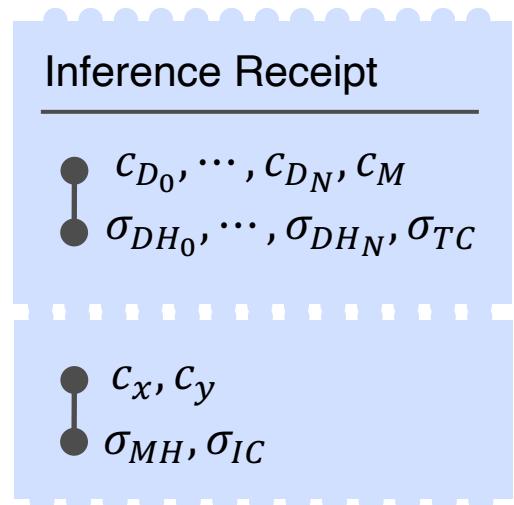
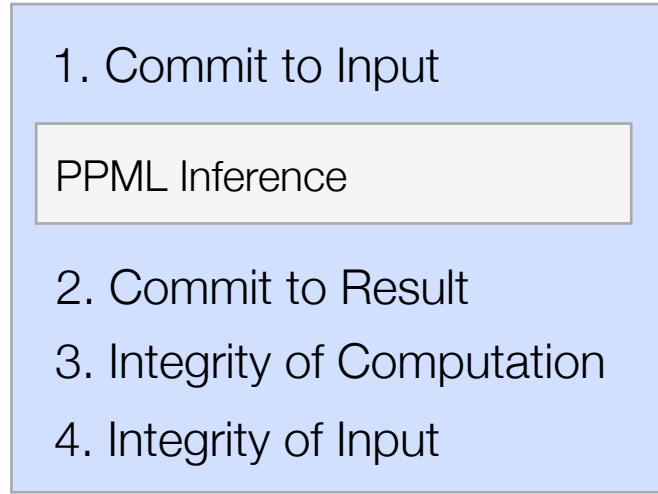
This Work: Arc, a Framework for End-to-End Auditing of PPML



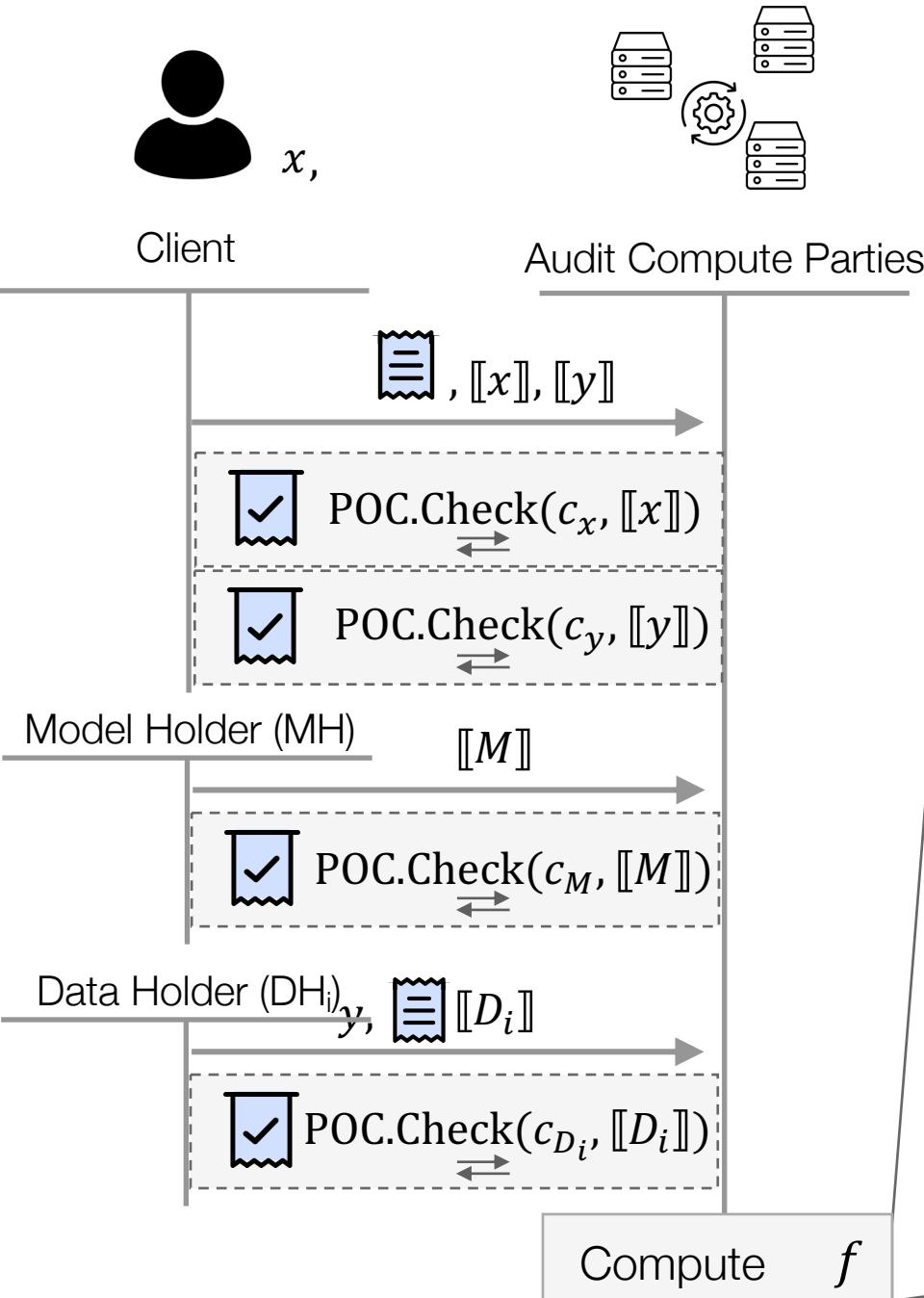
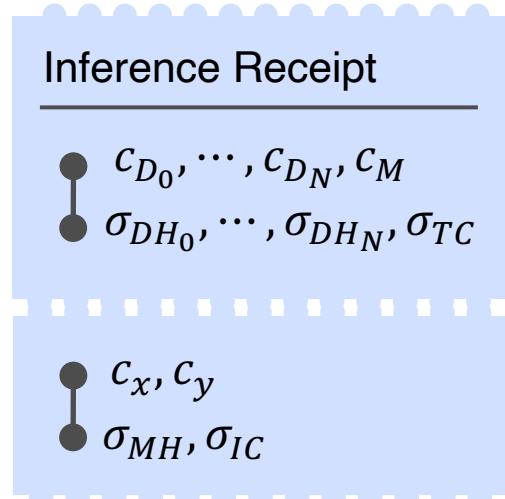
Arc Protocol: Training (MPC)



Arc Protocol: Inference (MPC)



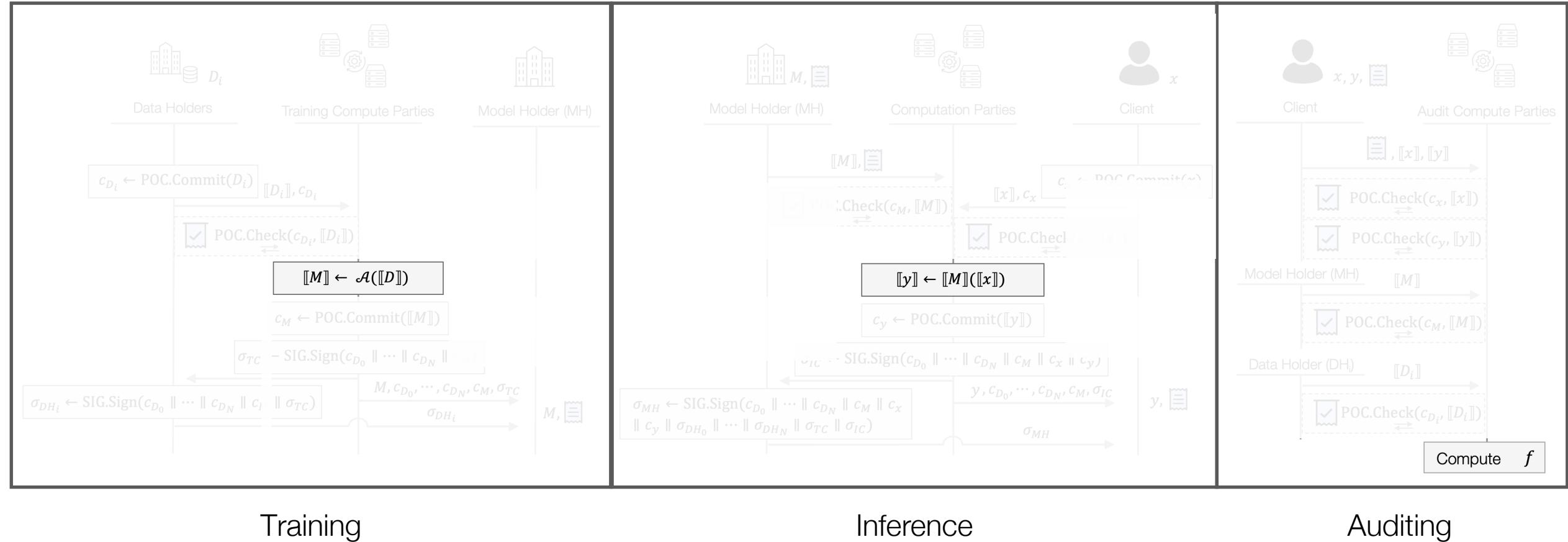
Arc Protocol: Auditing



Arithmetized Auditing Functions

- | | |
|----------------|---------------------|
| Fairness | f Indiv. Fairness |
| Safety | f Robustness |
| Transparency | f KernelSHAP |
| Accountability | f kNNShapley |

Arc Protocol

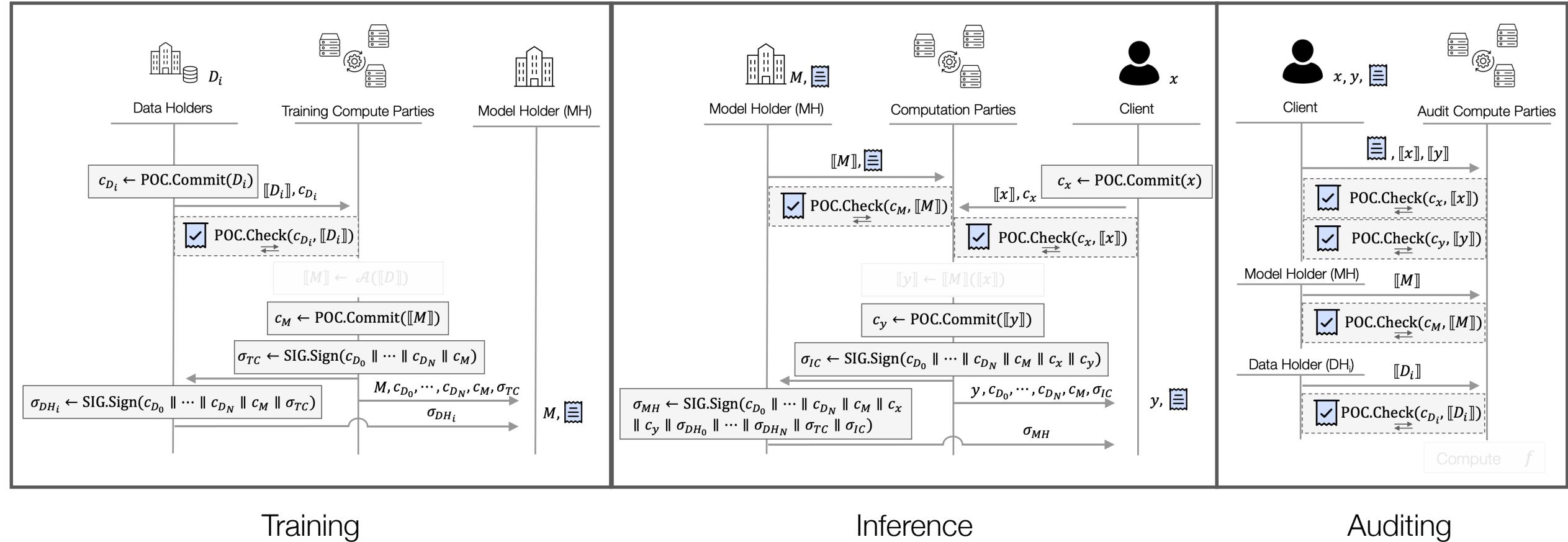


Training

Inference

Auditing

Arc Protocol

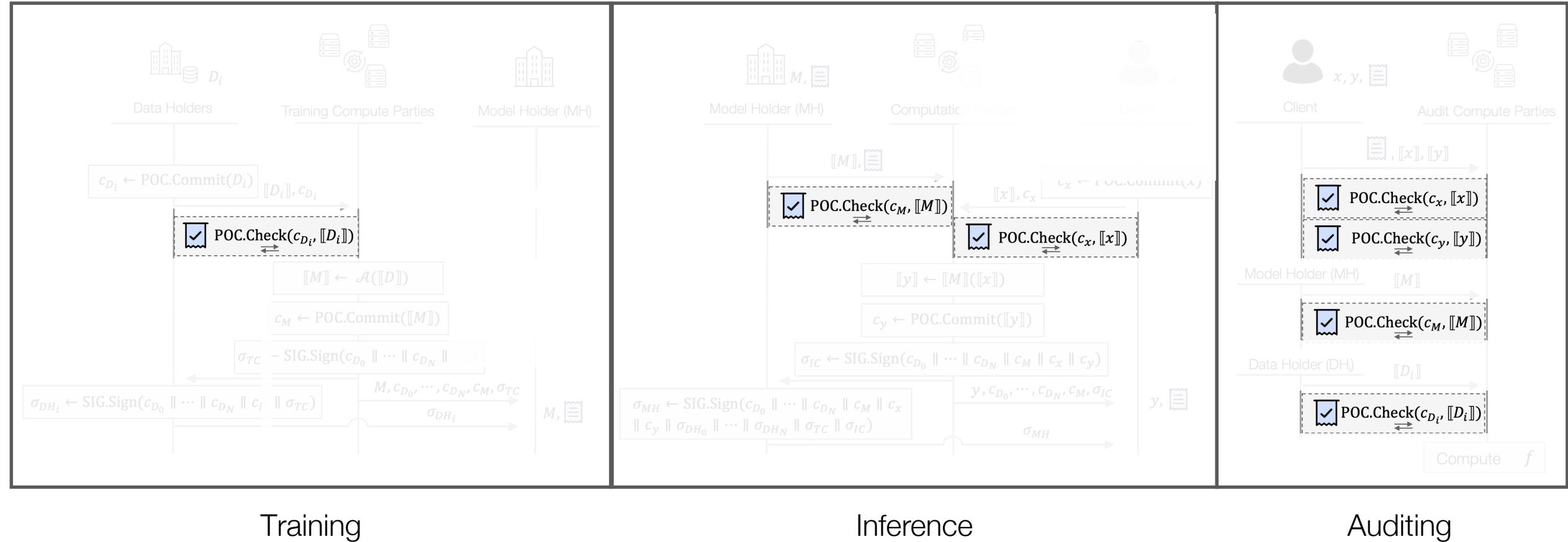


Training

Inference

Auditing

Arc Protocol

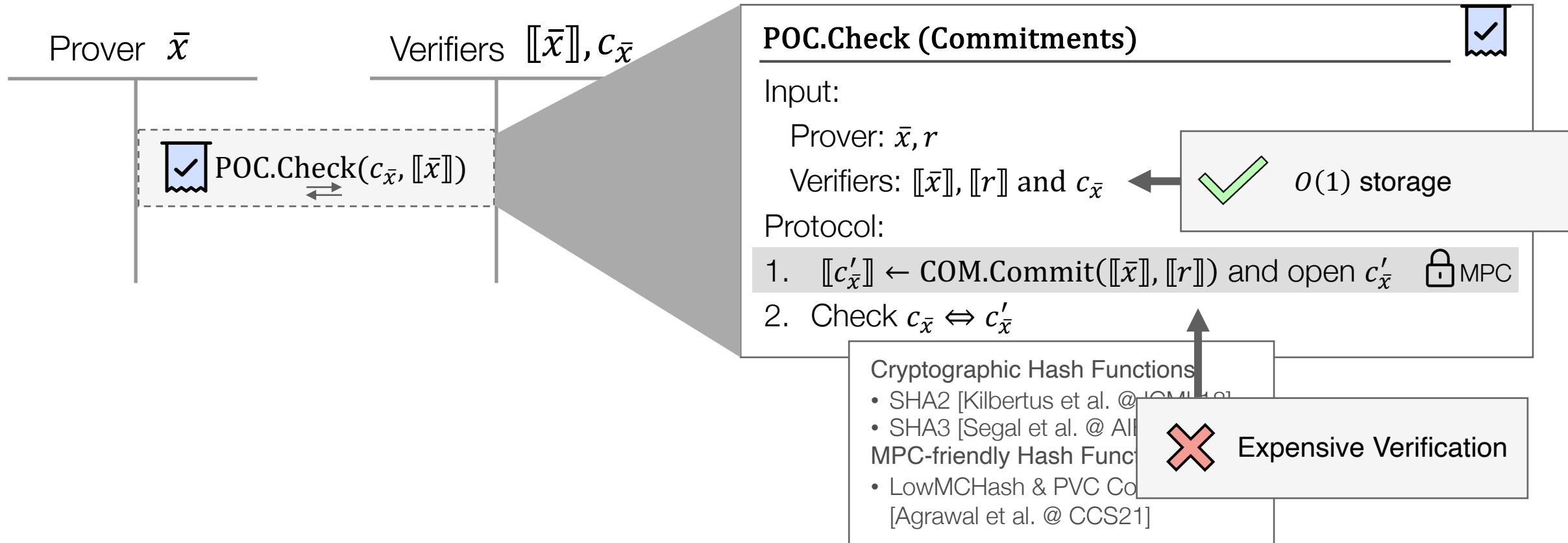


Training

Inference

Auditing

Proof of Consistency



Proof of Consistency

POC.Check (Homomorphic Commitments)

Input:

Prover: $\bar{x} = [x_1 \dots x_d], r = [r_1 \dots r_d]$

Verifiers: $\llbracket \bar{x} \rrbracket, \llbracket r \rrbracket$ and $c_{\bar{x}} = [c_{x_1} \dots c_{x_d}]$

Protocol:

1. Sample random challenge $\beta \stackrel{\$}{\leftarrow} \mathbb{F}_p$
2. Compute $\llbracket \tilde{x} \rrbracket = \sum_{i=0} \llbracket x_{i+1} \rrbracket \cdot \beta^i$,
 $\llbracket \tilde{r} \rrbracket = \sum_{i=0} \llbracket r_{i+1} \rrbracket \cdot \beta^i$
3. Compute $\llbracket \tilde{c}' \rrbracket \leftarrow \text{PED.Com}(\llbracket \tilde{x} \rrbracket, \llbracket \tilde{r} \rrbracket)$ and open \tilde{c}'
4. Check $\sum_{i=0} c_{x_{i+1}} \cdot \beta^i \Leftrightarrow \tilde{c}'$



$O(d)$ storage



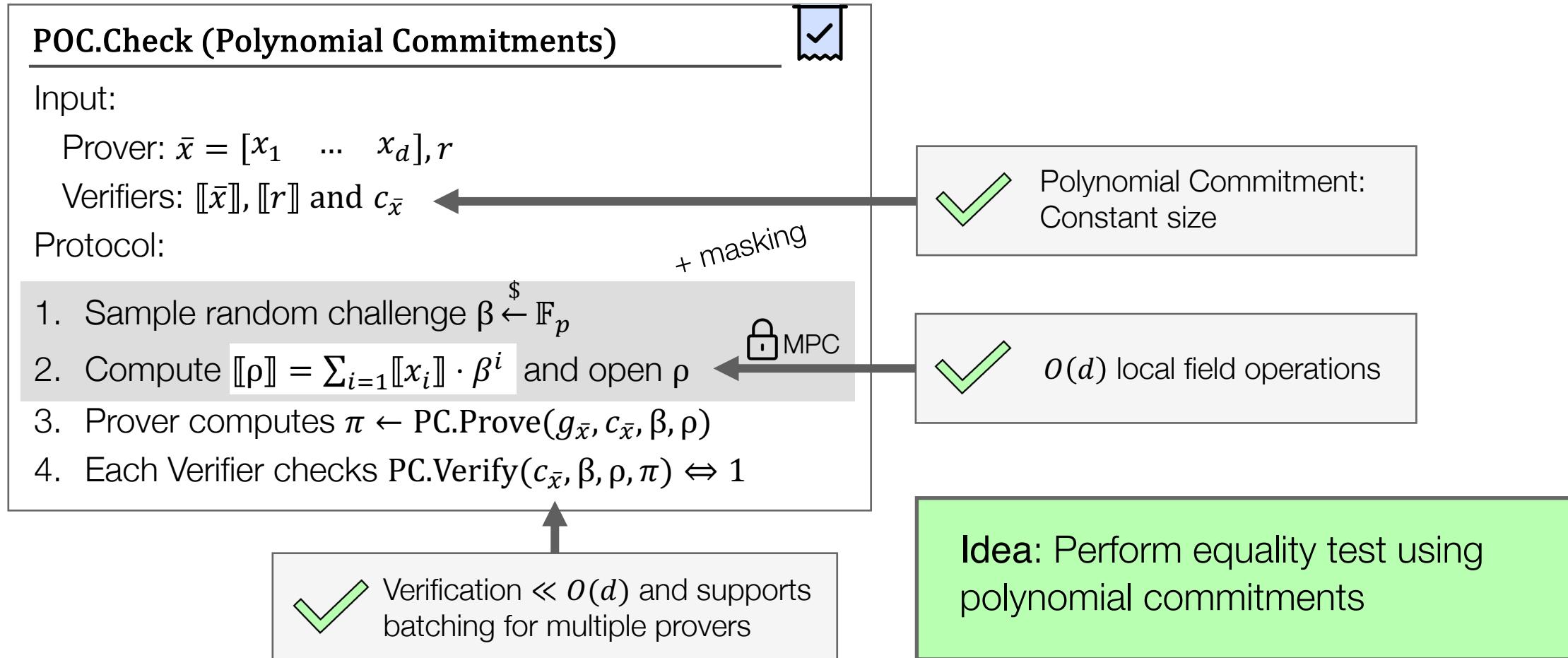
$O(1)$ EC-MPC computation



$O(d)$ local EC computation

Cerebro [Zheng et al. @ USENIX Sec'21]

Our Proof of Consistency

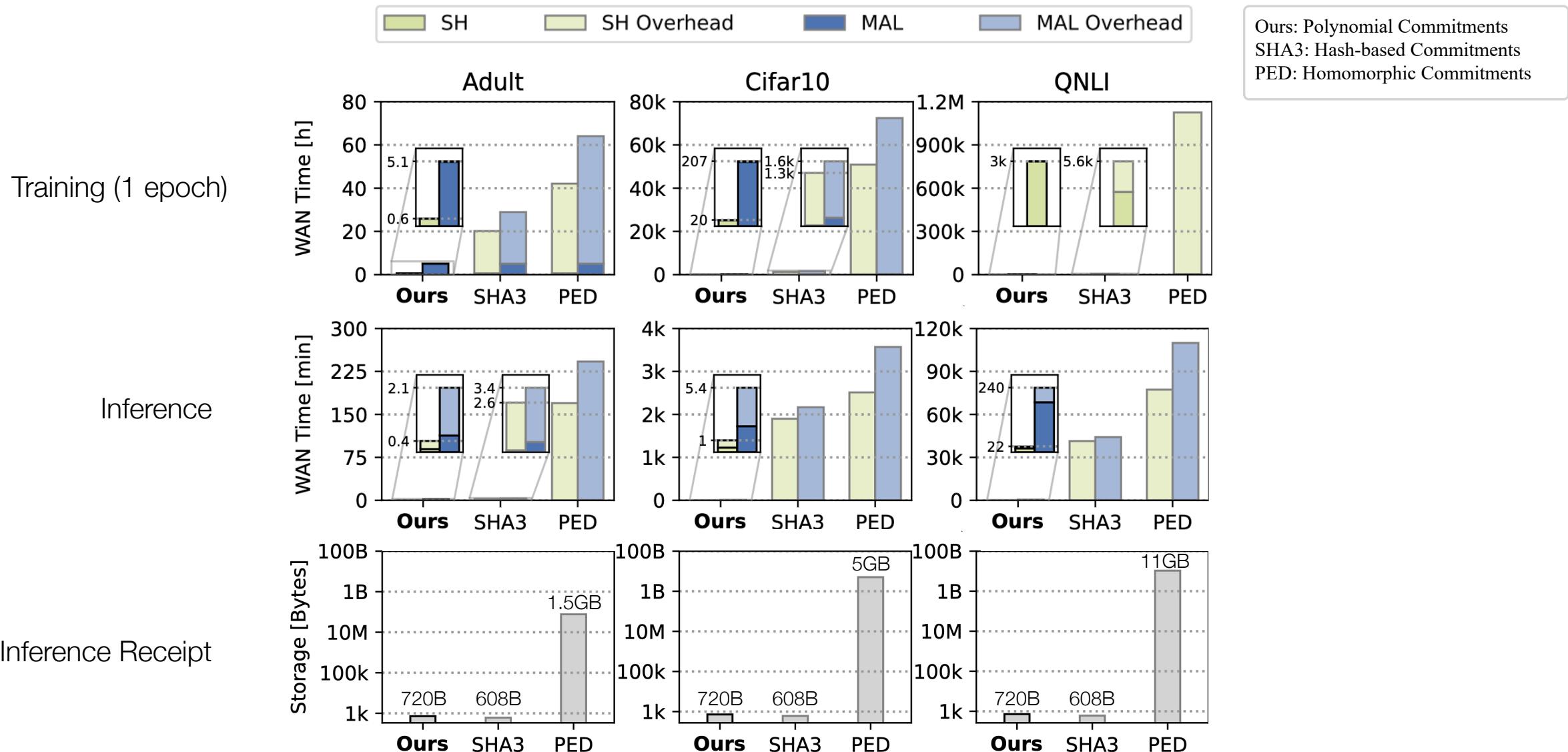


Evaluation

- End-to-end evaluation of Arc (integrated with MP-SPDZ), where POC is instantiated with:
 - Ours: Polynomial Commitments (KZG)
 - SHA3: Hash-based Commitments [Segal et al. @ AIES21]
 - PED: Pedersen Commitments [Cerebro, Zheng et al. @ USENIX Sec'21]
- Setting: 3PC semi-honest / malicious
- Auditing Functions: Indiv. Fairness, KernelSHAP, Robustness, kNNShapley

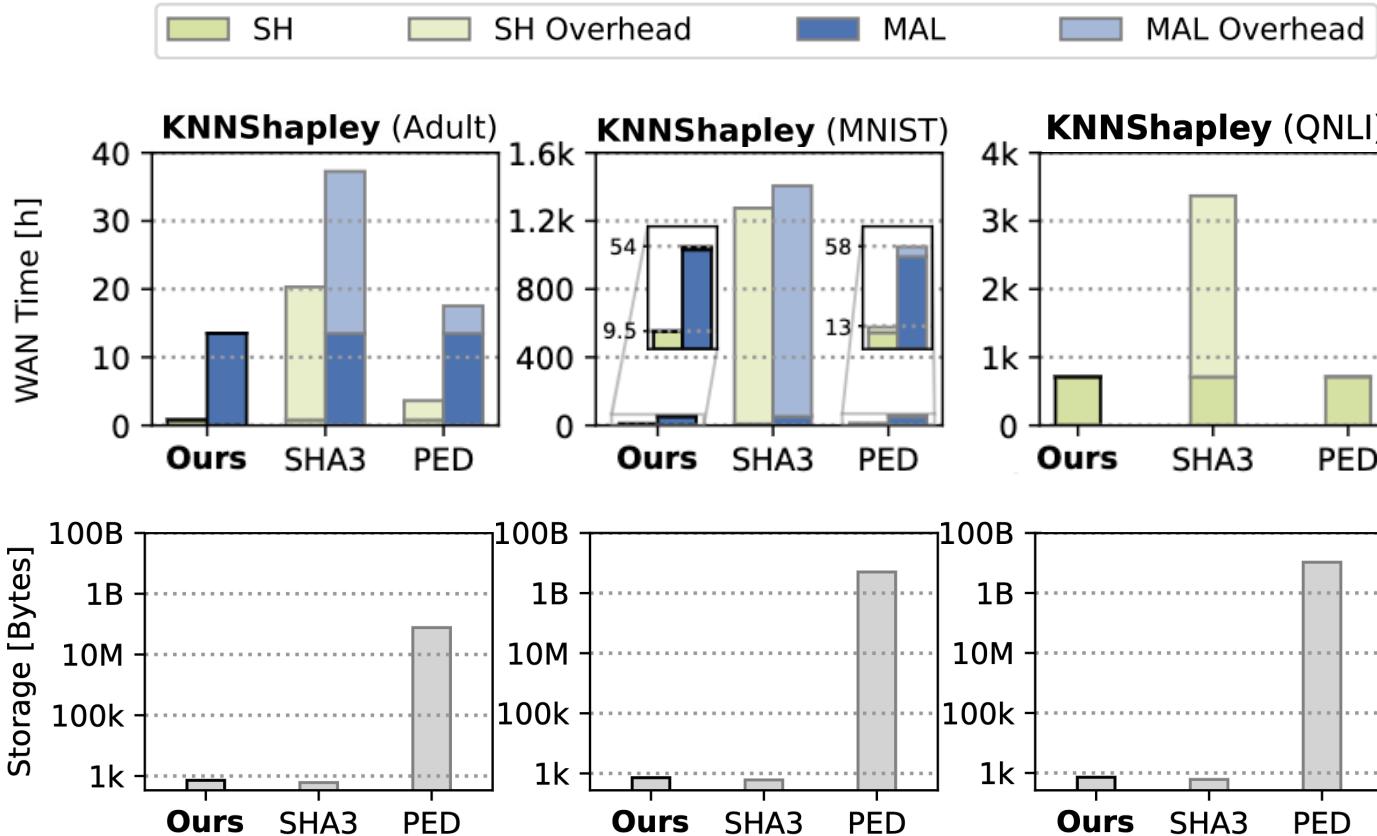
* The training overhead for malicious protocols and larger instances of the SHA3 and PED baselines are extrapolated

Evaluation: Training & Inference



Evaluation: Auditing

Auditing
Inference Receipt



Ours: Polynomial Commitments
SHA3: Hash-based Commitments
PED: Homomorphic Commitments

This Talk:

- Framework for Auditing PPML Pipelines
- New Protocol for Input Consistency in MPC

In the Paper:

- Protocol Extensions for Plaintext Training / Inference
- Arithmetic-to-Arithmetic Share Conversion
- MPC Arithmetizations of Auditing Functions

Future work:

- Robust Auditing Protocols
- Caching Auditing Inputs
- Verifiable Data Removal



pps-lab/arc



pps-lab.com/research/ml-sec

