

Paper-2

Enabling Execution Assurance of Federated Learning on Untrusted Participants.

Keywords:

Proof of work, Trusted Execution Environment (TEE),

Issues and Defenses:


1. Remove/weaken outliers from model update of participants, to prevent model poisoning attacks [7], [9]-[10]
2. ~~TEE~~ guarantee convergence of global model [8].
3. data poisoning attacks, stealthy backdoor attacks [12], [13].
4. prevent lazy participants.

TEE related issues.

1. only executing TEE requested training rounds.
2. ~~Defenses~~ freshness can't be guaranteed, as
3. ~~[7]-[9]~~ TEE's scheduling and I/O are controlled by host participant [20].

Issues Contd

5. old results corresponding correct proofs in the previous FL epoch to claim rewards.
6. sybil-based cheating attacks to ask for multiple rewards.

1. "Commit and prove" : Before exposing sampling decisions, the TEE requires to receive a commitment from the participant. Freshness:
2.  "Dynamic yet deterministic" data selection. Leverage dynamic input data to make training results different.
3. Bind selecting decision with each FL's epoch's identifier.
4. Tempering ~~the~~ free identifiers to participants ~~and~~ against sybil attacks.

Intel's Software Guard Extension (SGX)

1. generates a measurement with hash digest
2. program in enclave produces a report including the measurement and supplementary data (eg., output), and signs it as a quote with a private attestation key fused in the processor.
3. can be verified by trustworthy measurement of binary and Intel's Attestation Service (IAS)

→ attestation
 (SK_{TEE}, PK_{TEE}) ; $\Sigma \rightarrow$ unforgeable signature
 ↳ key

enclave generate attestation

$$\sigma = \Sigma \text{Sign}(\text{SKTEE}, \phi, \text{out})$$

output out is produced by program ϕ running in enclave.

$$\Sigma.\text{veri}(\text{PKTEE}, \sigma, \phi(\text{out})) \text{ outputs 1}$$

Threat Model

mitigating sybil attacks by diminishing gains of adversaries with increased computing efforts.

OVERVIEW

1. Solicitation

specifies training algo, rewards etc.

2. Registration

participants P and SGX-enabled platform can register with info like data description and identifier. Then, P reviews and downloads the program, launches the on enclave with the program, and generates an attestation report to demonstrate the correctness of program setup. $S.\text{verifies} \rightarrow \Sigma.\text{Veri}()$

3. Selection

S selects subset of participants [29], [33], for unselected S would instruct them when they can reconnect for participation.

4. Reporting: P gets the largest model Θ^t , performs local training, generates model update L^t and invokes inside program TEE to generate proof σ^t . TEE runs $(SK_{TEE}, DNN(1), D, \Theta^t, L^t) \rightarrow \sigma^t$ verifies using randomly selected training rounds.
5. Verification and Aggregation: S performs Verify $(PK_{TEE}, DNN(1), \Theta^t, L^t, \sigma^t) \rightarrow y_0$. S distributes rewards.

DESIGN DETAILS

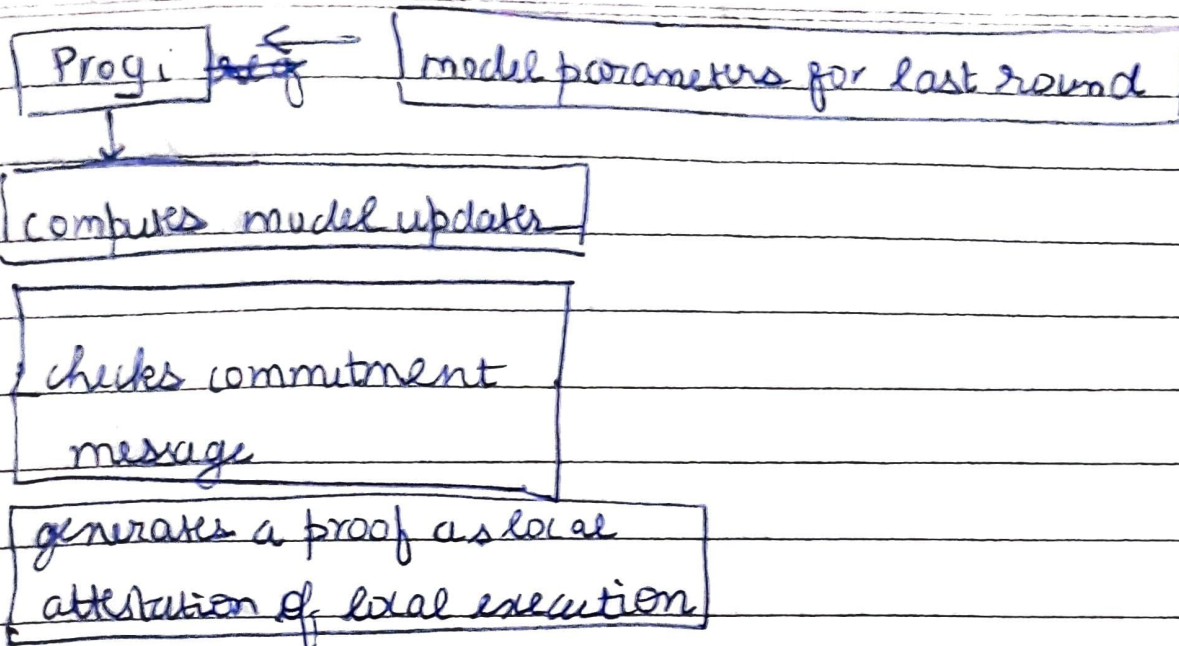
Progo - program outside enclave

Progi \rightarrow inside

A. Baseline Design

Progo stores the hash digest of all model parameters Θ s at the end of each training round and sends commit message to TEE when all training rounds finish.

Progi requests to check one random training round, Progo returns the input parameters, the hash value of output para, corresponding training data. Then progi checks whether they are consistent.



Design of Commitment Scheme

- 1) results of all training rounds and order of results
- 2) arbitrary result verified according to commitment efficiently.

Straightforward:

$$L = \{\underbrace{h_1, h_2, \dots, h_R}_{\text{hash from R rounds}}\}$$

Progi randomly fetches several pairs of two successive hash values, requests raw input model parameters from Progo, and recomputes hash values to insure coincidence.

For constant performance: we adopt Merkle hash tree [35] based commitment

MHT based: Progo builds Merkle Tree, leaf nodes are hash values. Internal nodes are concatenation of hash values of two children.

Progi requests ~~args~~ generates ~~nos~~ to identify the training rounds to be verified.

Input parameters, hash value of output para
auxiliary info (sibling nodes of nodes in
the path from the requested leaf node to
the tree root).

B. ~~How~~ Guaranteeing the Freshness of Results

Scheme is hard to guarantee the freshness of
single training round and final results.

$r \rightarrow$ rounds performed, $R \rightarrow$ required no. of rounds,

$L \rightarrow$ no. of randomly selected rounds

detected misbehaviour: $1 - \binom{R - (R_r) + 1}{L} / \binom{R}{L}$

Solution:

$$\binom{n}{r} = {}^nC_r$$

leverage dynamics of input training data, which
make epochs differentiable.

Progi negotiate with Progo about data organiza-
tion in advance.

Progi generates an index and HMAC for each
training data stored outside enclave
before training.

Progi $\xrightarrow{\text{seed}}$ Progo

$$I_{ij} = \text{PRF}(\text{seed} \times i + j) \bmod N$$

index of j th training data at round i

PRF \rightarrow pseudo random function.

$N \rightarrow$ total training data.

Checker checks integrity of outside data with HMAC, trains the DNN, verifies the output parameters
Bind training results with used global model
seed = Progi receive ~~light~~ digest of current global model (denoted as h^t) from Prog.

Using signature and model digest, the server can ensure ~~that~~ whether the participant used latest model and faithfully performed local training.

sybil-based, unique identifier T , tamper free identifier, even if an adversary registers multiple aliases, he must launch some no. of enclave instances to generate proofs.