

# SCALE-FL: Scalable Cryptography-based Aggregation with Lightweight Enclaves for Federated Learning

Micah Brody, Antonia Januszewicz, Jiachen Zhao, Nirajan Koirala, Taeho Jung  
University of Notre Dame

{mbrody, ajanusze, jzhao, nkoirala, tjung}@nd.edu

## Abstract

Privacy-Preserving Federated Learning (PPFL) emphasizes the security and privacy of contributors' data in scenarios such as healthcare, smart grids, and the Internet of Things. However, ensuring the security and privacy throughout PPFL can be challenging, given the complexities of maintaining relationships with many users across multiple epochs. Additionally, under a threat model in which the aggregating server and corrupted users are colluding adversaries, honest users' inputs and output data must be protected at all stages. Two common tools for enforcing privacy in federated learning are Private Stream Aggregation (PSA) and Trusted Execution Environments (TEE). However, PSA-only approaches still expose the raw aggregate to the server (and thus to colluding parties). TEE-only aggregation typically incurs non-negligible per-client per-epoch overhead at scale because the TEE must handle per-client communication and maintain per-client state/key material. This paper presents SCALE-FL, a novel solution for PPFL that maintains security while achieving near-plaintext performance using a state-of-the-art PSA protocol to collect user information and a TEE to hide information about the raw aggregate. By using a PSA protocol for aggregation, we can maintain the privacy of information on the untrusted server without requiring per-user key storage or use by the TEE. Then, the aggregate is securely processed by the TEE in plaintext, without the heavy encryption required on an untrusted server. Finally, we ensure the security of user inputs in the federated learning output by using Differential Privacy (DP). The additional overhead introduced by SCALE-FL is 1% of the overhead of the plain FL executions.

## Keywords

Federated Learning, Secure Aggregation, SGX, PPML, PPFL

## 1 Introduction

Privacy-preserving machine learning (PPML) has attracted sustained attention since the early development of machine learning. Modern machine learning systems often require aggregating large training datasets that contain sensitive personal, medical, or financial information [26]. Therefore, it is crucial to develop methods that preserve privacy while maintaining model accuracy and convergence behavior comparable to non-private training.

Federated learning (FL) offered a breakthrough that allowed users and model aggregators to cooperatively train a model without sharing their raw training data [31]. Consequently, FL has since been deployed at scale in settings such as mobile and other resource-constrained devices, where keeping data local and minimizing communication overhead are essential [3]. FL is also increasingly relevant in cross-silo settings, such as institutions that wish to train models jointly without centralizing sensitive records.

While each contributor's data remains local in FL, model updates themselves can reveal sensitive information, necessitating additional privacy mechanisms [47, 48]. Moreover, the client-server interaction in FL enables inference and reconstruction attacks in which a malicious server, or colluding and compromised clients, can extract information about an individual participant's data [2]. The study of privacy-preserving federated learning (PPFL) techniques aims to enable FL without compromising the privacy of contributors. However, privacy mechanisms often struggle to preserve the high-throughput, low-latency aggregation pipeline, since even small per-epoch overheads can become prohibitive at scale [46].

In PPFL, different mechanisms trade off efficiency/scalability and security. In particular, Trusted Execution Environments (TEEs), such as Intel SGX, are among the most efficient practical approaches for achieving strong confidentiality of aggregates while preserving near-plaintext computation, since they enable sensitive processing inside a hardware-isolated environment with memory encryption. Although TEE-only designs can keep aggregates confidential, they often rely on per-client secure channels and key management. Furthermore, enclave I/Os and system calls from networking and disk I/Os result in expensive context switches. At scale, these all result in non-negligible overhead that becomes problematic as the number of clients and training epochs increases [4, 21, 32]. This creates a tension between minimizing trusted online work and providing strong confidentiality guarantees for the aggregate.

*Secure aggregation* and *private stream aggregation* (PSA) are techniques that hide individual values while computing their sum. Thus, secure aggregation and PSA protocols can hide individual updates during collection. However, if the plaintext aggregate is revealed to the server, then in high-collusion settings the server can combine the aggregate with colluding clients' contributions to isolate and exactly recover a target user's update [1, 19, 29]. This risk is particularly salient in large-scale deployments, where participants cannot assume that no collusion is taking place. Therefore, we consider a scalable PPFL scheme in a strong collusion setting (a malicious server plus all but one corrupted clients), while avoiding common assumptions such as two non-colluding servers.

Despite these protections, PPFL must also address what is revealed during training. Even if individual updates are hidden in transit and during aggregation, the resulting aggregate (or the updated model parameters), which are released to the end users after each epoch, can leak information about a single participant, particularly if collusions between users and servers are allowed [11, 12, 45]. Differential privacy (DP) can mitigate this leakage, but enforcing DP in the stronger threat model while keeping the additional overhead negligible is nontrivial. Global (central) DP is meaningless if a malicious server can change the noise, and even small deviations can invalidate the claimed privacy parameters. At the same time,

local DP approaches shift the burden to clients and substantially degrade model utility due to the larger noise required [28, 33, 44]. We can enforce global DP through the usage of TEEs, specifically Intel SGX, but this would introduce the aforementioned non-negligible additional overhead.

We seek a PPFL design that keeps the aggregate confidential and enables enforceable global DP, while ensuring the additional overhead introduced by the secure and private computations is negligible compared to the plain FL (e.g.,  $<1\%$  additional overhead). No state-of-the-art approaches can achieve this level of efficiency while enforcing both confidentiality and global DP in the presence of collusive malicious adversaries. To fill this gap and push the limit of the current SGX-based approaches, we propose SCALE-FL, a hybrid PPFL architecture that achieves near-plaintext aggregation efficiency while providing strong privacy even under a collusive adversary model (Figure 1). The core idea is a heterogeneous secure computation architecture that partitions the aggregation pipeline to minimize enclave I/O and per-client trusted work, rather than treating PSA and SGX as black-box components. Clients encrypt their updates using a state-of-the-art PSA scheme, enabling high-throughput aggregation at near plain-computation speed while keeping updates private in transit and during aggregation [41]. An untrusted server performs only public, high-throughput aggregation over ciphertexts. A TEE then performs a small secret-dependent finalization step, namely completing decryption to recover the epoch aggregate. In our DP-enhanced variant, the enclave also enforces clipping and noise addition, and it releases only a DP-sanitized aggregate, so that the raw aggregate never leaves the enclave. This yields end-to-end privacy even when the server colludes with up to  $n-1$  clients. Concretely, SCALE-FL assigns throughput-critical work to the untrusted server and reserves the enclave for a minimal finalization role. The server receives  $n$  client updates and aggregates the PSA ciphertexts, thereby avoiding per-client enclave interactions. The enclave processes only the single aggregated ciphertext per epoch to finalize recovery of the aggregate, and it can also enforce clipping and noise addition.

Unlike prior TEE-centric PPFL designs that require per-client secure channels each epoch with significant overhead due to context switches from enclave/network I/Os [4, 21], SCALE-FL reduces enclave involvement to the minimum. Notably, we avoid per-client online sessions and per-client decryption within the enclave, ensuring the in-enclave overhead remains constant with respect to the number of clients (but grows linearly with the size of the model). The adversarial server collects  $n$  encrypted client messages per epoch and computes their coordinate-wise sum, while the enclave processes just one single aggregated ciphertext vector of length  $d$  per epoch to complete aggregate recovery and DP sanitization, where  $d$  is the number of model parameters in the update vector. We first present a base protocol that provides aggregate confidentiality, and then a DP-enhanced variant that enforces global DP within the enclave and provides formal end-to-end privacy guarantees under strong collusion.

In summary, our goal is a single-server PPFL system that scales to large client populations, keeps the aggregate confidential under strong collusive adversarial model, can enforce central DP without relying on client-side noise, and minimize the I/Os of enclaves to keep them constant. The contributions of this work are:

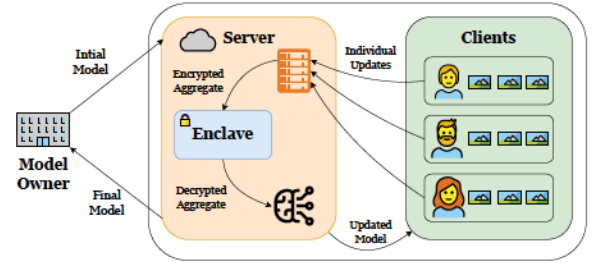


Figure 1: Overview of SCALE-FL and its trust split. Clients send TERSE-encrypted updates to an untrusted server, which performs public aggregation over ciphertexts. The enclave finalizes decryption to recover the epoch aggregate; in the base protocol  $\Pi$  it may release the raw aggregate, while in  $\Pi_{DP}$  it applies central DP and releases only a sanitized aggregate.

- **Near-plaintext PPFL via heterogeneous secure computation.** We introduce a hybrid PPFL architecture that partitions computation between cryptographic aggregation and TEE-based confidential computing. The untrusted server performs high-throughput aggregation over PSA ciphertexts, and the enclave performs only a small finalization step on the single aggregated ciphertext per epoch, plus DP enforcement in our DP-enhanced variant. This design avoids per-client enclave sessions and minimizes enclave I/O, enabling near-plaintext aggregation efficiency with  $<1\%$  overhead.
- **A scalable PPFL architecture for strong collusion without extra trust assumptions.** We introduce SCALE-FL, which combines TERSE-based PSA with a minimal TEE finalization step to keep client updates private in transit and during aggregation, even when the server colludes with up to  $n-1$  clients, while avoiding assumptions such as two non-colluding servers.
- **Enforceable global DP via a minimal trusted component.** We provide a DP-enhanced variant in which clipping and noise addition are enforced inside the enclave and only a sanitized aggregate is released, enabling end-to-end privacy guarantees in a malicious-server setting and avoiding utility costs of local DP perturbation.
- **Practical extensions and analysis of recent attacks.** We describe modifications to support client sampling, dropout handling, and weighted averaging, and we explain why recent IND-CPA<sup>D</sup> style attacks that rely on observing decryption outputs do not apply to our deployment model and parameterization.

## 2 Related Work

Existing PPFL systems often rely on either multi-server architectures or design assumptions that limit practical scalability. Moreover, our threat model extends beyond protecting the confidentiality of individual user inputs. We also consider the privacy risks arising from publicly released aggregates and model updates, which still leak sensitive information about user data [2, 23, 43].

Protocols in the two-server setting, such as AlphaFL [19] and ELSA [37], can improve efficiency by splitting trust and computation across two non-colluding servers, thereby reducing the amount of expensive secure computation required at any single party. However, these gains come at the cost of increased communication complexity. Each client must communicate with both servers, and

**Table 1: Comparison of SCALE-FL with state-of-the-art PPFL protocols in terms of trust assumptions (single server), aggregate privacy, scalability with the number of users, and overhead relative to plaintext.**

Work	Single Server	Aggregate Privacy	User scaling	Speed (vs. plaintext)
SCALE-FL	✓	✓	✓	✓
Anofel [1]	✓	✓	✗	✗
Olive [21]	✓	✓	✗	✗
SRFL [4]	✓	✗	✓	✗
AlphaFL [19]	✗	✓	✗	✗
ELSA [37]	✗	✓	✓	✗

the servers must additionally coordinate with one another during protocol execution. This dual communication requirement substantially increases overall communication overhead, which can become prohibitive as the number of participating clients grows.

Single-server solutions attempt to reduce trust assumptions but introduce different scalability challenges. For example, Olive [21] relies heavily on TEEs to provide security guarantees. While this design simplifies the trust model, it introduces significant overhead associated with TEE-client communication, remote attestation, and repeated key exchanges. These bottlenecks limit scalability as the client population increases. Our approach mitigates these constraints by combining TEE-based protection with PSA, thereby reducing per-client TEE interaction while preserving strong privacy guarantees. Similarly, SRFL [4] adopts a TEE-based approach but requires each client to operate its own TEE, effectively shifting the hardware and deployment burden to end users. This assumption may be impractical in many real-world FL deployments, particularly in heterogeneous or resource-constrained environments.

We summarize how these systems compare along key dimensions such as trust assumptions, aggregate confidentiality, scalability with the number of users, and performance relative to plaintext in Table 1. Here, *User scaling* indicates whether the protocol’s per-epoch computation and communication grow mildly with the number of clients (e.g., no per-client TEE sessions/attestation, no heavy cryptographic or communication costs scaling exponentially with  $n$ ), and *Speed* indicates whether the reported runtime overhead is close to plaintext execution, particularly in the aggregation step (i.e., near-plaintext aggregation throughput).

### 3 Preliminaries & Definitions

#### 3.1 Federated Learning

Federated learning (FL) is a distributed training framework in which multiple clients collaboratively train a shared machine learning model while keeping their data local. Rather than centralizing datasets, each client trains its own model on its own device and shares only model updates with a coordinating server. In the FedAvg algorithm, clients execute several local optimization steps and transmit their updated model parameters to the server, which aggregates them to produce a new global model that is redistributed for the next training epoch [40]. A related approach, FedSGD, instead communicates locally computed gradients at every epoch

using distributed stochastic gradient descent [31]. In both variants, raw training data never leave client devices, distinguishing FL from traditional centralized machine learning pipelines. Although this decentralized structure reduces exposure of sensitive data, it does not, by itself, provide formal privacy guarantees, as model updates may still leak information about local datasets.

We provide a brief summary of the basic notations used below:

- **FL.Update**( $W_\rho, \hat{g}_{\text{agg}}$ )  $\rightarrow W_{\rho+1}$ : Given the model params  $W_\rho$  and user data aggregate  $\hat{g}_{\text{agg}}$ , generate updated model params  $W_{\rho+1}$ .
- **FL.LocalTrain**( $W_\rho; D_i$ )  $\rightarrow u_{i,\rho}$ : Given the set of model parameters  $W_\rho$  and local data  $D_i$ , generate the user update  $u_{i,\rho}$ .

#### 3.2 Private Stream Aggregation

To collect client updates privately, we instantiate our Private Stream Aggregation (PSA) layer with TERSE [41]. PSA was originally introduced for privacy-preserving time-series analytics [39]. At each time step, each user encrypts its value and sends the ciphertext to an untrusted aggregator. The aggregator can compute only the sum over all users’ values for that time step, and it learns nothing about any individual input beyond what is implied by the aggregate [39]. TERSE is a highly optimized, lattice-based PSA construction built from *ring learning with errors (RLWE)*-based cryptography and an additively homomorphic encryption structure that enables efficient ciphertext-domain summation [41]. After a setup phase that distributes per-user encryption keys and an aggregation key, users encrypt their inputs and the server aggregates ciphertexts using homomorphic addition. Only the party holding the aggregation key can decrypt the final ciphertext to recover the plaintext aggregate.

*TERSE interface and our adaptations.* We summarize the TERSE routines we use, along with our fixed-point wrapper, below:

- **TERSE.Setup**( $\lambda, t, n$ )  $\rightarrow (s_0, \dots, s_{n-1}, s')$ : on input security parameter  $\lambda$ , plaintext space  $t$ , and number of users  $n$ , outputs per-user secret keys  $s_0, \dots, s_{n-1}$  and an aggregation key  $s'$ .
- **TERSE.Encode**( $S, u_{i,ts}$ )  $\rightarrow x_{i,ts}$ : fixed-point encodes a real-valued update  $u_{i,ts} \in \mathbb{R}^d$  into an integer vector  $x_{i,ts} \in \mathbb{Z}^d$  using a public scale factor  $S$ .
- **TERSE.Encrypt**( $s_i, ts, x_{i,ts}$ )  $\rightarrow ct_{i,ts}$ : encrypts the encoded vector under user key  $s_i$  to produce a ciphertext  $ct_{i,ts}$ .
- **TERSE.Add**( $ct_{0,ts}, \dots, ct_{n-1,ts}$ )  $\rightarrow y_{ts}$ : computes the ciphertext-domain sum  $y_{ts}$  of all users’ ciphertexts at timestamp  $ts$ .
- **TERSE.Decrypt**( $s', ts, y_{ts}$ )  $\rightarrow \text{agg}_{ts}^{\mathbb{Z}}$ : decrypts the aggregated ciphertext  $y_{ts}$  using aggregation key material  $s'$  to recover the plaintext integer aggregate  $\text{agg}_{ts}^{\mathbb{Z}}$ . In SCALE-FL, this step is executed inside the enclave.
- **TERSE.Decode**( $S, \text{agg}_{ts}^{\mathbb{Z}}$ )  $\rightarrow \text{agg}_{ts}^{\mathbb{R}}$ : fixed-point decodes the integer aggregate into a real-valued aggregate  $\text{agg}_{ts}^{\mathbb{R}} \in \mathbb{R}^d$ .

*Fixed-point encoding and decoding.* Since TERSE encrypts over an integer plaintext space, we encode each update  $u_{i,ts} \in \mathbb{R}^d$  using a scale factor  $S$ :

$$x_{i,ts} := \text{Encode}(S, u_{i,ts}) = \lfloor [S \cdot u_{i,ts}] \rfloor_t \in \mathbb{Z}_t^d,$$

where  $\lfloor \cdot \rfloor_t$  denotes coordinate-wise reduction modulo  $t$ . Fixed-point encoding introduces a bounded rounding error. Prior work shows that fixed-point representations can preserve model utility at common precisions and incur negligible accuracy impact [18, 30].

We also provide a matching Decode function to reverse the process after aggregation. Concretely, let  $\text{center}_t : \mathbb{Z}_t \rightarrow (-\frac{t}{2}, \frac{t}{2}]$  denote the centered representative map. The enclave sets

$$\text{agg}_{ts}^R \leftarrow \text{Decode}(S, \text{agg}_{ts}^Z) := \text{center}_t(\text{agg}_{ts}^Z)/S,$$

where  $\text{center}_t(\cdot)$  and division by  $S$  are applied coordinate-wise. Under our parameter choices for honest executions, the (unreduced) integer aggregate lies in the centered range  $(-\frac{t}{2}, \frac{t}{2}]^d$ , i.e., no coordinate wraps modulo  $t$ . Therefore, applying  $\text{center}_t(\cdot)$  and dividing by  $S$  recovers the intended real-valued aggregate up to the fixed-point rounding error introduced during encoding.

### 3.3 Trusted Execution Environment

A Trusted Execution Environment (TEE) provides hardware-backed isolation for code and data even when the host software stack is untrusted [38]. We use *Intel Software Guard Extensions (SGX)*, which supports isolated execution in enclaves with confidentiality and integrity guarantees enforced by the processor. We model these guarantees using the  $\mathcal{G}_{\text{att}}$ -hybrid model in the global universal composability framework [36].

SGX also supports remote attestation, which lets a client verify that it is communicating with a genuine enclave running a specific measured program and (optionally) establish an authenticated confidential channel for key exchange [7]. We denote actions occurring inside the TEE by **highlighting the text in blue**.

### 3.4 Differential Privacy

Differential Privacy (DP) is a formal notion of privacy that limits what can be inferred about any single participant from the output of a computation [9]. Informally, DP ensures that an adversary cannot reliably determine whether a particular user’s data was included in the dataset solely from the released output. It achieves this by requiring that the distribution of outputs changes only slightly when one user’s data is removed or modified.

**DEFINITION 1 (DIFFERENTIAL PRIVACY [9]).** A randomized mechanism  $\mathcal{M}$  with range  $\mathcal{R}$  satisfies  $(\epsilon, \delta)$ -differential privacy if for all adjacent datasets  $D$  and  $D'$  and for all measurable subsets  $O \subseteq \mathcal{R}$ ,

$$\Pr[\mathcal{M}(D) \in O] \leq e^\epsilon \Pr[\mathcal{M}(D') \in O] + \delta.$$

*Global (central) DP vs. local DP (where randomness is added).* DP can be enforced in two ways, depending on where randomness is introduced [9].

- *Global (central) DP:* a trusted curator first computes the exact aggregate and then adds calibrated noise before releasing the result.
- *Local DP:* each user randomizes its updates before sending anything. This removes the need for a trusted curator, but may incur a larger noise penalty because noise is added before aggregation [28].

*Terminology (global vs. local; user-level vs. record-level).* The terms *global vs. local DP* specify where randomization occurs (trusted curator vs. each user), whereas *user-level vs. record-level* specify the *adjacency relation*. Our design targets global (central) DP with *user-level adjacency*, which captures the goal that the output should not change much when a single user’s entire local dataset (or participation) is changed or removed.

**DEFINITION 2 (USER-LEVEL ADJACENCY [13]).** Let  $D = (D_1, \dots, D_n)$  and  $D' = (D'_1, \dots, D'_n)$  denote federated datasets, where  $D_i$  is the local dataset of user  $U_i$ . We say  $D$  and  $D'$  are user-adjacent if they differ in exactly one user’s dataset, meaning there exists  $i^*$  such that  $D_{i^*} \neq D'_{i^*}$  and  $D_j = D'_j$  for all  $j \neq i^*$ .

*DP Mechanism Interface.* Our protocol is parameterized by a global DP mechanism  $\mathcal{M}$ . We view  $\mathcal{M}$  as exposing two operations. First, users locally bound their updates to control sensitivity. Then, after the enclave decrypts and decodes the aggregate, it applies DP and releases only a sanitized aggregate.

- **M.Preprocess**( $u, C$ )  $\rightarrow \tilde{u}$ : bounds a user’s real-valued update  $u$  and returns preprocessed update  $\tilde{u}$ .

- **M.ApplyDP**( $\text{st}_{\mathcal{M}}, x_{\text{agg}}$ )  $\rightarrow (\tilde{x}_{\text{agg}}, \text{st}'_{\mathcal{M}})$ : adds calibrated noise to the plaintext aggregate and advances the mechanism’s internal state and returns  $(\tilde{x}_{\text{agg}}, \text{st}'_{\mathcal{M}})$ .

## 4 Definitions and Assumptions

### 4.1 System Model

Our FL system involves three parties: *users*, a *server*, and an *SGX enclave*. We consider synchronous training epochs in which all honest users attempt to contribute one update per epoch. Extensions to client sampling (where only a subset of users contribute each epoch) and dropout handling are discussed in Section 8.1 and Section 8.2, respectively.

*Users* hold private training data and participate in federated learning. In each epoch, user  $i$  computes a local update (e.g., a gradient) on the current global model, encrypts it using TERSE, and sends the resulting ciphertext to the server. Users do not communicate directly with each other.

The *server* coordinates training and aggregates encrypted user updates. In each epoch  $\rho$ , it receives ciphertexts from users, performs a public aggregation step, and forwards the aggregate ciphertext to the SGX enclave. The server does not possess the key material needed to recover the plaintext aggregate. After receiving the enclave’s output, the server applies the model update and broadcasts the updated model to users.

The SGX enclave is the trust boundary for decryption. On initialization, it generates TERSE parameters and keys, provisions each user’s secret key over an attested secure channel, and retains the aggregation key internally. Each epoch, it decrypts the server’s aggregate ciphertext and releases either the plaintext aggregate (base protocol) or a DP-sanitized aggregate (DP-enhanced protocol).

### 4.2 Assumptions

*TERSE cryptographic assumptions.* We rely on the security of TERSE, which is based on standard lattice assumptions (e.g., RLWE, per the TERSE construction and parameterization) [41]. We assume the public parameters are chosen so that TERSE decryption fails with at most negligible probability for honest executions.

*TEE (SGX) assumption.* We analyze our protocol in the  $\mathcal{G}_{\text{att}}$ -hybrid model [36], which idealizes an SGX-like attested execution environment. We assume confidentiality of the enclave state and integrity of the enclave execution and outputs. We also assume remote attestation that binds outputs and key-provisioning messages to a specific enclave program identity. Attestation lets each

user authenticate the intended enclave and derive a confidential, authenticated channel to it, for example, using an attested public key. Accordingly, we model user–enclave communication, including key provisioning, as occurring over secure channels that are not susceptible to attacker-in-the-middle attacks.

*Side-channel attacks.* We do not address availability since a malicious host can deny service, and we exclude microarchitectural side channels. Certain side-channel attacks can compromise parts of Intel SGX’s security [35]. Such attacks include: cache timing [14], speculative execution leakage [5], and hardware interposer attacks [8]. However, such attacks require code-specific modifications or physical defenses to prevent. Thus, these are outside the  $\mathcal{G}_{\text{att}}$ -hybrid model [36], and we consider such attacks orthogonal to our work.

### 4.3 Adversary and Threat Model

*Adversary Model* We consider a PPT adversary that corrupts and controls the server and a subset of users who are actively malicious. The adversary observes all communication strings and can deviate arbitrarily from the prescribed protocol. In particular, in each epoch of FL, it may omit, replay, or duplicate user ciphertext, as well as make any arbitrary changes to the final ciphertext  $\hat{y}_r$  sent to the enclave. We further assume the adversary/server is rational (or rational-covert); it wants to obtain a high-utility final model and is therefore disincentivized from manipulations that would materially degrade the performance of the final model.

The adversary statically corrupts a subset of users  $C \subseteq [n]$  with  $|C| \leq n - 1$  before the protocol execution begins. For each corrupted user  $U_j$  (for  $j \in C$ ), the adversary learns the user’s long-term secret key material (including the TERSE user key  $s_j$ ) and the user’s local state (e.g., its dataset  $D_j$  and any locally stored training state). During the execution, the adversary fully controls corrupted users’ actions and outgoing messages, and may choose their plaintext updates arbitrarily. Honest users  $\mathcal{H} = [n] \setminus C$  follow the protocol specification and do not reveal their secret keys or plaintext updates to the adversary. The enclave is not corrupted: in the  $\mathcal{G}_{\text{att}}$ -hybrid model, its internal state remains confidential and its execution/output integrity is protected, and remote attestation binds its behavior to the intended program identity [36].

*Security Goal: Input Confidentiality and Privacy* In the case of a malicious server that delivers manipulated ciphertexts to the enclave, we do *not* provide a cryptographic guarantee of model correctness or integrity. If training proceeds under such deviations, the model update computed from the enclave’s output may be adversarially biased. Such manipulation is often detectable after the fact through standard model evaluations and may trigger financial or reputational penalties. Our formal security goal is *privacy*: even under arbitrary active manipulation by the server, the adversary should learn nothing about honest users’ updates beyond what is implied by (i) the aggregate information explicitly released by the enclave and (ii) the internal state and inputs of corrupted users.

### 4.4 Protocol Definitions

*Protocol Syntax.* Let  $\Pi = (\text{Setup}, \text{Enc}, \text{Add}, \text{Dec}, \text{Update})$  be a protocol between users  $U_1, \dots, U_n$ , a server, and an SGX enclave:

**Setup**( $1^\lambda, n, R$ ): outputs public parameters  $pp$ , user secret keys  $\{k_i\}_{i \in [n]}$ , and enclave secret state  $\text{st}_{\text{enc}}$ .

**Parameters:** Let  $U_1, \dots, U_n$  denote the users and  $S$  the server. Let  $R$  be the number of epochs and  $d$  the update dimension. Let  $\mathcal{W}$  be the model space and  $W_0 \in \mathcal{W}$  the initial model. Let  $\text{Update} : \mathcal{W} \times \mathbb{Z}_t^d \rightarrow \mathcal{W}$  be the model-update rule.

**State:** A model  $W$  (initialized to  $W_0$ ); and an integer counter  $\rho \in \{0, \dots, R\}$  (initialized to 0).

**Inputs:**

- $U_i$ : updates  $(x_{i,0}, \dots, x_{i,R-1})$  with each  $x_{i,\rho} \in \mathbb{Z}_t^d$  (for each  $i \in [n]$ ).
- $S$ : for each epoch  $\rho$ , an offset  $\delta_\rho \in \mathbb{Z}_t^d$ .

**Outputs:** For each epoch  $\rho \in \{0, \dots, R-1\}$ :

- Output to  $S$ :

$$g_\rho := \left( \sum_{i \in [n]} x_{i,\rho} \right) + \delta_\rho \in \mathbb{Z}_t^d.$$

- Update and output:  $W \leftarrow \text{Update}(W, g_\rho)$ .

Figure 2: Ideal functionality  $\mathcal{F}_{\text{agg}}$ .

**Enc**( $k_i, ts, x$ ): on input a user key  $k_i$ , a timestamp identifier  $ts$ , and a plaintext value  $x$ , outputs a ciphertext  $c$ .

**Add**( $c_1, \dots, c_n$ ): on input a collection of ciphertexts for the same timestamp, outputs an aggregated ciphertext  $\hat{y}$ .

**Dec**( $\text{st}_{\text{enc}}, ts, \hat{y}$ ): on input enclave secret state  $\text{st}_{\text{enc}}$ , a timestamp identifier  $ts$ , and an aggregated ciphertext  $\hat{y}$ , outputs a plaintext aggregate value  $x_{\text{agg}}$ .

**Update**( $W_\rho, x_{\text{agg}}$ ): on input current model weights  $W_\rho$  and an aggregate update  $x_{\text{agg}}$ , outputs updated model weights  $W_{\rho+1}$ .

*Ideal Functionalities.* We define security via ideal functionalities for the base protocol  $\mathcal{F}_{\text{agg}}$  (Figure 2) and for the DP-enhanced protocol  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  (Figure 3). Both are multi-epoch and stateful, modeling the full FL process over  $R$  epochs. Both functionalities hide honest users’ individual updates from the adversary while allowing the adversary to arbitrarily bias the aggregate value delivered to the server in each epoch. This reflects our privacy-without-integrity guarantee: we protect the confidentiality of individual updates, but we do not prevent an actively malicious server from causing the training procedure to use a distorted aggregate.

$\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  extends  $\mathcal{F}_{\text{agg}}$  in one respect: the output to the adversary is mediated by the DP mechanism  $\mathcal{M}$ . Instead of revealing the delivered aggregate  $y_\rho$  directly,  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  applies  $\mathcal{M}$ . ApplyDP and reveals only the sanitized value  $\tilde{y}_\rho$ . The model update is then computed from  $\tilde{y}_\rho$  rather than  $y_\rho$ , so the DP noise propagates into the trained model. In both functionalities, honest users’ individual inputs are never revealed to the adversary. The only information released is the (possibly sanitized) aggregate and the model derived from it.

### 4.5 Security Definitions

We formalize two levels of guarantees for FL with secure aggregation. Definition 3 defines *aggregation security*: the protocol correctly hides individual users’ updates during the aggregation process. The base protocol  $\Pi$  satisfies this property (see Section 6.2).

**DEFINITION 3 (AGGREGATION SECURITY).** Protocol  $\Pi$  securely realizes  $\mathcal{F}_{\text{agg}}$  in the  $\mathcal{G}_{\text{att}}$ -hybrid model if for every PPT adversary  $\mathcal{A}$

**Parameters:** Let  $U_1, \dots, U_n$  denote the users and  $S$  the server. Let  $R$  be the number of epochs and  $d$  the update dimension. Let  $\mathcal{W}$  be the model space and  $W_0 \in \mathcal{W}$  the initial model. Let  $\text{Update} : \mathcal{W} \times \mathbb{R}^d \rightarrow \mathcal{W}$  be the model-update rule. Let  $\mathcal{M}$  be a stateful DP mechanism with an operation  $M$ .  $\text{ApplyDP}$  that takes inputs in  $\mathbb{Z}_t^d$  (internally applying a fixed public deterministic decoding to  $\mathbb{R}^d$ ).

**State:** A model  $W$  (initialized to  $W_0$ ); an integer counter  $\rho \in \{0, \dots, R\}$  (initialized to 0); and DP state  $st_{\mathcal{M}}$  (initialized appropriately).

**Inputs:**

- $U_i$ : updates  $(x_{i,0}, \dots, x_{i,R-1})$  with each  $x_{i,\rho} \in \mathbb{Z}_t^d$  (for each  $i \in [n]$ ).
- $S$ : for each epoch  $\rho$ , an offset  $\delta_\rho \in \mathbb{Z}_t^d$ .

**Outputs:** For each epoch  $\rho \in \{0, \dots, R-1\}$ :

- Output to  $S$ :

$$(\tilde{g}_\rho, st_{\mathcal{M}}) := M.\text{ApplyDP}\left(st_{\mathcal{M}}, \left(\sum_{i \in [n]} x_{i,\rho}\right) + \delta_\rho\right) \text{ with } \tilde{g}_\rho \in \mathbb{R}^d.$$

- Update and output to all parties:  $W \leftarrow \text{Update}(W, \tilde{g}_\rho)$ .

Figure 3: Ideal functionality  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$ .

controlling the server and up to  $n-1$  corrupted users, there exists a PPT simulator  $\mathcal{S}$  such that for all PPT environments  $\mathcal{Z}$ ,

$$\text{REAL}_{\Pi, \mathcal{A}, \mathcal{Z}}^{\mathcal{G}_{\text{att}}} (1^\lambda) \approx_c \text{IDEAL}_{\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}, \mathcal{S}, \mathcal{Z}} (1^\lambda).$$

For formal definitions of the execution experiments  $\text{REAL}_{\Pi, \mathcal{A}, \mathcal{Z}}^{\mathcal{G}_{\text{att}}}$  and  $\text{IDEAL}_{\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}, \mathcal{S}, \mathcal{Z}}$ , refer to Appendix B.

Informally, Definition 3 guarantees that the base protocol leaks nothing about honest users' individual updates beyond the aggregate released by the enclave. However, as discussed in Section 5.4, aggregation security alone does not prevent inference from the released aggregate or the trained model. We therefore define *FL privacy* (Definition 4), which additionally requires that the released transcript satisfies differential privacy, bounding the information any observer can extract about any individual user's data. The DP-enhanced protocol  $\Pi_{\text{DP}}$  satisfies this stronger property.

**DEFINITION 4 (FL PRIVACY).** Let  $\mathcal{M}$  be a DP mechanism satisfying  $(\epsilon, \delta)$ -differential privacy with user-level adjacency (Definition 2). Protocol  $\Pi_{\text{DP}}$  securely realizes  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  with  $(\epsilon, \delta)$ -DP in the  $\mathcal{G}_{\text{att}}$ -hybrid model if:

1. **(UC Security)** For every PPT adversary  $\mathcal{A}$  controlling the server and up to  $n-1$  corrupted users, there exists a PPT simulator  $\mathcal{S}$  such that for all PPT environments  $\mathcal{Z}$ :

$$\text{REAL}_{\Pi_{\text{DP}}, \mathcal{A}, \mathcal{Z}}^{\mathcal{G}_{\text{att}}} (1^\lambda) \approx_c \text{IDEAL}_{\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}, \mathcal{S}, \mathcal{Z}} (1^\lambda).$$

2. **(DP Guarantee for the Explicit Release)** Let  $\text{Rel}(D)$  denote the transcript of values explicitly released from the enclave when  $\Pi_{\text{DP}}$  is run on dataset  $D$ , i.e.,  $\text{Rel}(D) := (\tilde{y}_1, \dots, \tilde{y}_R)$ . Then for any two neighboring datasets  $D, D'$  (differing in one user's updates across all  $R$  epochs) and any measurable set  $S$ :

$$\Pr[\text{Rel}(D) \in S] \leq e^\epsilon \cdot \Pr[\text{Rel}(D') \in S] + \delta.$$

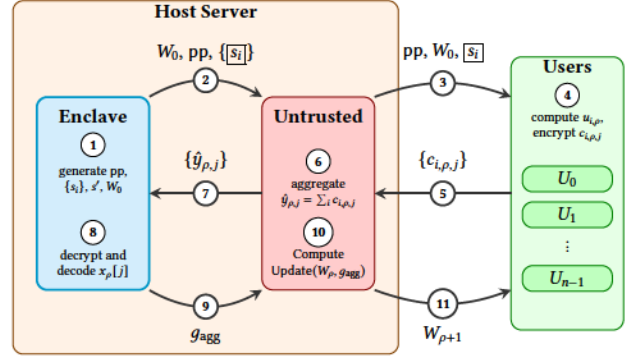


Figure 4: Protocol overview. **Setup (steps 1–3):** enclave generates pp, keys  $\{s_i\}, s'$ , and  $W_0$ ; server distributes public parameters and relays attested user keys  $\{s_i\}$  via secure channel. **Training rounds (steps 4–11, repeated  $R$  times):** users compute local updates  $u_{i,\rho}$  and encrypt as  $c_{i,\rho,j}$ ; server aggregates  $\hat{y}_{\rho,j}$  and forwards to enclave; enclave decrypts to recover  $x_\rho$ ; enclave outputs  $g_{\text{agg}}$ ; server prepares and broadcasts  $\Delta W_\rho$ .

The same bound holds for any (possibly adversarial) post-processing of  $\text{Rel}(D)$  together with arbitrary auxiliary information (e.g., corrupted users' states and public parameters).

Part 1 ensures that  $\Pi_{\text{DP}}$  is a faithful implementation of  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$ , guaranteeing no information leaks beyond the sanitized aggregate. Part 2 ensures that the sanitized aggregate itself is differentially private. Together, they provide end-to-end privacy: the cryptographic layer prevents leakage during aggregation, and the DP layer bounds the information content of the released output.

## 5 Protocol Description: SCALE-FL

### 5.1 Protocol Overview

*Protocol at a glance.* SCALE-FL runs in two phases. During setup, the SGX enclave generates TERSE public parameters and correlated key material, then provisions each user with its TERSE secret key over an attested secure channel. This setup is the only time users interact directly with the enclave, and it does not require any additional non-colluding server. During each training epoch, every user computes a local update, encodes it into  $d$  scalars, and encrypts each coordinate under an epoch-specific TERSE timestamp. The server collects the ciphertexts and aggregates them by coordinate-wise addition, producing a single aggregated ciphertext vector of length  $d$ , and forwards only this vector to the enclave. The enclave finalizes TERSE decryption to recover the plaintext aggregate, and in  $\Pi_{\text{DP}}$  it applies the DP mechanism to the aggregate before releasing the sanitized value to the server. The server then updates the model and broadcasts the new weights to users for the next epoch.

This division of labor yields a single-server protocol whose trusted component does not scale with the number of clients participating in an epoch. In each epoch, the server handles  $n$  inbound messages and performs only coordinate-wise ciphertext additions, while the enclave receives and processes one aggregated ciphertext vector of length  $d$ . In particular, the enclave's online work and network input per epoch scale with  $d$  and the number of epochs  $R$ , but not with  $n$ . Users also avoid per-epoch attestation or key exchange

**Construction  $\Pi$  – Setup**Setup( $1^\lambda, n, R$ ):

- (1)  $(pp, \{s_i\}_{i \in [n]}, s') \leftarrow \text{TERSE.Setup}(\lambda, t, n)$ .
- (2)  $st_{\text{enc}} \leftarrow (s', pp)$ .
- (3)  $W_0 \leftarrow \text{InitModel}()$ .
- (4) For each  $i \in [n]$ : provision  $s_i$  to user  $U_i$  via attested secure channel.
- (5) Publish  $pp$  and broadcast  $W_0$  to all parties.
- (6) Return  $st_{\text{enc}}$ .

**Figure 5: Base protocol  $\Pi$  – setup.**

with the enclave because all user–enclave interaction is confined to the one-time setup phase, and TERSE protects client updates during the server-mediated training phase.

$\Pi$  vs.  $\Pi_{\text{DP}}$  Finally, SCALE-FL comes in two variants that share the same aggregation pipeline. The architecture is demonstrated in  $\Pi$ , which uses SGX only to keep the TERSE aggregation key material off the server and recover the plaintext aggregate. The aggregate is then released for model training. This is useful as a stepping stone but, by itself, does not provide FL privacy in our threat model, as the final aggregate is still released to the adversary. The DP-enhanced protocol  $\Pi_{\text{DP}}$  addresses this by enforcing global differential privacy inside the enclave, so that only a sanitized aggregate is released and the raw aggregate never leaves the trusted boundary.

## 5.2 Setup

The Setup algorithm (Figure 5) is executed once before training begins. The enclave runs  $\text{TERSE.Setup}$  [41], which generates public parameters  $pp$  (including the hash function  $H$ , plaintext modulus  $t$ , ciphertext modulus  $q$ , and block size  $N$ ), a per-user secret key  $s_i$  for each user, and an aggregator secret  $s'$  used for decryption. The enclave retains  $s'$  in its sealed state, provisions each user with its  $s_i$  through an attested channel, and initializes the global model. Users do not communicate with the enclave during training, and only the server interacts with the enclave online.

The call to  $\text{TERSE.Setup}$  generates public parameters  $pp$  (including the hash function  $H$ , plaintext modulus  $t$ , ciphertext modulus  $q$ , and block size  $N$ ), per-user secret keys  $\{s_i\}_{i \in [n]}$ , and the aggregator secret  $s'$  used for decryption [41]. The enclave retains  $s'$  in its sealed state, and users receive only their individual  $s_i$ .

*Key provisioning and attestation.* Each user verifies the enclave program identity via remote attestation before accepting its key. Attestation binds the key-delivery payload to the intended enclave program. This ensures that the server, despite acting as a communication relay, cannot learn or modify provisioned secret keys [36].

## 5.3 Training

After setup, each client has remotely attested the enclave and received its per-user TERSE key via an attested secure channel, while the enclave retains the TERSE aggregation key material. Training then proceeds for  $R$  epochs. Importantly, during training clients communicate only with the untrusted server (i.e., no per-epoch client–enclave attestation or secure channel is required, unless the enclave is reinitialized and keys are reprovisioned).

*Local Training and Encoding* In a given epoch  $\rho$ , each client  $U_i$  performs local training on its private dataset  $D_i$  to compute an update vector  $u_{i,\rho} \in \mathbb{R}^d$  (e.g., a gradient or model delta). Because TERSE encrypts over an integer plaintext space, each client converts  $u_{i,\rho}$  to an integer vector using  $x_{i,\rho} = \text{TERSE.Encode}(S, u_{i,\rho})$ .

We instantiate TERSE’s timestamp interface with an epoch-and-coordinate indexing rule. TERSE encrypts scalar values indexed by a timestamp  $ts = (\theta, \tau)$  [41]. To support a  $d$ -dimensional update in each epoch, we map each (epoch  $\rho$ , coordinate  $j$ ) pair to a unique TERSE timestamp. We linearize to a stream index  $\ell := (\rho - 1) \cdot d + (j - 1)$  and set  $ts_{\rho,j} := (\lfloor \ell/N \rfloor, \ell \bmod N)$ . This ensures every encrypted coordinate uses a distinct timestamp.

Next, client  $U_i$  encrypts its encoded update coordinate-wise under TERSE. For each  $j \in \{1, \dots, d\}$ , it computes

$$c_{i,\rho,j} \leftarrow \text{TERSE.Encrypt}(s_i, ts_{\rho,j}, x_{i,\rho}[j]).$$

The client then transmits the ciphertext vector  $(c_{i,\rho,1}, \dots, c_{i,\rho,d})$  to the untrusted server. At this point, the client has crossed the trust boundary. The server learns only ciphertexts and public metadata such as the epoch identifier and timestamps. It does not have the TERSE aggregation key material needed to recover plaintext updates or the plaintext aggregate.

*Correctness conditions.* We assume parameters and encoding are chosen so that honest executions stay within TERSE’s correctness regime. First, for each coordinate, the integer aggregate of encoded updates remains in the centered range, so the plaintext sum does not wrap modulo  $t$ . Second, the TERSE parameters  $(q, t, n)$  satisfy the scheme’s stated correctness bounds, so that the accumulated encryption noise remains below the decoding threshold and  $\text{TERSE.Decrypt}$  recovers the intended plaintext aggregate [41]. A malicious server or corrupted clients can violate these conditions by submitting out-of-range updates or by modifying ciphertexts in transit. In that case, the recovered aggregate and the resulting model update may be incorrect. Our guarantees remain privacy-focused, and we do not provide integrity against such active manipulation.

*Server-Side Public Aggregation (outside SGX)* Upon receiving ciphertext vectors from clients in epoch  $\rho$ , the server aggregates them coordinate-wise using TERSE’s public addition procedure. For each coordinate  $j \in \{1, \dots, d\}$ , it computes an aggregated ciphertext

$$\hat{y}_{\rho,j} \leftarrow \text{TERSE.Add}(\{c_{i,\rho,j}\}_{i \in P_\rho}),$$

where  $P_\rho$  denotes the set of participating clients in epoch  $\rho$ . This step uses only public parameters and does not require any TERSE secret key material. Therefore, the server cannot decrypt either individual updates nor their aggregate. The server then forwards only the aggregated ciphertext vector  $(\hat{y}_{\rho,1}, \dots, \hat{y}_{\rho,d})$  to the enclave for aggregation recovery.

*Enclave Finalization and Decoding (inside SGX).* Upon receiving the aggregated ciphertext vector  $(\hat{y}_{\rho,1}, \dots, \hat{y}_{\rho,d})$  from the server, the enclave finalizes TERSE aggregation recovery using the aggregation key material that remains sealed inside SGX. For each coordinate  $j \in \{1, \dots, d\}$ , it computes

$$x_{\text{agg},\rho}[j] \leftarrow \text{TERSE.Decrypt}(s', ts_{\rho,j}, \hat{y}_{\rho,j}) \in \mathbb{Z}_t,$$

yielding an aggregate plaintext vector  $x_{\text{agg},\rho} \in \mathbb{Z}_t^d$ . The enclave then decodes this integer representation back to the real domain

by mapping each coordinate to its centered representative and rescaling by  $S$  using  $g_{\text{agg},\rho} = \text{TERSE.Decode}(S, x_{\text{agg},\rho})$ .

*Enclave Release (inside SGX).* In the base protocol  $\Pi$ , after finalizing aggregation recovery and decoding, the enclave releases the plaintext aggregate update  $g_{\text{agg},\rho}$  to the untrusted server. The server never obtains the TERSE aggregation key material, and it never observes any intermediate plaintext values. The enclave output is the only plaintext value derived from TERSE ciphertexts that leaves the trusted boundary during training.

*Server Update, Broadcast, and Termination (outside SGX).* Upon receiving  $g_{\text{agg},\rho}$  from the enclave, the server applies the learning rule to update the global model weights:

$$W_{\rho+1} \leftarrow \text{FL.Update}(W_{\rho}, g_{\text{agg},\rho}).$$

It then broadcasts  $W_{\rho+1}$  to all clients to begin the next epoch. After  $R$  epochs, the server outputs the final model  $W_R$  to the model owner (or publishes it according to the deployment).

*Correctness assumptions.* We assume parameters are chosen so that (i) the component-wise sum of all users' encoded updates stays within the plaintext range and therefore does not wrap modulo  $t$  (no plaintext overflow), and (ii) the TERSE parameters  $(q, t, n)$  satisfy the scheme's stated correctness bounds, which keep the effective noise small enough that reduction modulo  $t$  removes it and recovers the intended plaintext [41]. Under adversarial manipulation (e.g., corrupt users submitting out-of-range inputs or the server manipulating ciphertexts), correctness may be violated. In this case, we continue to provide privacy, but not any integrity guarantees.

Figure 6 defines the training-phase algorithms of the base protocol. Each is a direct invocation of the corresponding TERSE primitive [41], with the exception of Update, which applies the learning rule. By the correctness of TERSE decryption, Dec recovers exactly  $x_{\text{agg}} = \sum_{i=1}^n x_{i,\rho}$ , the coordinate-wise sum of all users' encoded updates. In the base protocol, the enclave releases  $x_{\text{agg}}$  directly to the server. We provide a protocol summary in Figure 7.

<p><b>Construction <math>\Pi</math> – Training</b></p> <p><math>\text{Enc}(k_i, ts, x)</math>: Return <math>c \leftarrow \text{TERSE.Encrypt}(k_i, ts, x)</math>.</p> <p><math>\text{Add}(c_1, \dots, c_n)</math>: Return <math>\hat{y} \leftarrow \text{TERSE.Add}(c_1, \dots, c_n)</math>.</p> <p><math>\text{Dec}(\text{st}_{\text{enc}}, ts, \hat{y})</math>: Return <math>x_{\text{agg}} \leftarrow \text{TERSE.Decrypt}(s', ts, \hat{y})</math>.</p> <p><math>\text{Update}(W_{\rho}, g_{\text{agg}})</math>: Return <math>W_{\rho+1} \leftarrow \text{FL.Update}(W_{\rho}, g_{\text{agg}})</math>.</p>
---

Figure 6: Base protocol  $\Pi$  – training algorithms.

After  $R$  epochs, the final model  $W_R$  is sent to the model owner.

## 5.4 Limitations of the Base Protocol

The base protocol  $\Pi$  guarantees aggregation security (Definition 3): the server cannot learn individual users' plaintext updates from the aggregation protocol itself. However, this guarantee is insufficient

for end-to-end privacy in FL. The enclave releases the plaintext aggregate  $x_{\text{agg},\rho} = \sum_{i=1}^n x_{i,\rho}$  every epoch, and after  $R$  epochs the final model  $W_R$  is delivered to the model owner. These releases can leak information about participants' private data even when secure aggregation protects the aggregation step.

*Worst-case leakage under  $(n-1)$ -corruption.* Under our threat model (Section 4.3), the adversary controls the server and up to  $n-1$  corrupted users  $C = [n] \setminus \{i^*\}$ . Since corrupted users' inputs are known to the adversary, the released aggregate reveals the remaining honest user's contribution exactly:

$$x_{i^*,\rho} = x_{\text{agg},\rho} - \sum_{j \in C} x_{j,\rho}.$$

Aggregation security prevents the server from computing this subtraction *during* the protocol (because it never sees  $x_{\text{agg},\rho}$  in the clear), but the released output makes it possible *after* decryption.

*Inference attacks on released aggregates and models.* Even without worst-case corruption, released aggregates and model parameters are known to enable a range of inference attacks. Aggregate deconstruction attacks can allow an adversary to infer information about a user's individual gradients, ultimately leading leakage about individual training samples [23, 43]. Membership inference attacks can determine whether a specific record was used in training [12]. Property inference attacks can extract unintended dataset attributes from model updates [22]. Model inversion attacks can recover representative samples of training classes [11, 45]. These vulnerabilities are well-documented in standard PPFL with secure aggregation.

*Architectural opportunity.* As discussed in Section 3.4, central DP achieves substantially better privacy–utility tradeoffs than local DP and remains meaningful under  $(n-1)$ -corruption when the curator is trusted. Our architecture provides a natural instantiation: the SGX enclave already recovers the plaintext aggregate but does not expose it to the server. Thus, a DP mechanism can be applied inside the enclave before any information is released to ensure privacy. Section 5.5 describes this extension.

## 5.5 DP Extension

We extend the base protocol  $\Pi$  to the DP-enhanced protocol  $\Pi_{\text{DP}}$ , summarized in Figure 7. The cryptographic aggregation primitives Enc, Add, and TERSE.Decrypt are unchanged. The DP extension modifies the message flow in two places:

- User-side preprocessing (Training Epoch step 2b):** each user applies  $\mathcal{M}.\text{Preprocess}$  to its real-valued update before encoding and encrypting.
- Enclave post-processing (Training Epoch step 4):** after decrypting and decoding the aggregate, the enclave applies  $\mathcal{M}.\text{ApplyDP}$  and releases only the *sanitized* aggregate.

*Setup modifications.* Setup follows the same structure as in  $\Pi$ , with two additions shown in Figure 7. First, the public parameters  $\text{pp}$  include the DP mechanism parameters (e.g., clipping bound  $C$ ). Second, the enclave initializes and stores the DP mechanism state  $\text{st}_{\mathcal{M}}$  alongside the TERSE key material.

*DP Mechanism.* The construction is parameterized by any DP mechanism  $\mathcal{M}$  satisfying the interface in Section 3.4. The enclave boundary ensures  $\mathcal{M}$  is applied correctly despite the server's behavior, and the server observes only  $\tilde{g}_{\text{agg}}$ , never the raw aggregate.

**Protocol  $\Pi_{\text{DP}}$  – Complete Protocol Summary****Setup.**

- (1) Enclave runs  $\text{TERSE.Setup}(\lambda, t, n)$  to generate pp, user keys  $\{s_i\}_{i \in [n]}$ , and aggregator key  $s'$  [41].
- (2) Enclave initializes model  $W_0$ , DP mechanism state  $\text{st}_{\mathcal{M}}$ , and stores  $\text{st}_{\text{enc}} \leftarrow (s', \text{pp}, \text{st}_{\mathcal{M}})$ .
- (3) Each user  $U_i$  receives  $s_i$  via attested secure channel (server relays only) [36].
- (4) Enclave publishes pp (including DP parameters) and broadcasts  $W_0$  to all parties.

**Training Epoch  $\rho = 1, \dots, R$ .**

- (1) Server broadcasts  $W_\rho$  to all users.
- (2) Each user  $U_i$ :
  - (a) Computes local update  $u_{i,\rho} \leftarrow \text{LocalTrain}(W_\rho; D_i)$ .
  - (b) Preprocesses  $\tilde{u}_{i,\rho} \leftarrow \mathcal{M}.\text{Preprocess}(u_{i,\rho}, \text{pp})$ .
  - (c) Encodes  $x_{i,\rho} \leftarrow (\tilde{u}_{i,\rho}) \in \mathbb{Z}_t^d$ .
  - (d) For each  $j = 1, \dots, d$ : computes  $c_{i,\rho,j} \leftarrow \text{Enc}(k_i, ts_{\rho,j}, x_{i,\rho}[j])$ .
 Sends  $\{c_{i,\rho,j}\}_{j=1}^d$  to the server.
- (3) Server computes  $\hat{y}_{\rho,j} \leftarrow \text{Add}(\{c_{i,\rho,j}\}_{i=1}^n)$  for each  $j$  and forwards  $\{\hat{y}_{\rho,j}\}_{j=1}^d$  to the enclave.
- (4) Enclave:
  - (a) For each  $j$ : recovers  $x_{\text{agg}}[j] \leftarrow \text{TERSE.Dec}(\text{st}_{\text{enc}}, ts_{\rho,j}, \hat{y}_{\rho,j})$ .
  - (b) Decodes  $g_{\text{agg}} \leftarrow \text{Decode}(x_{\text{agg}})$ .
  - (c) Computes  $(\tilde{g}_{\text{agg}}, \text{st}'_{\mathcal{M}}) \leftarrow \mathcal{M}.\text{ApplyDP}(\text{st}_{\mathcal{M}}, g_{\text{agg}})$ ; updates  $\text{st}_{\mathcal{M}} \leftarrow \text{st}'_{\mathcal{M}}$ .
  - (d) Releases  $\tilde{g}_{\text{agg}}$  to the server.
- (5) Server computes  $W_{\rho+1} \leftarrow \text{FL.Update}(W_\rho, \tilde{g}_{\text{agg}})$  and broadcasts  $W_{\rho+1}$  to all users.

**After Training.**

- (1) After  $R$  epochs,  $W_R$  is delivered to the model owner.

Figure 7: Complete protocol summary for  $\Pi_{\text{DP}}$ . The base protocol  $\Pi$  is identical except that step 2b is omitted and step 4 releases the raw aggregate  $g_{\text{agg}}$  without applying DP.

## 6 Security Analysis

### 6.1 Overview and assumptions

We state three security claims for our protocols. Theorem 1 establishes that the base protocol provides aggregation security. Theorem 2 extends this to the DP-enhanced protocol, and Theorem 3 establishes that the released transcript is differentially private. Theorems 2 and 3 combined create FL privacy for  $\Pi_{\text{DP}}$ . The cryptographic layer prevents unintended leakage during aggregation, and the DP layer bounds the information content of the released output.

*Aggregator obliviousness.* We use the TERSE notion of aggregator obliviousness in the encrypt-once model [41, Def. 1]. The adversary may corrupt users and/or the aggregator and obtain encryptions subject to the encrypt-once restriction, then is challenged on fresh ciphertexts at timestamp  $ts^*$ . If the aggregator is corrupted, the two challenge worlds must agree on the aggregated sum. By [41, Thm. 1], TERSE achieves this notion under RLWE.

### 6.2 Aggregation Security of $\Pi$

**THEOREM 1 (AGGREGATION SECURITY OF  $\Pi$ ).** *Using TERSE’s aggregator-obliviousness security [41, Theorem 1],  $\Pi$  securely realizes  $\mathcal{F}_{\text{agg}}$  in the  $\mathcal{G}_{\text{att}}$ -hybrid model (Definition 3).*

*Proof sketch.* Protocol  $\Pi$  combines TERSE as a black-box private stream aggregation scheme with an attested TEE modeled in the  $\mathcal{G}_{\text{att}}$ -hybrid model. Accordingly, the security argument reduces to TERSE aggregator-obliviousness and the fact that the adversary

only learns what the enclave explicitly outputs. In each epoch, honest users send only TERSE ciphertexts to the server, and the server performs only TERSE public aggregation. By TERSE aggregator-obliviousness (TERSE Theorem 1 [41], under RLWE), the server cannot distinguish encryptions of honest users’ updates from encryptions of unrelated values, even when it colludes with up to  $n-1$  corrupted users. Our timestamps ensure the encrypt-once condition required by TERSE. Since the enclave is modeled via  $\mathcal{G}_{\text{att}}$ , the adversary cannot access enclave secrets or intermediate plaintexts, so the only plaintext information it receives is the aggregate value released by the enclave. Finally,  $\mathcal{F}_{\text{agg}}$  allows a malicious server to bias the aggregate by an offset, and the real protocol matches this because a deviating server can cause the enclave to output some aggregate value. Therefore,  $\Pi$  realizes  $\mathcal{F}_{\text{agg}}$ .

### 6.3 Privacy of $\Pi_{\text{DP}}$

**THEOREM 2 (UC SECURITY OF  $\Pi_{\text{DP}}$ ).** *Using TERSE’s aggregator-obliviousness security (TERSE Theorem 1 [41]),  $\Pi_{\text{DP}}$  securely realizes  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  in the  $\mathcal{G}_{\text{att}}$ -hybrid model (Definition 4, Part 1).*

*Proof sketch.* The proof is a simulation argument that follows the same structure as Theorem 1. The main difference is that the enclave releases only the DP-sanitized aggregate  $\tilde{g}_\rho$ , so the simulator must reproduce the adversary’s view using only what  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  reveals. As in Theorem 1, TERSE aggregator-obliviousness lets the simulator replace honest users’ ciphertexts with encryptions of 0

without changing the adversary’s view, up to computational indistinguishability. Inside the enclave, any deviation by the malicious server can be represented as an additive offset  $\delta_\rho$  to the plaintext aggregate for epoch  $\rho$  over  $\mathbb{Z}_t^d$ . The simulator computes the corresponding  $\delta_\rho$ , provides  $\delta_\rho$  to  $\mathcal{F}_{\text{agg-DP}}^M$ , and forwards  $\mathcal{F}_{\text{agg-DP}}^M$ ’s output  $\tilde{g}_\rho$  to the adversary. By the  $\mathcal{G}_{\text{alt}}$ -hybrid model, the adversary learns only the enclave’s explicit outputs and cannot observe internal computations, so the simulated view matches the real view up to computational indistinguishability. Hence  $\Pi_{\text{DP}}$  realizes  $\mathcal{F}_{\text{agg-DP}}^M$ .

**THEOREM 3 (DP GUARANTEE OR  $\Pi_{\text{DP}}$ ).** *Let  $\mathcal{M}$  satisfy  $(\epsilon, \delta)$ -differential privacy with user-level adjacency. Then  $\Pi_{\text{DP}}$  satisfies Definition 4, Part 2: the transcript of values explicitly released from the enclave is  $(\epsilon, \delta)$ -differentially private.*

*Proof sketch.* In  $\Pi_{\text{DP}}$ , the only values explicitly released from the enclave are the DP-sanitized aggregates  $\text{Rel}(D) := (\tilde{g}_0, \dots, \tilde{g}_{R-1})$ . By construction,  $\text{Rel}(D)$  is exactly the output of the DP mechanism  $\mathcal{M}$  applied to the sequence of per-epoch aggregates computed from user updates. Since  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -DP with user-level adjacency for the full  $R$ -epoch transcript, it follows directly from the definition of differential privacy that  $\text{Rel}(D)$  is  $(\epsilon, \delta)$ -DP. Any additional values the server computes from  $\text{Rel}(D)$ , including any model parameters computed using Update and any other derived transcript, are post-processing of an  $(\epsilon, \delta)$ -DP output together with the adversary’s side information. By the post-processing property of differential privacy, these derived values remain  $(\epsilon, \delta)$ -DP. Therefore, the transcript of values explicitly released from the enclave is  $(\epsilon, \delta)$ -DP.

*Compositional DP.* The theorem above assumes a fixed-parameter mechanism  $\mathcal{M}$  calibrated once to provide an overall  $(\epsilon, \delta)$  guarantee for the  $R$ -epoch release transcript. If the enclave applies a per-epoch DP mechanism but targets a global  $(\epsilon, \delta)$  across epochs, then a privacy accountant must be specified and shown to bound the composed privacy loss (e.g., via standard composition, Rényi-DP, moments accounting). DP accounting is orthogonal to our protocol proof and follows from standard DP results once  $\mathcal{M}$  is fixed.

## 7 Evaluation

We evaluate SCALE-FL to quantify its computational overhead relative to plaintext aggregation, the scalability benefit of moving the  $O(n)$  fan-in outside the enclave, and the amortized preprocessing cost required by TERSE. Our primary performance results report *server-side per-epoch runtime*, defined as the runtime of the untrusted server plus enclave. We report client overhead separately.

*Reproducibility and artifacts.* We release our implementation as open source (anonymized) [34]. In addition to the cryptographic and SGX components, the release includes an integration with the Flower FL framework that allows custom FL workloads to use SCALE-FL with minimal changes to standard Flower training loops. We also include an end-to-end Flower simulation on CIFAR-10 based on a PyTorch tutorial model, which we use for utility experiments.

### 7.1 Experimental Setup

*Hardware platform.* All experiments run on a single machine and do not use a distributed deployment. The machine supports Intel SGX with a 128 GB EPC, has 384 GB RAM, and uses an Intel Xeon

Gold 5412U CPU with 24 physical cores and 48 hardware threads, up to 3.9 GHz. We execute enclave code using Gramine, a library OS for running applications inside Intel SGX enclaves. Because all protocol components run locally, we report computation time only and exclude network latency. [16]

*Benchmarking methodology.* We report wall-clock (elapsed) computation time throughout, including enclave timings measured inside Gramine. All client-side benchmarks were restricted to a single software thread pinned to one CPU core to approximate per-device latency. Server-side measurements ran with unrestricted host parallelism, and enclave execution used Gramine-SGX with `sgx.max_threads = 16` set in the manifest template [17].

*Measured components.* We measure and combine per-epoch costs for: (i) the SCALE-FL client-side fixed-point encoding and encryption of a  $d$ -dimensional update, (ii) the SCALE-FL untrusted server-side coordinate-wise public addition, and (iii) the SCALE-FL enclave-side finalization and decryption of a single aggregated ciphertext vector of length  $d$ . We exclude local training and differential privacy times since they are orthogonal to the aggregation mechanism and common to all schemes. We exclude local training and differential privacy times since they are orthogonal to the aggregation mechanism and common to all schemes.

*Baselines.* We compare against two aggregation-pipeline baselines (untrusted server plus enclave): a plaintext baseline that directly sums the  $n$  update vectors in the clear, and an AES-to-enclave baseline in which each user encrypts its update under a per-user AES key and sends it to the enclave, which manages keys, decrypts all  $n$  vectors, and sums them, incurring  $O(n)$  enclave work and enclave interactions per epoch.

We omit client-side baseline timings because plaintext has no cryptographic cost and AES encryption is standard. We focus instead on the server and enclave scaling behavior, and including baseline client timings would substantially increase the figure and table space.

### 7.2 SCALE-FL Per-epoch Overhead vs. Plaintext Aggregation and AES-to-enclave Baseline

We compare server-side per-epoch runtime against plaintext aggregation and the AES-to-enclave baseline. Client costs are recorded separately in Section 7.3. For each  $n \in \{10k, 50k, 100k\}$ , we measure server-side per-epoch runtime for plaintext summation of  $n$  update vectors, SCALE-FL (untrusted server ciphertext aggregation plus enclave finalization and decryption), and the AES-to-enclave baseline. We use update dimension  $d = 10k$  parameters. Figure 8 reports server-side per-epoch runtime as a function of  $n$ . SCALE-FL closely matches plaintext aggregation, with under 1% additional overhead across all  $n$ . In contrast, the AES-to-enclave baseline is substantially slower and scales linearly with  $n$  because it performs  $O(n)$  decryptions and aggregation inside the enclave.

### 7.3 Client-side Overhead

*Client precomputation (offline).* We measure the amortized per-epoch client precomputation time as a function of  $d$  (at  $n = 10k$  clients). As shown in Figure 9, this offline setup cost is small and grows roughly linearly with  $d$ , from 0.228 ms at  $d = 10k$  to 2.044 ms

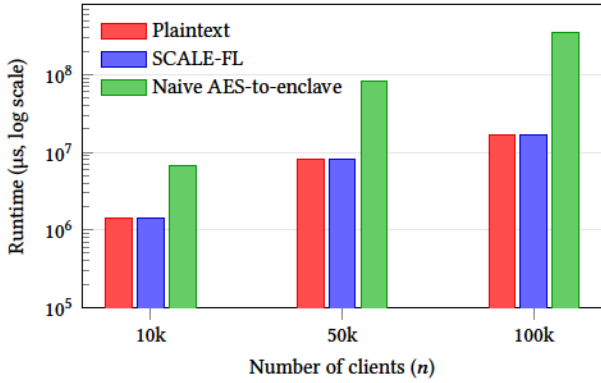


Figure 8: Server-side per-epoch runtime (untrusted server plus enclave) for plaintext aggregation, SCALE-FL, and a naive AES-to-enclave baseline, for  $n \in \{10k, 50k, 100k\}$  and  $d = 10k$ . The y-axis uses a logarithmic scale.

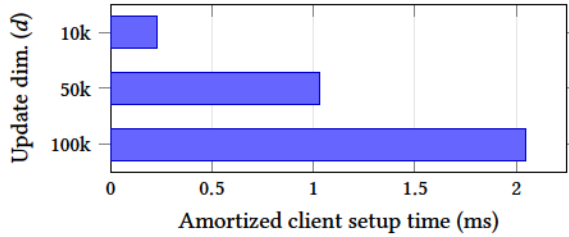


Figure 9: Average amortized per-client setup time at  $n = 10k$  and  $d \in \{10k, 50k, 100k\}$ .

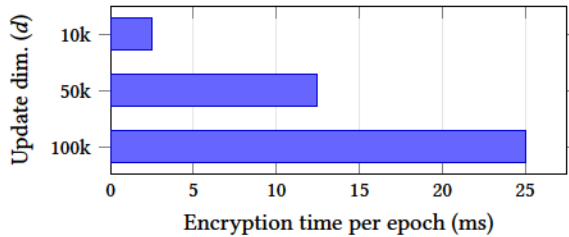


Figure 10: Average per-client encryption time per epoch for  $n = 10k$  and  $d \in \{10k, 50k, 100k\}$ .

at  $d = 100k$ . This work is not on the online aggregation critical path and can be computed ahead of time.

*Client encoding and encryption (online).* We measure the per-client time to encrypt a  $d$ -dimensional update (at  $n = 10k$  clients). Figure 10 shows that online overhead scales approximately linearly with  $d$ , increasing from 2.531 ms at  $d = 10k$  to 25.008 ms at  $d = 100k$ .

### 7.4 Enclave Precomputation Cost

TERSE uses key-dependent preprocessing that can be computed offline and reused across multiple epochs within a block. We report two amortized in-enclave preprocessing costs. The online aggregation pipeline is unchanged, since the untrusted server aggregates ciphertexts and the enclave finalizes and decrypts one aggregated ciphertext vector of length  $d$ . We measure amortized enclave preprocessing time per epoch for update dimension  $d = 10k$ . We

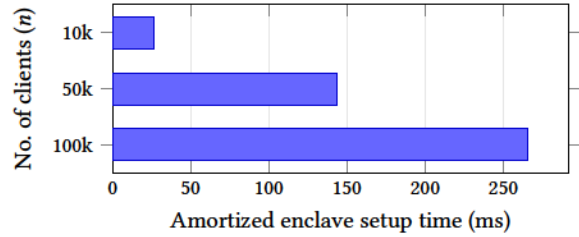


Figure 11: Amortized enclave setup time per epoch for  $d = 10k$  and  $n \in \{10k, 50k, 100k\}$ .

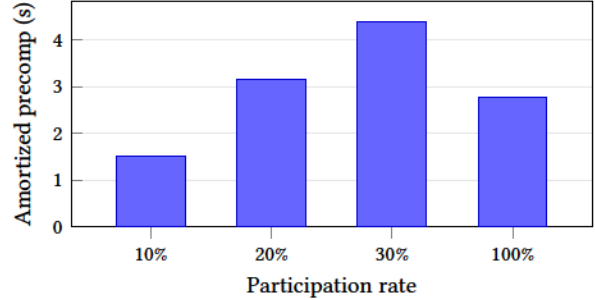


Figure 12: Amortized SGX enclave precomputation time per epoch under client sampling for  $n = 1,000,000$  and  $d = 10k$ .

measure (i) amortized enclave setup time per epoch as a function of  $n \in \{10k, 50k, 100k\}$  (Figure 11) and (ii) sampling-related enclave preprocessing for a fixed total population of  $n = 1,000,000$  clients under participation rates of 10%, 20%, 30%, and 100% (Figure 12).

Setup time grows from 26.3 ms at  $n = 10k$  to 265.5 ms at  $n = 100k$  (Figure 11). Under client sampling with  $n = 1,000,000$ , preprocessing ranges from 1.51 s to 4.38 s per epoch depending on participation (Figure 12). The dependence on participation rate is not strictly monotonic, which we attribute to differences in the amount of intermediate state stored and processed during preprocessing that affect enclave memory behavior and cache locality. Both costs are computed offline and are not on the online critical path.

### 7.5 Summary

SCALE-FL achieves near-plaintext server-side performance while keeping the enclave’s online work per epoch independent of the number of clients. For  $d = 10k$  and  $n \in \{10k, 50k, 100k\}$ , SCALE-FL adds  $< 1\%$  server-side overhead relative to plaintext aggregation. The AES-to-enclave baseline is  $5\times$  to  $21\times$  slower over the same range of  $n$  and scales linearly with  $n$  due to per-client enclave decryption and aggregation. Offline enclave preprocessing is modest. Amortized enclave setup grows from 26.3 ms to 265.5 ms as  $n$  increases from 10k to 100k, and sampling-related preprocessing ranges from 1.51 s to 4.38 s per epoch for  $n = 1,000,000$  total clients depending on participation.

## 8 Discussion

### 8.1 Client Sampling

In cross-device FL, communication constraints and intermittent availability can make full participation impractical at large scale. Many FedAvg-style analyses model this by sampling a subset of

clients each epoch independently across epochs [15, 27]. We include a simple modification to SCALE-FL to support client sampling.

*Implementation.* For each epoch  $\rho \in [R]$ , let  $P_\rho \subseteq [n]$  be the set of participating users, with  $m_\rho := |P_\rho|$ . Only users in  $P_\rho$  submit ciphertexts for epoch  $\rho$ . Recall TERSE timestamps are  $ts = (\theta, \tau)$ , and the scheme defines the aggregator key  $s' = -\left[\sum_{i=0}^{n-1} s_i\right]_q$  [41].

With client sampling, the enclave preserves the same cancellation property by deriving an epoch-specific effective aggregator key for the sampled subset:  $s'_\rho := -\left[\sum_{i \in P_\rho} s_i\right]_q$ . The server aggregates only ciphertexts from  $i \in P_\rho$ . The enclave then finalizes using  $s'_\rho$ , which is equivalent to using  $p'_{\rho,ts} := (A_\theta \cdot s'_\rho)[\tau]$  in the above TERSE aggregation equation. The recovered plaintext is therefore the sampled aggregate  $\sum_{i \in P_\rho} x_{i,\rho}$ , coordinate-wise, and no additional key material is revealed to users or the server.

*Efficiency and preprocessing.* Client sampling adds only light-weight enclave-side preprocessing. After selecting and publishing  $P_\rho$ , the enclave computes  $s'_\rho = -\left[\sum_{i \in P_\rho} s_i\right]_q$  using  $(m_\rho - 1)$  additions over stored user keys. It then recomputes the TERSE key-dependent finalization terms for each TERSE block index  $\theta$  that appears in epoch  $\rho$ 's timestamps (recall  $ts = (\theta, \tau)$ ). [41]

With our timestamp mapping, a  $d$ -dimensional update spans at most  $\lceil d/N \rceil$  distinct  $\theta$  values per epoch. Even if the same  $\theta$  value appears again in the following epoch, the enclave must recompute the corresponding precomputation once, because the effective aggregator key  $s'_\rho$  changes with  $P_\rho$ . Therefore, the enclave performs at most one extra key-dependent multiplication per epoch. Our evaluations of this overhead are shown in Figure 12.

*Interaction with server behavior.* Client sampling does not introduce new trust assumptions. The enclave selects and publishes  $P_\rho$ , and the server remains as in the base protocol.

## 8.2 Dropout Handling

To handle user dropouts, the protocol incorporates Lagrange interpolation as in [20]. For each user  $i \in 1, \dots, n$  and active user set  $U$ , the Lagrange coefficient is defined as

$$\mathcal{L}_{i,U}(x) = \prod_{j \in U \setminus i} \frac{x - j}{i - j}.$$

In practice, we use the simplified coefficient evaluated at zero,  $\mathcal{L}_{i,U}(0)$ . The aggregate secret polynomial  $q^d(x)$  is constructed as a degree- $d = n - 1$  polynomial uniquely determined by  $n$  points, where  $n$  is the initial number of users. Its remaining coefficients are sampled uniformly from the finite field  $\mathbb{F}_{q_{agg}}$ , and its constant term is fixed to zero (i.e.,  $q^d(0) = 0$ ).

During the setup, we then replace the user secret keys  $sk_i$  with a point on a polynomial multiplied by a Lagrange coefficient:

$$sk_i = q^d(i) \cdot \mathcal{L}_{i,U}$$

That is, during the key generation, each user  $i$  is provided with a set of polynomials  $q^{(2)}(i), \dots, q^{(n-1)}(i)$  and one data point for each secret polynomial  $q^d(i)$ , where the polynomial  $q^d(x)$  remains unknown to everyone because it is the sum of all users' random polynomials. Then, during the aggregation, the aggregator simply sums  $n$  ciphertexts  $c_{i,ts}$  received from  $n$  users such that the final aggregation value results in the sum of the users' inputs plus some

noise  $r_{i,ts}$  modulo the plaintext modulus. Since the constant terms in all the polynomials during aggregation are 0, the summation of the keys  $sk_i$  multiplied by the Lagrange coefficient becomes zero (or a multiple of  $q_{agg}$  where  $q_{agg}$  is the ciphertext modulus).

## 8.3 Weighted Averaging

The base protocol aggregates a uniform sum. To support weighted FedAvg (typically  $w_i = n_i$ ), we need the enclave to recover  $\sum_{i \in P_\rho} w_i \Delta_{i,\rho}$  and (optionally)  $\sum_{i \in P_\rho} w_i$ .

If weights are public and fixed, each client encrypts  $w_i \Delta_{i,\rho}$ . The enclave divides by  $W_{\Sigma,\rho} = \sum_{i \in P_\rho} w_i$  during post-decryption processing.

If weights are private or epoch-dependent, clients additionally encrypt  $w_i$  in a second TERSE scalar stream. The server aggregates both streams, and the enclave decrypts  $\sum_i w_i \Delta_{i,\rho}$  and  $\sum_i w_i$  and applies their ratio. This adds one scalar encryption per user per epoch (i.e.,  $d + 1$  instead of  $d$ ).

*Interaction with DP.* Weights change the sensitivity: clipping can be done before weighting (e.g., using  $C/w_{\max}$ ) to keep weighted contributions bounded by  $C$ . For private weights, noise must be calibrated to the worst-case admissible weight range, potentially enforced by the enclave [10].

## 8.4 IND-CPA<sup>D</sup> Attacks on RLWE-Based Schemes

IND-CPA<sup>D</sup> captures settings where an adversary can learn information from decryption behavior. Prior work gives key-recovery attacks when the adversary can observe exact decryptions or can adaptively probe a success or failure predicate that arises from decryption errors in practice [6, 24].

*No decryption oracle in our protocols.* In our use of TERSE, decryption does not expose a success or failure bit, and it does not provide an externally checkable correctness predicate. TERSE decryption always returns a plaintext modulo  $t$ , and any incorrectness would only appear as a normal aggregate that the server and corrupted clients cannot validate against an expected value. Thus the attack mechanism used in IND-CPA<sup>D</sup> attacks, (adaptive access to decryption outputs or decryption failure information) is not available in our setting.

Additionally, in  $\Pi_{DP}$ , only a sanitized aggregate is released to the server. This structurally removes the IND-CPA<sup>D</sup> attack surface, in the same spirit as approaches that restore security by releasing only noisy post-decryption outputs [25].

## 9 Conclusion

We presented SCALE-FL, a hybrid PPFL architecture that combines TERSE private stream aggregation with a minimal SGX finalization step, pushing  $O(n)$  aggregation to the untrusted server while keeping only  $O(d)$  work inside the enclave. This yields strong confidentiality of client updates against a malicious server and up to  $n - 1$  colluding clients, without additional non-colluding servers. Our DP-enhanced variant enforces central DP inside the enclave and releases only a sanitized aggregate. Experiments show near-plaintext aggregation performance at  $< 1\%$  server-side overhead.

## Acknowledgments

## References

- [1] Ghada Almashaqbeh and Zahra Ghodsi. 2025. AnoFel: Supporting Anonymity for Privacy-Preserving Federated Learning. *Proceedings on Privacy Enhancing Technologies* (2025).
- [2] Li Bai, Haibo Hu, Qingqing Ye, Haoyang Li, Leixia Wang, and Jianliang Xu. 2024. Membership Inference Attacks and Defenses in Federated Learning: A Survey. *ACM Comput. Surv.* 57, 4, Article 89 (Dec. 2024), 35 pages. <https://doi.org/10.1145/3704633>
- [3] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloe Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan, et al. 2019. Towards federated learning at scale: System design. *Proceedings of machine learning and systems* 1 (2019), 374–388.
- [4] Yihao Cao, Jianbiao Zhang, Yaru Zhao, Pengchong Su, and Haoxiang Huang. 2024. SRFL: A Secure & Robust Federated Learning framework for IoT with trusted execution environments. *Expert Systems with Applications* 239 (2024), 122410. <https://doi.org/10.1016/j.eswa.2023.122410>
- [5] Guoxing Chen, Sanchuan Chen, Yuan Xiao, Yinqian Zhang, Zhiqiang Lin, and Ten H.Lai. 2019. Sgxpectre: Stealing intel secrets from sgx enclaves via speculative execution. In *2019 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 142–157.
- [6] Jung Hee Cheon, Hyeongmin Choe, Alain Passelègue, Damien Stehlé, and Elias Suvanto. 2024. Attacks against the IND-CPAD security of exact FHE schemes. In *Proceedings of the 2024 ACM SIGSAC Conference on Computer and Communications Security*. 2505–2519.
- [7] Victor Costan and Srinivas Devadas. 2016. Intel SGX explained. *Cryptology ePrint Archive* (2016).
- [8] Jesse De Meulemeester, David Oswald, Ingrid Verbauwhede, and Jo Van Bulck. 2026. Battering RAM: Lowcost interposer attacks on confidential computing via dynamic memory aliasing. In *47th IEEE Symposium on Security and Privacy (S&P)*.
- [9] Cynthia Dwork. 2006. Differential privacy. In *International colloquium on automata, languages, and programming*. Springer, 1–12.
- [10] Cynthia Dwork, Aaron Roth, et al. 2014. The algorithmic foundations of differential privacy. *Foundations and trends® in theoretical computer science* 9, 3–4 (2014), 211–407.
- [11] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. 2015. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*. 1322–1333.
- [12] Jiqiang Gao, Boyu Hou, Xiaojie Guo, Zheli Liu, Ying Zhang, Kai Chen, and Jin Li. 2021. Secure aggregation is insecure: Category inference attack on federated learning. *IEEE Transactions on Dependable and Secure Computing* 20, 1 (2021), 147–160.
- [13] Robin C Geyer, Tassilo Klein, and Moin Nabi. 2017. Differentially private federated learning: A client level perspective. *arXiv preprint arXiv:1712.07557* (2017).
- [14] Johannes Götzfried, Moritz Eckert, Sebastian Schinzel, and Tilo Müller. 2017. Cache attacks on Intel SGX. In *Proceedings of the 10th European Workshop on Systems Security*. 1–6.
- [15] Ala Gouissem, Zina Chkirbene, and Ridha Hamila. 2024. A comprehensive survey on client selections in federated learning. *Innovation and Technological Advances for Sustainability* (2024), 417–428.
- [16] Gramine Project. 2023. *Gramine Documentation*. <https://gramine.readthedocs.io/> Accessed: Feb 27, 2026.
- [17] Gramine Project. 2026. *Manifest syntax*. <https://gramine.readthedocs.io/en/stable/manifest-syntax.html> Accessed: Feb 27, 2026.
- [18] Peishan Huang, Dong Li, and Zhigang Yan. 2023. Wireless federated learning with asynchronous and quantized updates. *IEEE Communications Letters* 27, 9 (2023), 2393–2397.
- [19] Yufan Jiang, Maryam Zarezadeh, Tianxiang Dai, and Stefan Köpsell. 2025. Al-phaFL: Secure Aggregation with Malicious Security for Federated Learning against Dishonest Majority. *Proceedings on Privacy Enhancing Technologies* (2025).
- [20] Taeho Jung, Junze Han, and Xiang-Yang Li. 2016. PDA: semantically secure time-series data analytics with dynamic user groups. *IEEE Transactions on Dependable and Secure Computing* 15, 2 (2016), 260–274.
- [21] Fumiyuki Kato, Yang Cao, and Masatoshi Yoshikawa. 2023. Olive: Oblivious Federated Learning on Trusted Execution Environment against the Risk of Sparsification. *Proc. VLDB Endow.* 16, 10 (June 2023), 2404–2417. <https://doi.org/10.14778/3603581.3603583>
- [22] Raouf Kerkouche, Gergely Acs, and Mario Fritz. 2023. Client-specific property inference against secure aggregation in federated learning. In *Proceedings of the 22nd Workshop on Privacy in the Electronic Society*. 45–60.
- [23] Maximilian Lam, Gu-Yeon Wei, David Brooks, Vijay Janapa Reddi, and Michael Mitzenmacher. 2021. Gradient disaggregation: Breaking privacy in federated learning by reconstructing the user participant matrix. In *International Conference on Machine Learning*. PMLR, 5959–5968.
- [24] Baiyu Li and Daniele Micciancio. 2021. On the security of homomorphic encryption on approximate numbers. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 648–677.
- [25] Baiyu Li, Daniele Micciancio, Mark Schultz-Wu, and Jessica Sorrell. 2022. Securing approximate homomorphic encryption using differential privacy. In *Annual International Cryptology Conference*. Springer, 560–589.
- [26] Li Li, Yuxi Fan, Mike Tse, and Kuo-Yi Lin. 2020. A review of applications in federated learning. *Computers & Industrial Engineering* 149 (2020), 106854. <https://doi.org/10.1016/j.cie.2020.106854>
- [27] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. 2020. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine* 37, 3 (2020), 50–60.
- [28] Hai Liu, Changgen Peng, Youliang Tian, Shigong Long, Feng Tian, and Zhenqiang Wu. 2022. Gdp vs. ldp: A survey from the perspective of information-theoretic channel. *Entropy* 24, 3 (2022), 430.
- [29] Ziyao Liu, Jiale Guo, Wenzhuo Yang, Jiani Fan, Kwok-Yan Lam, and Jun Zhao. 2022. Privacy-preserving aggregation in federated learning: A survey. *IEEE Transactions on Big Data* (2022).
- [30] Grigory Malinovsky, Dmitry Kovalev, Elmur Gasanov, Laurent Condat, and Peter Richtarik. 2020. From local SGD to local fixed-point methods for federated learning. In *International Conference on Machine Learning*. PMLR, 6692–6701.
- [31] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*. Pmlr, 1273–1282.
- [32] Fan Mo, Hamed Haddadi, Kleomenis Katevas, Eduard Marin, Diego Perino, and Nicolas Kourtellis. 2021. PPFL: Privacy-preserving federated learning with trusted execution environments. In *Proceedings of the 19th annual international conference on mobile systems, applications, and services*. 94–108.
- [33] Mohammad Naseri, Jamie Hayes, and Emiliano De Cristofaro. 2020. Local and central differential privacy for robustness and privacy in federated learning. *arXiv preprint arXiv:2009.03561* (2020).
- [34] nd-dsp-lab. 2026. aggregation-sgx-fl: SCALE-FL (Secure Aggregation + SGX + Federated Learning). <https://github.com/nd-dsp-lab/aggregation-sgx-fl>
- [35] Alexander Nilsson, Pegah Nikbakht Bideh, and Joakim Brorsson. 2020. A survey of published attacks on Intel SGX. *arXiv preprint arXiv:2006.13598* (2020).
- [36] Rafael Pass, Elaine Shi, and Florian Tramer. 2017. Formal Abstractions for Attested Execution Secure Processors. In *EUROCRYPT*.
- [37] Mayank Rathee, Conghao Shen, Sameer Wagh, and Raluca Ada Popa. 2023. ELSA: Secure Aggregation for Federated Learning with Malicious Actors. In *2023 IEEE Symposium on Security and Privacy (SP)*. 1961–1979. <https://doi.org/10.1109/SP46215.2023.10179468>
- [38] Mohamed Sabt, Mohammed Achemlal, and Abdelmadjid Bouabdallah. 2015. Trusted Execution Environment: What It is, and What It is Not. In *2015 IEEE Trustcom/BigDataSE/ISPA*, Vol. 1. 57–64. <https://doi.org/10.1109/Trustcom.2015.357>
- [39] Elaine Shi, HTH Chan, Eleanor Rieffel, Richard Chow, and Dawn Song. 2011. Privacy-preserving aggregation of time-series data. In *Annual Network & Distributed System Security Symposium (NDSS)*. Internet Society.
- [40] Tao Sun, Dongsheng Li, and Bao Wang. 2022. Decentralized federated averaging. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 4 (2022), 4289–4301.
- [41] Jonathan Takeshita, Zachariah Carmichael, Ryan Karl, and Taeho Jung. 2022. TERSE: tiny encryptions and really speedy execution for post-quantum private stream aggregation. In *International Conference on Security and Privacy in Communication Systems*. Springer, 331–352.
- [42] Hui-Po Wang, Dingfan Chen, Raouf Kerkouche, and Mario Fritz. 2024. FedLAP-DP: Federated Learning by Sharing Differentially Private Loss Approximations. *Proceedings on Privacy Enhancing Technologies* (2024).
- [43] Zhibo Wang, Zhiwei Chang, Jiahui Hu, Xiaoyi Pang, Jiacheng Du, Yongle Chen, and Kui Ren. 2024. Breaking secure aggregation: Label leakage from aggregated gradients in federated learning. In *IEEE INFOCOM 2024-IEEE Conference on Computer Communications*. IEEE, 151–160.
- [44] Dawood Wasif, Dian Chen, Sindhuja Madabushi, Nithin Alluru, Terrence J Moore, and Jin-Hee Cho. 2025. Empirical analysis of privacy-fairness-accuracy trade-offs in federated learning: a step towards responsible AI. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, Vol. 8. 2666–2677.
- [45] Wencheng Yang, Song Wang, Di Wu, Taotao Cai, Yanming Zhu, Shicheng Wei, Yiyang Zhang, Xu Yang, Zhaohui Tang, and Yan Li. 2025. Deep learning model inversion attacks and defenses: a comprehensive survey. *Artificial Intelligence Review* 58, 8 (2025), 242.
- [46] Xuefei Yin, Yanming Zhu, and Jiankun Hu. 2021. A comprehensive survey of privacy-preserving federated learning: A taxonomy, review, and future directions. *ACM Computing Surveys (CSUR)* 54, 6 (2021), 1–36.
- [47] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. 2020. idlg: Improved deep leakage from gradients. *arXiv preprint arXiv:2001.02610* (2020).
- [48] Ligeng Zhu, Zhijian Liu, and Song Han. 2019. Deep leakage from gradients. *Advances in neural information processing systems* 32 (2019).

## A More on DP

**THEOREM 4 (POST-PROCESSING [9, 42]).** *If  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -DP, then  $G \circ \mathcal{M}$  will also satisfy  $(\epsilon, \delta)$ -DP for any data-independent function  $G$ .*

**Bounding sensitivity via clipping.** To apply DP to aggregated model updates, we must bound how much a single user’s contribution can affect the released statistic. Let  $u_{i,\rho} \in \mathbb{R}^d$  denote user  $i$ ’s (real-valued) update in epoch  $\rho$ . Each user applies norm-bounding (e.g.,  $\ell_2$  clipping) to enforce  $\|u_{i,\rho}\|_2 \leq C$ , yielding  $\tilde{u}_{i,\rho}$ . The enclave then forms the (plaintext) aggregate  $g_\rho = \sum_{i=1}^n \tilde{u}_{i,\rho}$  (or an average).

Under the add/remove notion of user adjacency, the  $\ell_2$ -sensitivity of the sum is at most  $C$  per epoch; under replace-one adjacency it is at most  $2C$ .

**DP mechanism interface.** We parameterize the protocol by a DP mechanism  $\mathcal{M}$  with:

**M.Preprocess**( $u, C$ ): deterministically bounds a user’s update (e.g.,  $\ell_2$  clipping to threshold  $C$ ), returning  $\tilde{u}$ .

**M.ApplyDP**( $\text{st}_{\mathcal{M}}, g$ ): adds calibrated noise to the aggregate  $g$  and updates internal state (e.g., a privacy accountant), returning  $(\tilde{g}, \text{st}'_{\mathcal{M}})$ .

In  $\Pi_{\text{DP}}$ ,  $\mathcal{M}$ .ApplyDP is executed inside the enclave, and only  $\tilde{g}$  is released to the server.

## B Security Proofs

### B.1 Proof of Theorem 1 (Aggregation Security of $\Pi$ )

**DEFINITION 5 (REAL EXECUTION EXPERIMENT ( $H_0$ )).** *Fix a protocol  $\Pi = (\text{Setup}, \text{Enc}, \text{Add}, \text{Dec}, \text{Update})$ . The experiment  $\text{REAL}_{\Pi, \mathcal{A}, \mathcal{Z}}^{\mathcal{G}_{\text{att}}}(1^\lambda)$  is the output bit of  $\mathcal{Z}$  in the following interaction:*

- (1) **Initialization.** *The environment  $\mathcal{Z}(1^\lambda)$  chooses public parameters (e.g.,  $n, R, d, t$ ), a static corruption set  $C \subseteq [n]$ , and inputs for all parties. For each user  $U_i$ ,  $\mathcal{Z}$  specifies an update sequence  $(x_{i,0}, \dots, x_{i,R-1})$  with  $x_{i,\rho} \in \mathbb{Z}_t^d$ . The adversary  $\mathcal{A}$  receives the full internal state/inputs of  $S$  and of corrupted users  $\{U_i\}_{i \in C}$ .*
- (2) **Setup (in the  $\mathcal{G}_{\text{att}}$ -hybrid model).** *The enclave party  $E$  executes  $\text{Setup}(1^\lambda, n, R)$ . Using  $\mathcal{G}_{\text{att}}$ ,  $E$  provisions each user secret key  $k_i$  to  $U_i$  over an attested secure channel, retains enclave secret state  $\text{st}_{\text{enc}}$ , publishes public parameters  $\text{pp}$ , and broadcasts  $W_0$ . The adversary controls the network and sees all messages except those protected by attested secure channels.*
- (3) **Training rounds.** *For each round  $\rho = 0, \dots, R-1$ :*
  - (a) *Each honest user  $U_i$  ( $i \in \mathcal{H}$ ) computes ciphertexts according to  $\Pi$ . For vector updates, encryption is coordinate-wise: for each  $j \in [d]$ ,*

$$c_{i,\rho,j} \leftarrow \text{Enc}(k_i, ts_{\rho,j}, x_{i,\rho}[j]).$$

*Each honest user sends its ciphertexts to  $S$  over the adversary-controlled network. Corrupted users’ outgoing messages are generated arbitrarily by  $\mathcal{A}$ .*

- (b) *The server  $S$  (controlled by  $\mathcal{A}$ ) computes any aggregated ciphertext(s) it wishes (not necessarily honestly) and sends an aggregate ciphertext vector to the enclave  $E$ .*

*We write the server-provided aggregate ciphertext vector as  $\hat{y}_\rho := (\hat{y}_{\rho,1}, \dots, \hat{y}_{\rho,d})$ .*

- (c) *The enclave  $E$  computes the plaintext aggregate coordinate-wise:*

$$\begin{aligned} \text{for each } j \in [d]: \quad & x_{\text{agg},\rho}[j] \leftarrow \text{Dec}(\text{st}_{\text{enc}}, ts_{\rho,j}, \hat{y}_{\rho,j}), \\ & x_{\text{agg},\rho} := (x_{\text{agg},\rho}[1], \dots, x_{\text{agg},\rho}[d]). \end{aligned}$$

*The enclave returns  $x_{\text{agg},\rho}$  to  $S$ . By  $\mathcal{G}_{\text{att}}$ ,  $\mathcal{A}$  does not learn enclave secrets or intermediate state beyond this prescribed output.*

- (d) *The server updates the model via Update and broadcasts  $W_{\rho+1}$  (if included in the transcript).*

- (4) **Output.** *The experiment outputs whatever bit  $\mathcal{Z}$  outputs at the end of the interaction.*

**DEFINITION 6 (HYBRID  $H_1$  (HIDDEN PLAINTEXT SIDE-CHANNEL)).** *The experiment  $\mathcal{H}_1(1^\lambda)$  is identical to  $\mathcal{H}_0(1^\lambda)$ , except that for each round  $\rho = 0, \dots, R-1$  and each honest user  $i \in \mathcal{H}$ , user  $U_i$  additionally sends its plaintext update  $x_{i,\rho} \in \mathbb{Z}_t^d$  to the enclave  $E$  over an attested secure channel provided by  $\mathcal{G}_{\text{att}}$ . The enclave stores these values in an auxiliary internal buffer (or ignores them).*

*All adversary-visible messages and all outputs (in particular, the value released by the enclave to the server each round) are unchanged from  $\mathcal{H}_0$ .*

**LEMMA 1.**  $\mathcal{H}_0 \equiv \mathcal{H}_1$  (perfect indistinguishability).

**PROOF.** The only difference between  $\mathcal{H}_0$  and  $\mathcal{H}_1$  is that in  $\mathcal{H}_1$  each honest user  $U_i$  additionally transmits its plaintext update  $x_{i,\rho}$  to the enclave  $E$  over a  $\mathcal{G}_{\text{att}}$ -attested secure channel in each round  $\rho$ .

By the definition of the  $\mathcal{G}_{\text{att}}$ -hybrid model, communication over such an attested secure channel enjoys confidentiality and authenticity against the adversary: the adversary  $\mathcal{A}$  (and hence the environment  $\mathcal{Z}$ ) does not see the contents of these additional messages and cannot modify them. Moreover, the enclave’s externally visible behavior in  $\mathcal{H}_1$  is unchanged: it runs the same code and releases the same plaintext aggregate value to the server as in  $\mathcal{H}_0$ .

Therefore, the joint distribution of all values visible to  $\mathcal{A}$  and  $\mathcal{Z}$  (public parameters, network messages outside attested secure channels, and all party outputs) is identical in  $\mathcal{H}_0$  and  $\mathcal{H}_1$ . Hence,  $\mathcal{H}_0 \equiv \mathcal{H}_1$ .  $\square$

**DEFINITION 7 (HYBRID  $H_2$  (RE-EXPRESS ENCLAVE OUTPUT AS FULL SUM  $+\delta_\rho$ )).** *The experiment  $\mathcal{H}_2(1^\lambda)$  is identical to  $\mathcal{H}_1(1^\lambda)$ , except that in each round  $\rho = 0, \dots, R-1$  the enclave computes an explicit offset  $\delta_\rho$  and releases an equivalent value expressed as a full sum plus  $\delta_\rho$ .*

*Formally, in round  $\rho$ , let  $y_\rho \in \mathbb{Z}_t^d$  denote the plaintext vector that the enclave would obtain by running the decryption/finalization procedure on the server-provided aggregate ciphertexts (i.e., the value it would have released in  $\mathcal{H}_1$ ). In  $\mathcal{H}_2$ , the enclave computes this same  $y_\rho$  and then:*

- (1) *For each honest user  $i \in \mathcal{H}$ , the enclave reads  $x_{i,\rho}$  from the auxiliary buffer populated in  $\mathcal{H}_1$ .*
- (2) *For each corrupted user  $i \in C$ , let  $x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d$  denote the plaintext update chosen by the adversary/corrupted user for this round.*

(3) *The enclave sets and releases:*

$$\delta_\rho := y_\rho - \sum_{i \in \mathcal{H}} x_{i,\rho} - \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d,$$

$$g_\rho := \left( \sum_{i \in \mathcal{H}} x_{i,\rho} \right) + \left( \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) + \delta_\rho \in \mathbb{Z}_t^d.$$

All other aspects of the experiment are unchanged.

LEMMA 2.  $\mathcal{H}_1 \equiv \mathcal{H}_2$  (perfect indistinguishability).

PROOF. For every round  $\rho$ , by construction,

$$g_\rho = \left( \sum_{i \in \mathcal{H}} x_{i,\rho} \right) + \left( \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) + \delta_\rho = y_\rho.$$

Thus, in every round, the value released by the enclave to the server in  $\mathcal{H}_2$  is exactly the same as the value released in  $\mathcal{H}_1$ . No other adversary-visible message is changed. Therefore the adversary/environment view is identically distributed in  $\mathcal{H}_1$  and  $\mathcal{H}_2$ , and hence  $\mathcal{H}_1 \equiv \mathcal{H}_2$ .  $\square$

DEFINITION 8 (HYBRID  $\mathcal{H}_3$  (HONEST USERS ENCRYPT 0)). *The experiment  $\mathcal{H}_3(1^\lambda)$  is identical to  $\mathcal{H}_2(1^\lambda)$ , except for the following change in each round  $\rho = 0, \dots, R-1$ :*

(1) **Honest-user ciphertexts.** *For each honest user  $i \in \mathcal{H}$  (and each coordinate  $j \in [d]$ ), instead of encrypting the true update entry  $x_{i,\rho}[j]$ , user  $U_i$  encrypts 0:*

$$c_{i,\rho,j} \leftarrow \text{Enc}(k_i, ts_{\rho,j}, 0).$$

*Corrupted users' messages are unchanged.*

(2) **Enclave output compensation.** *Let  $y_\rho^{(0)} \in \mathbb{Z}_t^d$  denote the plaintext vector obtained by the enclave by decrypting the server-provided aggregate ciphertexts in this hybrid:*

$$\text{for each } j \in [d]: \quad y_\rho^{(0)}[j] \leftarrow \text{Dec}(\text{st}_{\text{enc}}, ts_{\rho,j}, \hat{y}_{\rho,j}^{(0)}),$$

$$y_\rho^{(0)} := (y_\rho^{(0)}[1], \dots, y_\rho^{(0)}[d]).$$

*The enclave then reads the honest plaintext updates  $\{x_{i,\rho}\}_{i \in \mathcal{H}}$  from the auxiliary buffer populated via the  $\mathcal{G}_{\text{att}}$ -protected side-channel in  $\mathcal{H}_1$ , and releases to the server*

$$g_\rho := y_\rho^{(0)} + \sum_{i \in \mathcal{H}} x_{i,\rho} \in \mathbb{Z}_t^d.$$

All other aspects of the experiment are unchanged.

We rely on TERSE's privacy notion of *aggregator obliviousness* in the *encrypt-once model* [41, Definition 1].

LEMMA 3. *Assume TERSE is aggregator oblivious in the encrypt-once model [41, Definition 1]. Then  $\mathcal{H}_2 \approx_c \mathcal{H}_3$ .*

PROOF. The only difference between  $\mathcal{H}_2$  and  $\mathcal{H}_3$  is the distribution of ciphertexts sent by honest users: in  $\mathcal{H}_2$  they encrypt the true updates, while in  $\mathcal{H}_3$  they encrypt 0. (The enclave's compensation term in  $\mathcal{H}_3$  is a deterministic post-processing of the same honest plaintext updates  $\{x_{i,\rho}\}_{i \in \mathcal{H}}$ , which are fixed by  $\mathcal{Z}$  and hidden from  $\mathcal{A}$  by  $\mathcal{G}_{\text{att}}$ .)

Since TERSE encryption is applied independently and identically to each coordinate  $j \in [d]$ , a standard hybrid over coordinates reduces the  $d$ -dimensional case to  $d$  independent scalar instances; we therefore prove security for  $d = 1$  without loss of generality.

We proceed via a hybrid argument over *timestamps*. Let

$$T := \{ts_{\rho,j} \mid \rho \in \{0, \dots, R-1\}, j \in [d]\}$$

be the set of timestamps used by the protocol. By construction, timestamps in  $T$  are fresh/non-repeating, hence the encrypt-once condition holds (each user encrypts at most once per timestamp).

Fix an arbitrary ordering of the timestamps in  $T$  and denote it by  $ts^{(1)}, \dots, ts^{(|T|)}$ . Define intermediate hybrids  $\mathcal{G}_0, \dots, \mathcal{G}_{|T|}$  as follows:  $\mathcal{G}_0 := \mathcal{H}_2$  and  $\mathcal{G}_{|T|} := \mathcal{H}_3$ , and for each  $q \in \{0, \dots, |T|\}$ , the hybrid  $\mathcal{G}_q$  is identical to  $\mathcal{H}_2$  except that for every honest user  $i \in \mathcal{H}$  and every timestamp among  $ts^{(1)}, \dots, ts^{(q)}$ , the ciphertext that  $U_i$  sends for that timestamp is an encryption of 0 (rather than an encryption of the true message for that timestamp). All other code (including the enclave's behavior) is as in the corresponding endpoint hybrid.

It suffices to show that for every  $q \in \{0, \dots, |T| - 1\}$ ,

$$\mathcal{G}_q \approx_c \mathcal{G}_{q+1}.$$

Suppose towards a contradiction that there exist a PPT distinguisher  $\mathcal{D}$  and an index  $q^*$  such that  $\mathcal{D}$  distinguishes  $\mathcal{G}_{q^*}$  from  $\mathcal{G}_{q^*+1}$  with non-negligible advantage. We build a PPT adversary  $\mathcal{B}$  that breaks TERSE aggregator obliviousness in the encrypt-once model [41, Definition 1].

**Adversary  $\mathcal{B}$ .**  $\mathcal{B}$  interacts with the TERSE challenger and internally runs  $\mathcal{A}$  and  $\mathcal{Z}$ , thereby simulating the entire execution given to  $\mathcal{D}$ . (In the TERSE game notation, we make explicit the noise/one-time-pad values  $r_{i,ts} \in \mathbb{Z}_t$  that are implicit in calls to  $\text{Enc}(\cdot)$ .)

- *Setup:*  $\mathcal{B}$  receives public parameters from the TERSE challenger and forwards them to the internal simulation.
- *Compromise queries (corrupted users and aggregator for simulation):* For each  $i \in \mathcal{C}$ ,  $\mathcal{B}$  issues  $\text{Compromise}(i)$  to obtain the user secret key and gives it to  $\mathcal{A}$  as required by the corruption model. In addition,  $\mathcal{B}$  issues  $\text{Compromise}(\square)$  once to obtain the aggregator/enclave decryption key  $s'$ , which  $\mathcal{B}$  uses *only internally* to simulate the enclave's decryption procedure and thus the plaintext aggregate outputs delivered to  $\mathcal{A}$ .
- *Encrypting honest users at non-challenge timestamps:*  $\mathcal{B}$  fixes  $ts^* := ts^{(q^*+1)}$  and will use the TERSE challenge query only for timestamp  $ts^*$ . For every honest user  $i \in \mathcal{H}$  and every timestamp  $ts \neq ts^*$ , whenever the simulation requires the ciphertext for user  $i$  at timestamp  $ts$ ,  $\mathcal{B}$  answers it using the TERSE Encrypt oracle:
  - If  $ts \in \{ts^{(1)}, \dots, ts^{(q^*)}\}$ , then  $\mathcal{B}$  queries  $\text{Encrypt}(i, ts, 0, r_{i,ts})$  and uses the returned ciphertext.
  - If  $ts \in \{ts^{(q^*+2)}, \dots, ts^{(|T|)}\}$ , then  $\mathcal{B}$  queries  $\text{Encrypt}(i, ts, m_{i,ts}, r_{i,ts})$  and uses the returned ciphertext, where  $m_{i,ts}$  is the honest message that  $U_i$  would encrypt at timestamp  $ts$  in  $\mathcal{G}_{q^*}$  (i.e., the appropriate scalar coordinate of  $x_{i,\rho}$ ).

Here  $r_{i,ts}$  is sampled as specified by TERSE.

- *Challenge timestamp  $ts^*$ :* When the simulation reaches the point where honest users must produce ciphertexts for timestamp  $ts^*$ ,  $\mathcal{B}$  makes its single TERSE Challenge query with participants  $U := \mathcal{H}$  and time  $ts^*$ .

Let  $i^* \in \mathcal{H}$  be an arbitrary fixed honest user. For each  $i \in U$ , let  $m_{i,ts^*}$  denote the honest message that  $U_i$  would encrypt at  $ts^*$  in  $\mathcal{G}_{q^*}$ .  $\mathcal{B}$  chooses challenge pairs  $(x_{i,ts^*}^0, r_{i,ts^*}^0)$  and  $(x_{i,ts^*}^1, r_{i,ts^*}^1)$  as:

$$\begin{aligned} (x_{i,ts^*}^0, r_{i,ts^*}^0) &:= (m_{i,ts^*}, r_{i,ts^*}) \text{ for all } i \in U, \\ (x_{i,ts^*}^1, r_{i,ts^*}^1) &:= (0, r_{i,ts^*}) \text{ for all } i \in U \setminus \{i^*\}, \\ (x_{i^*,ts^*}^1, r_{i^*,ts^*}^1) &:= \left(0, r_{i^*,ts^*} + \sum_{i \in U} m_{i,ts^*}\right) \pmod{t}. \end{aligned}$$

This ensures the equal-aggregate condition required by [41, Definition 1] in the aggregator-compromised case:

$$\sum_{i \in U} (x_{i,ts^*}^0 + r_{i,ts^*}^0) = \sum_{i \in U} (x_{i,ts^*}^1 + r_{i,ts^*}^1) \text{ in } \mathbb{Z}_t.$$

The challenger samples  $b \in \{0, 1\}$  and returns  $\{c_{i,ts^*}\}_{i \in U}$ ;  $\mathcal{B}$  uses these ciphertexts as the honest users' outgoing ciphertexts for timestamp  $ts^*$ .

- *All other messages:* Corrupted users' outgoing messages are produced by  $\mathcal{A}$  (who knows their keys). The adversarial server behavior and all other parts of the protocol execution are simulated faithfully, using  $s'$  to simulate the enclave's decryption outputs where needed.

**Correctness of the reduction.** By construction, the simulated execution is distributed exactly as  $\mathcal{G}_{q^*}$  if  $b = 0$ , and exactly as  $\mathcal{G}_{q^*+1}$  if  $b = 1$ . Therefore,

$$\mathcal{A}_{\mathcal{B}} = |\Pr[\mathcal{D}(\mathcal{G}_{q^*}) = 1] - \Pr[\mathcal{D}(\mathcal{G}_{q^*+1}) = 1]|,$$

which is non-negligible by assumption, contradicting TERSE aggregator obliviousness.

Hence  $\mathcal{G}_q \approx_c \mathcal{G}_{q+1}$  for all  $q$ , and by transitivity,

$$\mathcal{H}_2 = \mathcal{G}_0 \approx_c \mathcal{G}_{|T|} = \mathcal{H}_3. \quad \square$$

**DEFINITION 9 (HYBRID  $\mathcal{H}_4$  (DEFINE  $\delta_\rho$  EXPLICITLY IN THE ENCRYPT-0 WORLD)).** The experiment  $\mathcal{H}_4(1^\lambda)$  is identical to  $\mathcal{H}_3(1^\lambda)$ , except that in each round  $\rho = 0, \dots, R-1$  the enclave computes an explicit offset  $\delta_\rho$  and releases the aggregate in the ideal-functionality form  $\sum_{i \in [n]} x_{i,\rho} + \delta_\rho$ .

Formally, in round  $\rho$ , let  $y_\rho^{(0)} \in \mathbb{Z}_t^d$  be as computed in  $\mathcal{H}_3$ . The enclave then:

- (1) Reads the honest plaintext updates  $\{x_{i,\rho}\}_{i \in \mathcal{H}}$  from the auxiliary buffer populated via the  $\mathcal{G}_{\text{att}}$ -protected side-channel in  $\mathcal{H}_1$ .
- (2) For each corrupted user  $i \in \mathcal{C}$ , let  $x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d$  denote the plaintext update chosen by the adversary/corrupted user for this round.
- (3) Sets and releases:

$$\begin{aligned} \delta_\rho &:= y_\rho^{(0)} - \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d, \\ g_\rho &:= \left( \sum_{i \in \mathcal{H}} x_{i,\rho} \right) + \left( \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) + \delta_\rho \\ &= \left( \sum_{i \in [n]} x_{i,\rho} \right) + \delta_\rho \in \mathbb{Z}_t^d. \end{aligned}$$

All other aspects of the experiment are unchanged.

**LEMMA 4.**  $\mathcal{H}_3 \equiv \mathcal{H}_4$  (perfect indistinguishability).

**PROOF.** Fix any execution and any round  $\rho$ . In  $\mathcal{H}_3$ , the enclave releases

$$g_\rho^{(3)} := y_\rho^{(0)} + \sum_{i \in \mathcal{H}} x_{i,\rho}.$$

In  $\mathcal{H}_4$ , the enclave sets  $\delta_\rho := y_\rho^{(0)} - \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}}$  and releases

$$\begin{aligned} g_\rho^{(4)} &:= \left( \sum_{i \in \mathcal{H}} x_{i,\rho} \right) + \left( \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) + \delta_\rho \\ &= \left( \sum_{i \in \mathcal{H}} x_{i,\rho} \right) + \left( \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) + \left( y_\rho^{(0)} - \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) \\ &= y_\rho^{(0)} + \sum_{i \in \mathcal{H}} x_{i,\rho} = g_\rho^{(3)}. \end{aligned}$$

Thus, in every round, the value released to the server is identical in  $\mathcal{H}_3$  and  $\mathcal{H}_4$ , and no other adversary-visible message is changed. Hence  $\mathcal{H}_3 \equiv \mathcal{H}_4$ .  $\square$

**DEFINITION 10 (HYBRID  $\mathcal{H}_5$  (IDEAL EXECUTION WITH  $\mathcal{F}_{\text{agg}}$ )).** The experiment  $\mathcal{H}_5(1^\lambda)$  is defined as  $\text{IDEAL}_{\mathcal{F}_{\text{agg}}, \mathcal{S}, \mathcal{Z}}(1^\lambda)$  for the functionality  $\mathcal{F}_{\text{agg}}$  in Figure 2, with simulator  $\mathcal{S}$  specified as follows.

**Simulator  $\mathcal{S}$ .**  $\mathcal{S}$  runs  $\mathcal{A}$  internally and simulates the view of  $\mathcal{A}$  as in  $\mathcal{H}_4$ :

- (1) **Setup simulation.**  $\mathcal{S}$  samples TERSE keys for all users and the enclave/aggregator as in the real setup (these keys are used only for simulation). It simulates the  $\mathcal{G}_{\text{att}}$ -attested key provisioning to honest users and hands the keys of corrupted users to  $\mathcal{A}$ .
- (2) **Ciphertext simulation.** For each round  $\rho$  and timestamps  $ts_{\rho,j}$ ,  $\mathcal{S}$  simulates honest users' ciphertexts as encryptions of 0:

$$c_{i,\rho,j} \leftarrow \text{Enc}(k_i, ts_{\rho,j}, 0) \quad \text{for all } i \in \mathcal{H}, j \in [d].$$

Corrupted users' ciphertexts are whatever  $\mathcal{A}$  produces.

- (3) **Extract  $\delta_\rho$  from the adversary's aggregate ciphertext.** In each round  $\rho$ , after  $\mathcal{A}$  sends an aggregate ciphertext vector  $\hat{y}_\rho^{(0)} := (\hat{y}_{\rho,1}^{(0)}, \dots, \hat{y}_{\rho,d}^{(0)})$ ,  $\mathcal{S}$  decrypts it (coordinate-wise) using the simulated enclave secret key to obtain  $y_\rho^{(0)} \in \mathbb{Z}_t^d$ . Using the corrupted users' plaintext updates  $x_{i,\rho}^{\mathcal{A}}$  (known from corrupted parties' internal state),  $\mathcal{S}$  sets

$$\delta_\rho := y_\rho^{(0)} - \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d,$$

and submits  $\delta_\rho$  to  $\mathcal{F}_{\text{agg}}$  on behalf of the corrupted server.

- (4) **Deliver the ideal aggregate.** Upon receiving  $g_\rho$  from  $\mathcal{F}_{\text{agg}}$ ,  $\mathcal{S}$  forwards  $g_\rho$  to  $\mathcal{A}$  as the enclave's output for round  $\rho$ .

**LEMMA 5.**  $\mathcal{H}_4 \equiv \mathcal{H}_5$  (perfect indistinguishability).

**PROOF.** We compare the adversary-visible transcripts round by round.

In  $\mathcal{H}_4$ , honest users' ciphertexts are encryptions of 0, and the enclave computes  $y_\rho^{(0)}$  by decrypting the adversary-provided aggregate ciphertext  $\hat{y}_\rho^{(0)}$ . It then defines  $\delta_\rho := y_\rho^{(0)} - \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}}$  and

releases

$$g_\rho = \left( \sum_{i \in [n]} x_{i,\rho} \right) + \delta_\rho.$$

In  $\mathcal{H}_5$ , the simulator  $\mathcal{S}$  produces the same distribution of honest ciphertexts (encryptions of 0), decrypts the same adversary-provided aggregate ciphertext to obtain the same  $y_\rho^{(0)}$ , and hence computes the same  $\delta_\rho$ . The functionality  $\mathcal{F}_{\text{agg}}$  then outputs

$$g_\rho := \left( \sum_{i \in [n]} x_{i,\rho} \right) + \delta_\rho,$$

which  $\mathcal{S}$  forwards to  $\mathcal{A}$ . Since all other adversary-visible messages are simulated to match  $\mathcal{H}_4$ , the overall view distributions are identical. Hence  $\mathcal{H}_4 \equiv \mathcal{H}_5$ .  $\square$

**PROOF OF THEOREM 1.** Consider the sequence of experiments  $\mathcal{H}_0, \dots, \mathcal{H}_5$  defined in Definitions 5 to 10. We have:

- $\mathcal{H}_0 \equiv \mathcal{H}_1$  by Lemma 1.
- $\mathcal{H}_1 \equiv \mathcal{H}_2$  by Lemma 2.
- $\mathcal{H}_2 \approx_c \mathcal{H}_3$  by Lemma 3.
- $\mathcal{H}_3 \equiv \mathcal{H}_4$  by Lemma 4.
- $\mathcal{H}_4 \equiv \mathcal{H}_5$  by Lemma 5.

By transitivity of (computational) indistinguishability, it follows that

$$\mathcal{H}_0 \approx_c \mathcal{H}_5,$$

i.e.,

$$\text{REAL}_{\Pi, \mathcal{A}, \mathcal{Z}}^{\mathcal{G}_{\text{att}}} (1^\lambda) \approx_c \text{IDEAL}_{\mathcal{F}_{\text{agg}}, \mathcal{S}, \mathcal{Z}} (1^\lambda).$$

This proves that  $\Pi$  securely realizes  $\mathcal{F}_{\text{agg}}$  in the  $\mathcal{G}_{\text{att}}$ -hybrid model.  $\square$

## B.2 Proof of Theorem 2 (UC Security of $\Pi_{\text{DP}}$ )

**DEFINITION 11 (REAL EXECUTION EXPERIMENT ( $H_0$ )).** Fix a protocol  $\Pi_{\text{DP}} = (\text{Setup}, \text{Enc}, \text{Add}, \text{Dec}, \text{Update}, \text{M. ApplyDP})$ . The experiment  $\text{REAL}_{\Pi_{\text{DP}}, \mathcal{A}, \mathcal{Z}}^{\mathcal{G}_{\text{att}}} (1^\lambda)$  is the output bit of  $\mathcal{Z}$  in the following interaction:

- (1) **Initialization.** The environment  $\mathcal{Z} (1^\lambda)$  chooses public parameters (e.g.,  $n, R, d, t$ ), a static corruption set  $C \subseteq [n]$ , and inputs for all parties. For each user  $U_i$ ,  $\mathcal{Z}$  specifies an update sequence  $(x_{i,0}, \dots, x_{i,R-1})$  with  $x_{i,\rho} \in \mathbb{Z}_t^d$ . The adversary  $\mathcal{A}$  receives the full internal state/inputs of  $\mathcal{S}$  and of corrupted users  $\{U_i\}_{i \in C}$ .
- (2) **Setup (in the  $\mathcal{G}_{\text{att}}$ -hybrid model).** The enclave party  $E$  executes  $\text{Setup}(1^\lambda, n, R)$ . Using  $\mathcal{G}_{\text{att}}$ ,  $E$  provisions each user secret key  $k_i$  to  $U_i$  over an attested secure channel, retains enclave secret state  $\text{st}_{\text{enc}}$ , initializes DP state  $\text{st}_{\mathcal{M}}$ , publishes public parameters  $\text{pp}$ , and broadcasts  $W_0$ . The adversary controls the network and sees all messages except those protected by attested secure channels.
- (3) **Training rounds.** For each round  $\rho = 0, \dots, R-1$ :
  - (a) Each honest user  $U_i$  ( $i \in \mathcal{H}$ ) computes ciphertexts according to  $\Pi_{\text{DP}}$ . For vector updates, encryption is coordinate-wise: for each  $j \in [d]$ ,

$$c_{i,\rho,j} \leftarrow \text{Enc}(k_i, ts_{\rho,j}, x_{i,\rho}[j]).$$

Each honest user sends its ciphertexts to  $\mathcal{S}$  over the adversary-controlled network. Corrupted users' outgoing messages are generated arbitrarily by  $\mathcal{A}$ .

- (b) The server  $\mathcal{S}$  (controlled by  $\mathcal{A}$ ) computes any aggregated ciphertext(s) it wishes (not necessarily honestly) and sends an aggregate ciphertext vector to the enclave  $E$ . We write the server-provided aggregate ciphertext vector as  $\hat{y}_\rho := (\hat{y}_{\rho,1}, \dots, \hat{y}_{\rho,d})$ .
- (c) The enclave  $E$  computes the plaintext aggregate coordinate-wise:

$$\begin{aligned} \text{for each } j \in [d]: \quad & x_{\text{agg},\rho}[j] \leftarrow \text{Dec}(\text{st}_{\text{enc}}, ts_{\rho,j}, \hat{y}_{\rho,j}), \\ & x_{\text{agg},\rho} := (x_{\text{agg},\rho}[1], \dots, x_{\text{agg},\rho}[d]). \end{aligned}$$

It then applies the DP mechanism (updating DP state) and returns the sanitized value:

$$(\tilde{g}_\rho, \text{st}_{\mathcal{M}}) \leftarrow \text{M. ApplyDP}(\text{st}_{\mathcal{M}}, x_{\text{agg},\rho}), \quad \tilde{g}_\rho \in \mathbb{R}^d.$$

The enclave returns  $\tilde{g}_\rho$  to  $\mathcal{S}$ . By  $\mathcal{G}_{\text{att}}$ ,  $\mathcal{A}$  does not learn enclave secrets or intermediate state beyond this prescribed output.

- (d) The server updates the model via  $\text{Update}$  and broadcasts  $W_{\rho+1}$  (if included in the transcript).
- (4) **Output.** The experiment outputs whatever bit  $\mathcal{Z}$  outputs at the end of the interaction.

**DEFINITION 12 (HYBRID  $H_1$  (HIDDEN PLAINTEXT SIDE-CHANNEL)).** The experiment  $\mathcal{H}_1(1^\lambda)$  is identical to  $\mathcal{H}_0(1^\lambda)$ , except that for each round  $\rho = 0, \dots, R-1$  and each honest user  $i \in \mathcal{H}$ , user  $U_i$  additionally sends its plaintext update  $x_{i,\rho} \in \mathbb{Z}_t^d$  to the enclave  $E$  over an attested secure channel provided by  $\mathcal{G}_{\text{att}}$ . The enclave stores these values in an auxiliary internal buffer (or ignores them).

All adversary-visible messages and all outputs (in particular, the value released by the enclave to the server each round) are unchanged from  $\mathcal{H}_0$ .

**LEMMA 6.**  $\mathcal{H}_0 \equiv \mathcal{H}_1$  (perfect indistinguishability).

**PROOF.** Identical to the proof of Lemma 1: the only change is an additional message sent from each honest user to the enclave over a  $\mathcal{G}_{\text{att}}$ -attested secure channel, which is confidential and authenticated; the enclave's adversary-visible output remains the same.  $\square$

**DEFINITION 13 (HYBRID  $H_2$  (RE-EXPRESS ENCLAVE INPUT AS FULL SUM  $+\delta_\rho$ )).** The experiment  $\mathcal{H}_2(1^\lambda)$  is identical to  $\mathcal{H}_1(1^\lambda)$ , except that in each round  $\rho = 0, \dots, R-1$  the enclave computes an explicit offset  $\delta_\rho$  and then applies DP to an equivalent value expressed as a full sum plus  $\delta_\rho$ .

Formally, in round  $\rho$ , let  $y_\rho \in \mathbb{Z}_t^d$  denote the plaintext vector that the enclave would obtain by running the decryption/finalization procedure on the server-provided aggregate ciphertexts (i.e., the unsanitized value it would pass to  $\text{M. ApplyDP}$  in  $\mathcal{H}_1$ ). In  $\mathcal{H}_2$ , the enclave computes this same  $y_\rho$  and then:

- (1) For each honest user  $i \in \mathcal{H}$ , the enclave reads  $x_{i,\rho}$  from the auxiliary buffer populated in  $\mathcal{H}_1$ .
- (2) For each corrupted user  $i \in C$ , let  $x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d$  denote the plaintext update chosen by the adversary/corrupted user for this round.

(3) *The enclave sets*

$$\delta_\rho := y_\rho - \sum_{i \in \mathcal{H}} x_{i,\rho} - \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d,$$

$$g_\rho := \left( \sum_{i \in \mathcal{H}} x_{i,\rho} \right) + \left( \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) + \delta_\rho \in \mathbb{Z}_t^d,$$

and releases to the server the DP-sanitized value

$$(\tilde{g}_\rho, \text{st}_M) \leftarrow \text{M.ApplyDP}(\text{st}_M, g_\rho), \quad \tilde{g}_\rho \in \mathbb{R}^d.$$

All other aspects of the experiment are unchanged.

LEMMA 7.  $\mathcal{H}_1 \equiv \mathcal{H}_2$  (perfect indistinguishability).

PROOF. For every round  $\rho$ , by construction,

$$g_\rho = \left( \sum_{i \in \mathcal{H}} x_{i,\rho} \right) + \left( \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) + \delta_\rho = y_\rho.$$

Thus, in every round,  $\text{M.ApplyDP}$  is invoked on the *same* input in  $\mathcal{H}_2$  as it is in  $\mathcal{H}_1$ , and the DP state  $\text{st}_M$  evolves identically. No other adversary-visible message is changed. Therefore the adversary/environment view is identically distributed in  $\mathcal{H}_1$  and  $\mathcal{H}_2$ , and hence  $\mathcal{H}_1 \equiv \mathcal{H}_2$ .  $\square$

DEFINITION 14 (HYBRID  $\mathcal{H}_3$  (HONEST USERS ENCRYPT 0)). *The experiment  $\mathcal{H}_3(1^\lambda)$  is identical to  $\mathcal{H}_2(1^\lambda)$ , except for the following change in each round  $\rho = 0, \dots, R-1$ :*

- (1) **Honest-user ciphertexts.** For each honest user  $i \in \mathcal{H}$  (and each coordinate  $j \in [d]$ ), instead of encrypting the true update entry  $x_{i,\rho}[j]$ , user  $U_i$  encrypts 0:

$$c_{i,\rho,j} \leftarrow \text{Enc}(k_i, ts_{\rho,j}, 0).$$

Corrupted users' messages are unchanged.

- (2) **Enclave output compensation and DP release.** Let  $y_\rho^{(0)} \in \mathbb{Z}_t^d$  denote the plaintext vector obtained by the enclave by decrypting the server-provided aggregate ciphertexts in this hybrid:

$$\text{for each } j \in [d]: \quad y_\rho^{(0)}[j] \leftarrow \text{Dec}(\text{st}_{\text{enc}}, ts_{\rho,j}, \hat{y}_{\rho,j}^{(0)}),$$

$$y_\rho^{(0)} := (y_\rho^{(0)}[1], \dots, y_\rho^{(0)}[d]).$$

The enclave then reads the honest plaintext updates  $\{x_{i,\rho}\}_{i \in \mathcal{H}}$  from the auxiliary buffer populated via the  $\mathcal{G}_{\text{att}}$ -protected side-channel in  $\mathcal{H}_1$ , defines the compensated aggregate

$$g_\rho := y_\rho^{(0)} + \sum_{i \in \mathcal{H}} x_{i,\rho} \in \mathbb{Z}_t^d,$$

and releases the DP-sanitized value

$$(\tilde{g}_\rho, \text{st}_M) \leftarrow \text{M.ApplyDP}(\text{st}_M, g_\rho), \quad \tilde{g}_\rho \in \mathbb{R}^d.$$

All other aspects of the experiment are unchanged.

We rely on TERSE's privacy notion of *aggregator obliviousness in the encrypt-once model* [41, Definition 1].

LEMMA 8. Assume TERSE is aggregator oblivious in the encrypt-once model [41, Definition 1]. Then  $\mathcal{H}_2 \approx_c \mathcal{H}_3$ .

PROOF. The only difference between  $\mathcal{H}_2$  and  $\mathcal{H}_3$  is the distribution of ciphertexts sent by honest users: in  $\mathcal{H}_2$  they encrypt the true updates, while in  $\mathcal{H}_3$  they encrypt 0. The enclave's additional steps in  $\mathcal{H}_3$  (adding the honest plaintext updates from the  $\mathcal{G}_{\text{att}}$ -protected side-channel and then applying  $\text{M.ApplyDP}$ ) are a post-processing of values fixed by  $\mathcal{Z}$  and hidden from  $\mathcal{A}$  by  $\mathcal{G}_{\text{att}}$ .

Since TERSE encryption is applied independently and identically to each coordinate  $j \in [d]$ , a standard hybrid over coordinates reduces the  $d$ -dimensional case to  $d$  independent scalar instances; we therefore prove security for  $d = 1$  without loss of generality.

We proceed via a hybrid argument over *timestamps*. Let

$$T := \{ts_{\rho,j} \mid \rho \in \{0, \dots, R-1\}, j \in [d]\}$$

be the set of timestamps used by the protocol. By construction, timestamps in  $T$  are fresh/non-repeating, hence the encrypt-once condition holds (each user encrypts at most once per timestamp).

Fix an arbitrary ordering of the timestamps in  $T$  and denote it by  $ts^{(1)}, \dots, ts^{(|T|)}$ . Define intermediate hybrids  $\mathcal{G}_0, \dots, \mathcal{G}_{|T|}$  as follows:  $\mathcal{G}_0 := \mathcal{H}_2$  and  $\mathcal{G}_{|T|} := \mathcal{H}_3$ , and for each  $q \in \{0, \dots, |T|\}$ , the hybrid  $\mathcal{G}_q$  is identical to  $\mathcal{H}_2$  except that for every honest user  $i \in \mathcal{H}$  and every timestamp among  $ts^{(1)}, \dots, ts^{(q)}$ , the ciphertext that  $U_i$  sends for that timestamp is an encryption of 0 (rather than an encryption of the true message for that timestamp). All other code (including the enclave's behavior) is as in the corresponding endpoint hybrid.

It suffices to show that for every  $q \in \{0, \dots, |T| - 1\}$ ,

$$\mathcal{G}_q \approx_c \mathcal{G}_{q+1}.$$

Suppose towards a contradiction that there exist a PPT distinguisher  $\mathcal{D}$  and an index  $q^*$  such that  $\mathcal{D}$  distinguishes  $\mathcal{G}_{q^*}$  from  $\mathcal{G}_{q^*+1}$  with non-negligible advantage. We build a PPT adversary  $\mathcal{B}$  that breaks TERSE aggregator obliviousness in the encrypt-once model [41, Definition 1].

**Adversary  $\mathcal{B}$ .**  $\mathcal{B}$  interacts with the TERSE challenger and internally runs  $\mathcal{A}$  and  $\mathcal{Z}$ , thereby simulating the entire execution given to  $\mathcal{D}$ . (In the TERSE game notation, we make explicit the noise/one-time-pad values  $r_{i,ts} \in \mathbb{Z}_t$  that are implicit in calls to  $\text{Enc}(\cdot)$ .) In addition,  $\mathcal{B}$  maintains an internal DP state  $\text{st}_M$  and simulates the sanitized releases by running  $\text{M.ApplyDP}$  on the same plaintext inputs as the enclave would.

- *Setup:*  $\mathcal{B}$  receives public parameters from the TERSE challenger and forwards them to the internal simulation.
- *Compromise queries (corrupted users and aggregator for simulation):* For each  $i \in \mathcal{C}$ ,  $\mathcal{B}$  issues  $\text{Compromise}(i)$  to obtain the user secret key and gives it to  $\mathcal{A}$  as required by the corruption model. In addition,  $\mathcal{B}$  issues  $\text{Compromise}(\square)$  once to obtain the aggregator/enclave decryption key  $s'$ , which  $\mathcal{B}$  uses *only internally* to simulate the enclave's decryption procedure.
- *Encrypting honest users at non-challenge timestamps:*  $\mathcal{B}$  fixes  $ts^* := ts^{(q^*+1)}$  and will use the TERSE challenge query only for timestamp  $ts^*$ . For every honest user  $i \in \mathcal{H}$  and every timestamp  $ts \neq ts^*$ , whenever the simulation requires the ciphertext for user  $i$  at timestamp  $ts$ ,  $\mathcal{B}$  answers it using the TERSE Encrypt oracle:
  - If  $ts \in \{ts^{(1)}, \dots, ts^{(q^*)}\}$ , then  $\mathcal{B}$  queries  $\text{Encrypt}(i, ts, 0, r_{i,ts})$  and uses the returned ciphertext.

- If  $ts \in \{ts^{(q^*+2)}, \dots, ts^{(|T|)}\}$ , then  $\mathcal{B}$  queries  $\text{Encrypt}(i, ts, m_{i,ts}, r_{i,ts})$  and uses the returned ciphertext, where  $m_{i,ts}$  is the honest message that  $U_i$  would encrypt at timestamp  $ts$  in  $\mathcal{G}_{q^*}$  (i.e., the appropriate scalar coordinate of  $x_{i,\rho}$ ).

Here  $r_{i,ts}$  is sampled as specified by TERSE.

- **Challenge timestamp  $ts^*$ :** When the simulation reaches the point where honest users must produce ciphertexts for timestamp  $ts^*$ ,  $\mathcal{B}$  makes its single TERSE Challenge query with participants  $U := \mathcal{H}$  and time  $ts^*$ .

Let  $i^* \in \mathcal{H}$  be an arbitrary fixed honest user. For each  $i \in U$ , let  $m_{i,ts^*}$  denote the honest message that  $U_i$  would encrypt at  $ts^*$  in  $\mathcal{G}_{q^*}$ .  $\mathcal{B}$  chooses challenge pairs  $(x_{i,ts^*}^0, r_{i,ts^*}^0)$  and  $(x_{i,ts^*}^1, r_{i,ts^*}^1)$  as:

$$\text{for all } i \in U : (x_{i,ts^*}^0, r_{i,ts^*}^0) = (m_{i,ts^*}, r_{i,ts^*}),$$

$$\text{for all } i \in U \setminus \{i^*\} : (x_{i,ts^*}^1, r_{i,ts^*}^1) = (0, r_{i,ts^*}),$$

$$(x_{i^*,ts^*}^1, r_{i^*,ts^*}^1) = \left(0, r_{i^*,ts^*} + \sum_{i \in U} m_{i,ts^*}\right) \pmod{t}.$$

This ensures the equal-aggregate condition required by [41, Definition 1] in the aggregator-compromised case:

$$\sum_{i \in U} (x_{i,ts^*}^0 + r_{i,ts^*}^0) = \sum_{i \in U} (x_{i,ts^*}^1 + r_{i,ts^*}^1) \text{ in } \mathbb{Z}_t.$$

The challenger samples  $b \in \{0, 1\}$  and returns  $\{c_{i,ts^*}\}_{i \in U}$ ;  $\mathcal{B}$  uses these ciphertexts as the honest users' outgoing ciphertexts for timestamp  $ts^*$ .

- **All other messages:** Corrupted users' outgoing messages are produced by  $\mathcal{A}$ . The adversarial server behavior and all other parts of the protocol execution are simulated faithfully, using  $s'$  to simulate the enclave's decryption outputs where needed. Whenever the simulation reaches the enclave's release step in some round  $\rho$ ,  $\mathcal{B}$  computes the corresponding plaintext input  $g_\rho$  for  $\text{M.ApplyDP}$  (as prescribed by the hybrid being simulated), runs

$$(\tilde{g}_\rho, \text{st}_M) \leftarrow \text{M.ApplyDP}(\text{st}_M, g_\rho),$$

and delivers  $\tilde{g}_\rho$  to  $\mathcal{A}$ .

**Correctness of the reduction.** By construction, the simulated execution is distributed exactly as  $\mathcal{G}_{q^*}$  if  $b = 0$ , and exactly as  $\mathcal{G}_{q^*+1}$  if  $b = 1$ . Therefore,

$$\mathcal{A}_{\mathcal{B}} = |\Pr[\mathcal{D}(\mathcal{G}_{q^*}) = 1] - \Pr[\mathcal{D}(\mathcal{G}_{q^*+1}) = 1]|,$$

which is non-negligible by assumption, contradicting TERSE aggregator obliviousness.

Hence  $\mathcal{G}_q \approx_c \mathcal{G}_{q+1}$  for all  $q$ , and by transitivity,

$$\mathcal{H}_2 = \mathcal{G}_0 \approx_c \mathcal{G}_{|T|} = \mathcal{H}_3. \quad \square$$

**DEFINITION 15 (HYBRID  $\mathcal{H}_4$  (DEFINE  $\delta_\rho$  EXPLICITLY IN THE ENCRYPT-0 WORLD)).** The experiment  $\mathcal{H}_4(1^\lambda)$  is identical to  $\mathcal{H}_3(1^\lambda)$ , except that in each round  $\rho = 0, \dots, R-1$  the enclave computes an explicit offset  $\delta_\rho$  and applies DP to the aggregate in the ideal-functionality form  $\sum_{i \in [n]} x_{i,\rho} + \delta_\rho$ .

Formally, in round  $\rho$ , let  $y_\rho^{(0)} \in \mathbb{Z}_t^d$  be as computed in  $\mathcal{H}_3$ . The enclave then:

- (1) Reads the honest plaintext updates  $\{x_{i,\rho}\}_{i \in \mathcal{H}}$  from the auxiliary buffer populated via the  $\mathcal{G}_{\text{att}}$ -protected side-channel in  $\mathcal{H}_1$ .
- (2) For each corrupted user  $i \in \mathcal{C}$ , let  $x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d$  denote the plaintext update chosen by the adversary/corrupted user for this round.
- (3) Sets

$$\delta_\rho := y_\rho^{(0)} - \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d,$$

$$\begin{aligned} g_\rho &:= \left( \sum_{i \in \mathcal{H}} x_{i,\rho} \right) + \left( \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) + \delta_\rho \\ &= \left( \sum_{i \in [n]} x_{i,\rho} \right) + \delta_\rho \in \mathbb{Z}_t^d, \end{aligned}$$

and releases the DP-sanitized value

$$(\tilde{g}_\rho, \text{st}_M) \leftarrow \text{M.ApplyDP}(\text{st}_M, g_\rho).$$

All other aspects of the experiment are unchanged.

**LEMMA 9.**  $\mathcal{H}_3 \equiv \mathcal{H}_4$  (perfect indistinguishability).

**PROOF.** Fix any execution and any round  $\rho$ . In  $\mathcal{H}_3$ , the enclave defines

$$g_\rho^{(3)} := y_\rho^{(0)} + \sum_{i \in \mathcal{H}} x_{i,\rho}.$$

In  $\mathcal{H}_4$ , the enclave sets  $\delta_\rho := y_\rho^{(0)} - \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}}$  and defines

$$\begin{aligned} g_\rho^{(4)} &:= \left( \sum_{i \in \mathcal{H}} x_{i,\rho} \right) + \left( \sum_{i \in \mathcal{C}} x_{i,\rho}^{\mathcal{A}} \right) + \delta_\rho \\ &= y_\rho^{(0)} + \sum_{i \in \mathcal{H}} x_{i,\rho} = g_\rho^{(3)}. \end{aligned}$$

Therefore  $\text{M.ApplyDP}$  is invoked on the same input in both hybrids in round  $\rho$ , and thus the released  $\tilde{g}_\rho$  (and the updated DP state) are identically distributed. No other adversary-visible message is changed. Hence  $\mathcal{H}_3 \equiv \mathcal{H}_4$ .  $\square$

**DEFINITION 16 (HYBRID  $\mathcal{H}_5$  (IDEAL EXECUTION WITH  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$ )).** The experiment  $\mathcal{H}_5(1^\lambda)$  is defined as  $\text{IDEAL}_{\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}, \mathcal{S}, \mathcal{Z}}(1^\lambda)$  for the functionality  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  in Figure 3, with simulator  $\mathcal{S}$  specified as follows.

**Simulator  $\mathcal{S}$ .**  $\mathcal{S}$  runs  $\mathcal{A}$  internally and simulates the view of  $\mathcal{A}$  as in  $\mathcal{H}_4$ :

- (1) **Setup simulation.**  $\mathcal{S}$  samples TERSE keys for all users and the enclave/aggregator as in the real setup (these keys are used only for simulation). It simulates the  $\mathcal{G}_{\text{att}}$ -attested key provisioning to honest users and hands the keys of corrupted users to  $\mathcal{A}$ .
- (2) **Ciphertext simulation.** For each round  $\rho$  and timestamps  $ts_{\rho,j}$ ,  $\mathcal{S}$  simulates honest users' ciphertexts as encryptions of 0:

$$c_{i,\rho,j} \leftarrow \text{Enc}(k_i, ts_{\rho,j}, 0) \quad \text{for all } i \in \mathcal{H}, j \in [d].$$

Corrupted users' ciphertexts are whatever  $\mathcal{A}$  produces.

- (3) **Extract  $\delta_\rho$  from the adversary's aggregate ciphertext.** In each round  $\rho$ , after  $\mathcal{A}$  sends an aggregate ciphertext vector  $\hat{y}_\rho^{(0)} := (\hat{y}_{\rho,1}^{(0)}, \dots, \hat{y}_{\rho,d}^{(0)})$ ,  $\mathcal{S}$  decrypts it (coordinate-wise) using the simulated enclave secret key to obtain  $y_\rho^{(0)} \in \mathbb{Z}_t^d$ . Using the corrupted users' plaintext updates  $x_{i,\rho}^{\mathcal{A}}$  (known from corrupted parties' internal state),  $\mathcal{S}$  sets

$$\delta_\rho := y_\rho^{(0)} - \sum_{i \in C} x_{i,\rho}^{\mathcal{A}} \in \mathbb{Z}_t^d,$$

and submits  $\delta_\rho$  to  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  on behalf of the corrupted server.

- (4) **Deliver the ideal DP-sanitized aggregate.** Upon receiving  $\tilde{g}_\rho$  from  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$ ,  $\mathcal{S}$  forwards  $\tilde{g}_\rho$  to  $\mathcal{A}$  as the enclave's output for round  $\rho$ .

LEMMA 10.  $\mathcal{H}_4 \equiv \mathcal{H}_5$  (perfect indistinguishability).

PROOF. We compare the adversary-visible transcripts round by round.

In  $\mathcal{H}_4$ , honest users' ciphertexts are encryptions of 0, and the enclave computes  $y_\rho^{(0)}$  by decrypting the adversary-provided aggregate ciphertext  $\hat{y}_\rho^{(0)}$ . It then defines  $\delta_\rho := y_\rho^{(0)} - \sum_{i \in C} x_{i,\rho}^{\mathcal{A}}$  and thus

$$g_\rho = \left( \sum_{i \in [n]} x_{i,\rho} \right) + \delta_\rho,$$

and releases  $(\tilde{g}_\rho, st_{\mathcal{M}}) \leftarrow \text{M.ApplyDP}(st_{\mathcal{M}}, g_\rho)$ .

In  $\mathcal{H}_5$ , the simulator  $\mathcal{S}$  produces the same distribution of honest ciphertexts (encryptions of 0), decrypts the same adversary-provided aggregate ciphertext to obtain the same  $y_\rho^{(0)}$ , and hence computes the same  $\delta_\rho$ , which it submits to  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$ . The functionality  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  then computes the same  $g_\rho := (\sum_{i \in [n]} x_{i,\rho}) + \delta_\rho$  and releases the corresponding DP-sanitized value  $\tilde{g}_\rho$ , which  $\mathcal{S}$  forwards to  $\mathcal{A}$ . Since all other adversary-visible messages are simulated to match  $\mathcal{H}_4$ , the overall view distributions are identical. Hence  $\mathcal{H}_4 \equiv \mathcal{H}_5$ .  $\square$

PROOF OF THEOREM 2. Consider the sequence of experiments  $\mathcal{H}_0, \dots, \mathcal{H}_5$  defined in Definitions 11 to 16. We have:

- $\mathcal{H}_0 \equiv \mathcal{H}_1$  by Lemma 6.
- $\mathcal{H}_1 \equiv \mathcal{H}_2$  by Lemma 7.
- $\mathcal{H}_2 \approx_c \mathcal{H}_3$  by Lemma 8.
- $\mathcal{H}_3 \equiv \mathcal{H}_4$  by Lemma 9.
- $\mathcal{H}_4 \equiv \mathcal{H}_5$  by Lemma 10.

By transitivity of (computational) indistinguishability, it follows that

$$\mathcal{H}_0 \approx_c \mathcal{H}_5,$$

i.e.,

$$\text{REAL}_{\Pi_{\text{DP}}, \mathcal{A}, \mathcal{Z}}^{\mathcal{G}_{\text{att}}} (1^\lambda) \approx_c \text{IDEAL}_{\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}, \mathcal{S}, \mathcal{Z}} (1^\lambda).$$

This proves that  $\Pi_{\text{DP}}$  securely realizes  $\mathcal{F}_{\text{agg-DP}}^{\mathcal{M}}$  in the  $\mathcal{G}_{\text{att}}$ -hybrid model.  $\square$

### B.3 Proof of Theorem 3 (DP Guarantee or $\Pi_{\text{DP}}$ )

PROOF. Let  $D$  and  $D'$  be user-adjacent datasets (i.e., they differ in exactly one user's full contribution across all rounds). Let  $\text{Rel}(D)$  denote the sequence of values released by the enclave in  $\Pi_{\text{DP}}$  across all rounds. By assumption, there exists a deterministic function  $F$  (collecting the enclave's true aggregates across rounds) and input-independent randomness  $U$  such that

$$\text{Rel}(D) = \mathcal{M}(F(D); U) \quad \text{and} \quad \text{Rel}(D') = \mathcal{M}(F(D'); U).$$

Since  $\mathcal{M}$  is  $(\epsilon, \delta)$ -DP under user-level adjacency, for any measurable set of transcripts  $S$ ,

$$\begin{aligned} \Pr[\text{Rel}(D) \in S] &= \Pr[\mathcal{M}(F(D); U) \in S] \\ &\leq e^\epsilon \Pr[\mathcal{M}(F(D'); U) \in S] + \delta \\ &= e^\epsilon \Pr[\text{Rel}(D') \in S] + \delta. \end{aligned}$$

Thus, the enclave's released transcript is  $(\epsilon, \delta)$ -DP at user level.

Finally, the adversary's full view (e.g., any derived model updates or final model) is a (possibly randomized) function of  $\text{Rel}(D)$  together with data-independent information (public parameters and the adversary's own randomness). By closure of DP under post-processing, any such derived view remains  $(\epsilon, \delta)$ -DP. This proves the claim.  $\square$