# FedML-HE: An Efficient Homomorphic-Encryption-Based Privacy-Preserving Federated Learning System

**Weizhao Jin**[*]
weizhaoj@usc.edu
University of Southern California

**Yuhang Yao**[*]
yuhangya@andrew.cmu.edu
Carnegie Mellon University

**Shanshan Han**
shanshan.han@uci.edu
University of California Irvine

**Carlee Joe-Wong**
cjoewong@andrew.cmu.edu
Carnegie Mellon University

**Srivatsan Ravi**
sravi@isi.edu
University of Southern California

**Salman Avestimehr**
avestimehr@fedml.ai
FedML Inc.

**Chaoyang He**[†]
ch@fedml.ai
FedML Inc.

## Abstract

Federated Learning (FL) enables machine learning model training on distributed edge devices by aggregating local model updates rather than local data. However, privacy concerns arise as the FL server's access to local model updates can potentially reveal sensitive personal information by performing attacks like gradient inversion recovery. To address these concerns, privacy-preserving methods, such as Homomorphic Encryption (HE)-based approaches, have been proposed. Despite HE's post-quantum security advantages, its applications suffer from impractical overheads. In this paper, we present FedML-HE, *the first practical system for efficient HE-based secure federated aggregation* that provides a user/device-friendly deployment platform. FedML-HE utilizes a novel universal overhead optimization scheme, significantly reducing both computation and communication overheads during deployment while providing customizable privacy guarantees. Our optimized system demonstrates considerable overhead reduction, particularly for large models (e.g., ~10x reduction for HE-federated training of ResNet-50 and ~40x reduction for BERT), demonstrating the potential for scalable HE-based FL deployment.
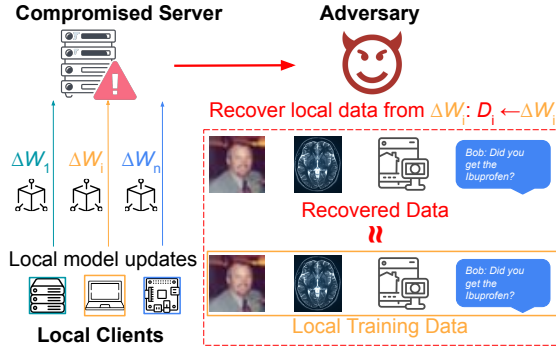
## 1 Introduction

Federated learning enables model training on decentralized edge devices without aggregating data to a central location by aggregating locally-trained models, which should ideally ensure data privacy. Privacy preservation in standard federated learning systems depends on straightforward aggregation functions, such as federated averaging (FedAvg) [33]. Rather than uploading data to a central server for training, clients process their local data and share the local models with the server. However, this approach exposes individual local models to the aggregation server and potentially other parties within the system. Since local models are typically trained on limited-size datasets, sharing local models in plaintext during aggregation poses a privacy vulnerability. Various attacks have been proposed to extract private training data information and compromise individual privacy, including membership inference attacks [36, 50, 53] and sensitive user data recovery attacks [7, 12, 21, 24, 58]. As illustrated in Figure 1, a compromised server can conduct attacks to recover personal data used for local training, which can expose sensitive personal information, e.g., if FL is used on identifying information, medical records, smart home (e.g., security camera) deployment, or private messages.

Existing defense mechanisms to prevent the FL server from learning plaintext local models include Multi-Party Computation (MPC) secure aggregation protocols [8, 45] and noise-based Differential Privacy (DP) solutions [10, 49]. MPC protocols mask private inputs, while DP solutions add privacy noise to original inputs. However, these methods have limitations: MPC requires additional interactive synchronization steps and is susceptible to

---

[*]Equal contribution
[†]Corresponding author

**Figure 1: FL Inversion Attacks**: an adversary server can recover local training data from unprotected local model updates.

failures like client dropout, necessitating careful engineering during deployment, while DP can result in model performance degradation due to privacy noises. Homomorphic Encryption (HE) [9, 11, 16, 18, 39] offers a robust post-quantum solution that protects local models against attacks and ensures local model privacy while maintaining adequate utility. HE-based federated aggregation encrypts local models during aggregation and performs model aggregation over ciphertexts. This approach can enable secure federated learning deployments and has been adopted by several FL systems [15, 26, 38, 57] and a few domain-specific applications have been implemented [28, 47, 56].

| Features | IBMFL | Nvidia FLARE | FedML-HE |
|---|---|---|---|
| Homomorphic Encryption | ✓ | ✓ | ✓ |
| Key Management | ◯ | ✓ | ✓ |
| Efficient Encrypted Computation | ✗ | ◯ | ✓ |
| HE Multi-Party Functionalities | ✗ | ✗ | ✓ |

**Table 1: Comparison with Existing HE-Based FL Systems**: ◯ implies limited support.

Despite advances in the Homomorphic Encryption research community to enhance HE performance in theory and practice, HE remains a powerful but complex cryptographic foundation with impractical overheads for most real-world applications. Prior HE solutions mainly employ existing generic HE methods without sufficient optimization for large-scale privacy-preserving federated learning deployment. The scalability of encrypted computation and communication during federated training then becomes a bottleneck, restricting its feasibility for real-world edge computing scenarios. This HE overhead limitation is rather noticeable (commonly ~15x

increase in both computation and communication [19]) when training large federated models across resource-constrained mobile devices, where encrypted computing and communication of large models might take considerably longer than the actual model training.

To address these challenges, we propose FedML-HE, an efficient Homomorphic Encryption-based privacy-preserving FL system with a universal optimization scheme, designed for practical deployment across distributed edge devices[1]. Our system significantly reduces overheads, enabling HE-based federated learning to be more accessible and efficient in real-world scenarios.

**Key contributions**:

- We implement our HE-based FL system on a user/device-friendly deployment platform that includes key management, system profiling, grouping, and edge binding. The platform also enables easy monitoring of HE overhead distribution across edge devices.
- We benchmark our HE-FL framework under various aspects of FL training, which allows us to pinpoint the HE bottlenecks. Our comparison with other HE-based solutions reveals that our vanilla implementation already has the lowest computational overhead while offering more functionalities (see Table 1 and Table 5).
- We propose a universal HE-FL optimization scheme, **Parameter Efficiency** x **Parameter Selection**, that minimizes the size of model updates for encrypted computation while preserving privacy guarantees. This approach mitigates both computational and communication overhead increases from HE operations. With optimization, our framework achieves significant overhead reduction, particularly for large models (e.g., ~10x reduction for HE-federated training ResNet-50 and ~40x reduction for BERT), demonstrating the potential for practical HE-based FL deployments.

## 2 Preliminaries And Related Work

### 2.1 Federated Learning

Federated learning is first proposed in [33], which builds distributed machine learning models while keeping personal data on clients. Instead of uploading data to the server for centralized training, clients process their local data and share updated local models with the server. Model parameters from a large population of clients are aggregated by the server and combined to create an improved global model.

The FedAvg [33] is commonly used on the server to combine client updates and produce a new global model. At each round, a global model $\mathbf{W}_{glob}$ is sent to $N$ client devices. Each client $i$ performs gradient descent on its

---

[1]Our code is open-source at https://github.com/FedML-AI/FedML.

local data with $E$ local iterations to update the model $\mathbf{W}_i$. The server then does a weighted aggregation of the local models to obtain a new global model, $\mathbf{W}_{\text{glob}} = \sum_{i=1}^{N} \alpha_i \mathbf{W}_i$, where $\alpha_i$ is the weighting factor for client $i$.

Typically, the aggregation runs using plaintext model parameters through a central server (in some cases, via a decentralized protocol), giving the server visibility of each local client's model in plaintext.

## 2.2 Homomorphic Encryption

- HE.$KeyGen(\lambda)$: given the security parameter $\lambda$, the key generation algorithm outputs a key pair $(pk, sk)$ and the related cryptographic context.
- HE.$Enc(pk, m)$:the encryption algorithm takes in $pk$ and a plaintext message $m$, then outputs the ciphertext $c$.
- HE.$Eval(c, f)$:the encrypted evaluation algorithm takes in a ciphertext message $c$ and a function $f$, then outputs the computation result $c'$.
- HE.$Dec(sk, c')$:the encryption algorithm takes in $sk$ and a ciphertext message $c'$, then outputs the plaintext $m'$.

**Figure 2:** General Scheme of Homomorphic Encryption

Homomorphic Encryption is a cryptographic primitive that allows computation to be performed on encrypted data without revealing the underlying plaintext. It usually serves as a foundation for privacy-preserving outsourcing computing models. HE has generally four algorithms ($KeyGen$, $Enc$, $Eval$, $Dec$) as defined in Figure 2. The fundamental concept is to encrypt data prior to computation, perform the computation on the encrypted data without decryption, and then decrypt the resulting ciphertext to obtain the final plaintext.

Since FL model parameters are usually not integers, our method is built on the Cheon-Kim-Kim-Song (CKKS) scheme [11], a (leveled) HE variant that can work with approximate numbers.

## 2.3 Previous Work

**Existing Privacy Attacks On FL** Threats and attacks on privacy in the domain of Federated Learning have been studied in recent years [35]. General FL privacy attacks can be categorized into two types: inference attacks [36, 50, 53] and data leakage/reconstruction [7, 12, 24]. Attacks are usually carried out on the models to retrieve certain properties of data providers or even reconstruct the data in the training datasets. With direct access to more fine-grained local models trained on a smaller dataset [53], the adversary can have a higher chance of a successful attack. Moreover, further attacks can be performed using GAN-based attacks to even fully
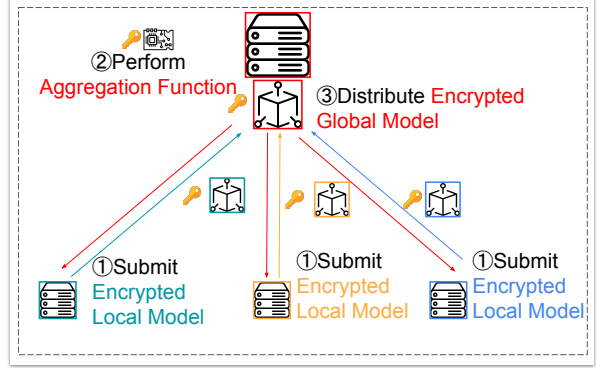


**Figure 3:** FedML-HE System Illustration: local models are encrypted upon federated aggregation and the server acts as a computing service without access to local models.

recover the original data [24]. The majority of the privacy attacks can be traced back to the direct exposure of plaintext accesses to local models to other parties (usually the aggregation server) within the system.

**Existing Non-HE Defense Mechanism** Compared to other existing solutions providing privacy protection in FL, HE is non-interactive and dropout-resilient (vs. general Multi-Party Computation protocols [8, 45]); it introduces negligible model performance degradation (vs. noise-based Differential Privacy solutions [10, 49]).

**Existing HE-based FL Work** Existing HE-based FL work either apply restricted HE schemes (e.g. additive scheme Paillier) [17, 27, 57] without extensibility to further FL aggregation functions as well as sufficient performance and security guarantee (due to Paillier) or provide a generic HE implementation on FL aggregation [15, 26, 27, 32, 43]. However, previous work still leaves the HE overhead increase issue as an open question. In our work, we propose a universal optimization scheme to largely reduce the overhead while providing promised privacy guarantees in a both systematic and algorithmic fashion, which makes HE-based FL viable in practical deployments.

## 3 System Design

In this section, we first define the problem and describe the algorithmic design of HE-based aggregation; then illustrate our HE implementation with reasoning about the HE libraries. Lastly, we explain the realization of the FedML-HE deployment platform.

### 3.1 Design Goals

**Adversary Definition** We define a semi-honest adversary $\mathcal{A}$ that can corrupt the aggregation server or any subset of local clients. $\mathcal{A}$ follows the protocol but tries to learn as much information as possible. Loosely speaking, in the presence of such an adversary, the security

definition requires that only the private information in local models from the corrupted clients will be learned when $\mathcal{A}$ corrupts a subset of clients; no private information from local models nor global models will be learned by $\mathcal{A}$ when $\mathcal{A}$ corrupts the aggregation server.

When $\mathcal{A}$ corrupts both the aggregation server and a number of clients, the default setup where the private key is shared with all clients (also with corrupted clients) will enable $\mathcal{A}$ to decrypt local models from benign clients (by combining encrypted local models received by the corrupted server and the private key received by any corrupted client). This issue can be solved by adopting the threshold or multi-key variant of HE where decryption has to be collaboratively performed by a certain amount of clients [4, 15, 32].

**Security And Privacy Goal** Our framework should guarantee that in the presence of $\mathcal{A}$, a set of clients wishes to collaboratively train a global model on each individual client's local dataset in the federated setting, but only the aggregated model will be shared among clients. No information about any individual training data will be learned by $\mathcal{A}$ neither directly nor indirectly via attacks such as gradient inversion.

**Efficient Federated Training Goal** Encrypted models by Homomorphic Encryption are around $15\times$ larger than the plaintext models and usually incur $10\times$ computation overhead, which limits HE's scalability on large models across edge devices.

With overhead optimization strategies, the encrypted FL system can efficiently train a good model with fast overall training time with low communication overhead. The system should have similar overall training time and communication cost compared to the plaintext federated training. The system also needs to support different overhead optimization strategies and support for different encryption libraries for providing different levels of overhead reduction and privacy preservation for dynamic user requirements.

## 3.2 Algorithm for HE-Based Federated Learning

Our privacy-preserving federated learning system utilizes Homomorphic Encryption to enable the aggregation server to combine local model parameters without viewing them in their unencrypted form.

As FedAvg has been proven as still one of the most robust federated aggregation strategies while maintaining computational simplicity [52], we primarily implement FedAvg in our system. However, as shown in §3.3 and Figure 4, our system can be easily extended to support more FL aggregation functions with HE by encrypting

and computing the new parameters in these algorithms (e.g. FedProx [30], FedOpt [42], and Scaffold [29]).

Our HE-based secure aggregation algorithm, as shown in Algorithm 1 and illustrated in Figure 3, can be summarized as: $[[\mathbf{W}_{\text{glob}}]] = \sum_{i=1}^{N} \alpha_i [[\mathbf{W}_i]]$, where $[[\mathbf{W}_i]]$ is the $i$th encrypted local model, $\alpha_i$ is weighting factor for client $i$ and $[[\mathbf{W}_{\text{glob}}]]$ is the encrypted global model. Note that weighting factors can be either encrypted or in plaintext depending on whether the aggregation server is trustworthy enough to obtain that information. In our system, we set weighting factors to be plaintext by default. We only need one depth of HE multiplication in our algorithm for weighting, which is preferred to reduce HE multiplication operations.

---

**Algorithm 1** HE-Based Federated Aggregation

---

- Aggregation server $\mathcal{S}$ and $N$ clients;
- Key authority server generates a key pair $(pk, sk)$ and the crypto context, then distributes it to clients and server. (except server does not get $sk$);
- Client $i \in [N]$ owns a local dataset $\mathcal{D}_i$ and initializes a local model $\mathbf{W}_i$ with the aggregation weighing factor $\alpha_i$; $[[\mathbf{W}]]$ is the encrypted model;
- $T$ is the number of communication rounds;

**for** $t = 1, 2, \ldots, T$ **do**
    **for** *each client* $i \in [N]$ **do in parallel**
        **if** $t > 1$ **then**
            Receive $[[\mathbf{W}_{\text{glob}}]]$ from $\mathcal{S}$;
            $\mathbf{W}_i \leftarrow Dec(sk, [[\mathbf{W}_{\text{glob}}]])$;
        **end**
        $\mathbf{W}_i \leftarrow Train(\mathbf{W}_i, \mathcal{D}_i)$;
        $[[\mathbf{W}_i]] \leftarrow Enc(pk, \mathbf{W}_i)$;
        Send $[[\mathbf{W}_i]]$ to server $\mathcal{S}$;
    **end**
    // Server Aggregation
    $[[\mathbf{W}_{\text{glob}}]] \leftarrow \sum_{i=1}^{N} \alpha_i [[\mathbf{W}_i]]$;
**end**

---

## 3.3 HE Implementation

In this part, we will explain in detail how we implement our HE-based aggregation.

Figure 4 provides a high-level design of our framework. Our framework consists of three major layers:

- **Crypto Foundation**: the foundation layer where python wrappers are built to realize HE functions including key generation, encryption/decryption, secure aggregation, and ciphertext serialization using open-sourced HE libraries;
- **ML Bridge**: the bridging layer to connect the FL system orchestration and cryptographic functions. Specifically, we have ML processing APIs to process inputs

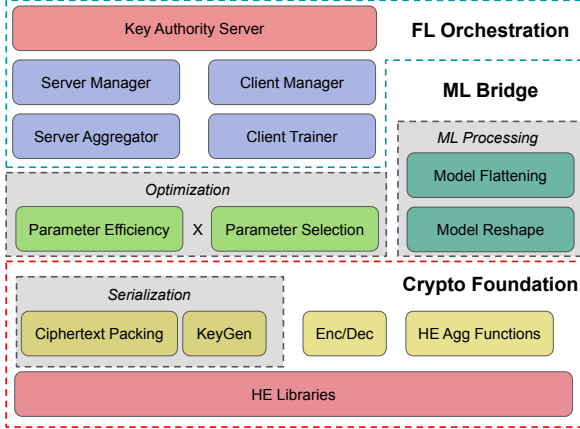| API Name | Description |
|---|---|
| *pk, sk* = **key_gen**(*params*) | Generate a pair of HE keys (public key and private key) |
| *1d_local_model* = **flatten**(*local_model*) | Flatten local trained model tensors into a 1D local model |
| *enc_local_model* = **enc**(*pk, 1d_model*) | Encrypt the 1D model |
| *enc_global_model* = **he_aggregate**(<br>*enc_models[n], weight_factors[n]*) | Homomorphically aggregate a list of 1D local models |
| *dec_global_model* = **dec**(*sk, enc_global_model*) | Decrypt the 1D global model |
| *global_model* = **reshape**(<br>*dec_global_model, model_shape*) | Reshape the 1D global model back to the original shape |

**Table 2:** HE Framework APIs



**Figure 4:** Framework Overview: our framework consists of a three-layer structure including Crypto Foundation to support basic HE building blocks, ML Bridge to connect crypto tools with ML functions, and FL Orchestration to coordinate different parties during a task.

to HE functions from local training processes and outputs vice versa. Additionally, we realize the optimization module here to mitigate the HE overheads;

- **FL Orchestration**: the FL system layer where the key authority server manages the key distribution as well as (server/client) managers and task executors orchestrate participants.

Our layered design makes the HE crypto foundation and the optimization module *semi-independent*, allowing different HE libraries to be easily switched into FedML-HE and further FL optimization techniques to be conveniently added to the system.

Detailed key APIs are listed in Table 2, which helps illustrate the HE-based FL procedure:

(1) **Before local training**: if it is before the first FL round, each client receives the same pair of keys ($pk, sk$) along with related crypto contexts from the key authority server using **key_gen**; otherwise, each client **dec** then **reshape** the encrypted global model.

(2) **After local training**: each client **flatten** then **enc** its local model and sends it to the aggregation server.

(3) **On aggregation**: the server **he_aggregate** a list of encrypted local models with a weight factor list and sends back the encrypted global model.

Note that models are regarded as generic tensors and the HE functions are universally applied to all models without the need for special modifications for each different type of model. It is also worth mentioning that instead of directly applying HE functions on tensors, we **flatten** tensors into the 1D shape. This is because directly encrypting tensors will yield large ciphertexts (e.g. on TenSEAL [6], a small tensor with a size of 5k numbers results in a ciphertext of 1.7 GB), which would easily crash memory on most machines with a reasonable number of local ML models that are also reasonably-sized. However, encrypting a one-dimensional vector largely reduces the ciphertext size.

**Choices of HE libraries** Our Crypto Foundation layer can universally support different HE realizations. Currently, several open-source HE libraries are available [23, 40, 44]. Among these libraries, PALISADE has one of the fastest HE implementations [20] (verified by our experimental results in §5.2.5) and also supports several decentralized multi-party functionalities, such as Proxy Re-Encryption [5, 28] and Threshold Homomorphic Encryption [4] for serverless or decentralized FL scenarios [22]. In our framework, we prioritize PALISADE as our HE core. We use pybind11 [41] to pythonize PALISADE functions (C++) in our framework.

## 3.4 Deploy Anywhere: An FedML-HE Deployment Platform MLOps For Edges/Cloud

We implement our deployment-friendly platform[2] such that FedML-HE can be easily deployed across cloud and edge devices as shown in Figure 5. Before the training starts, a user uploads the configured server package and the local client package to the web platform. The server package defines the operations on the FL server,

---

[2]Our platform is available at https://doc.fedml.ai/mlops/user_guide.html.
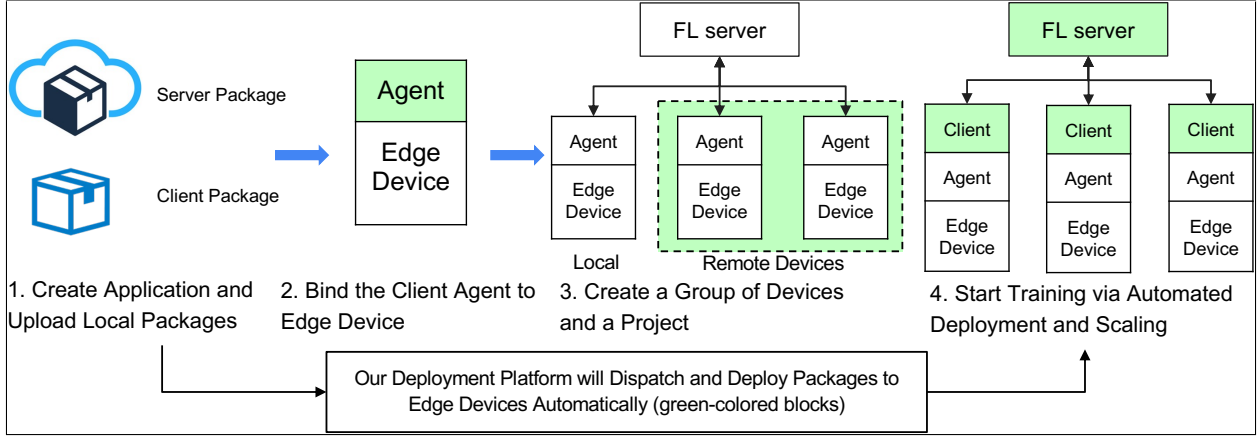
**Figure 5:** Our Deployment Platform: local simulation can be easily deployed to real-world edge-cloud scenarios without code modification

such as the aggregation function and client sampling function; the local client package defines the customized model architecture to be trained (model files will be distributed to edge devices in the first round of the training). Both packages are written in Python. The platform then builds and runs the docker image with the uploaded server package to operate as the server for the training with edge devices configured using the client package.

As shown in Figure 6, during the training, users can also keep tracking the learning procedure including device status, training progress/model performance, and FedML-HE system overheads (e.g., training time, communication time, CPU/GPU utilization, and memory utilization) via the web interface. Our platform keeps close track of overheads, which allows users to in real-time pinpoint HE overhead bottlenecks if any.
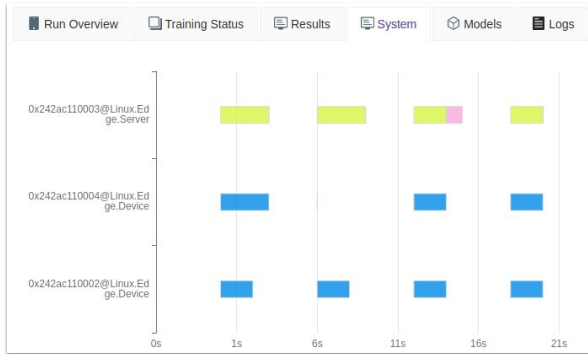


**Figure 6:** Deployment Interface Example of FedML-HE: Overhead distribution monitoring on each edge device (e.g. Desktop (Ubuntu), Laptop (MacBook), and Raspberry Pi 4), which can be used to pinpoint HE overhead bottlenecks and guide optimization.
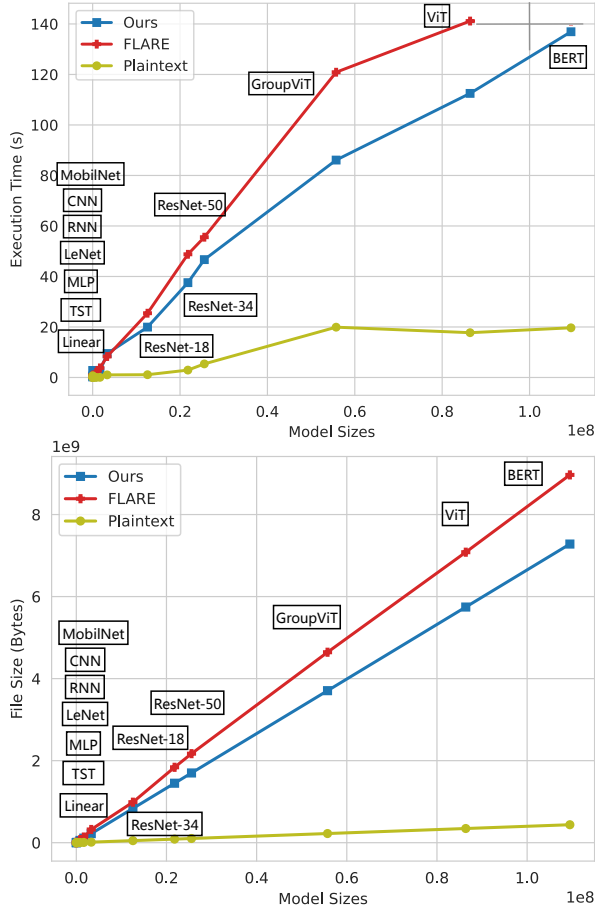
## 4 Optimizations

It is widely acknowledged that HE inevitably introduces large overheads regarding both computation and communication [19].

*Observation*: The computational and communicational (package size) overheads introduced by HE are $O(n)$, both growing linearly with the input size $n$, which in our case the sizes of the models for aggregation. This is also shown by the evaluation results in Figure 7 (also Table 3 in §5.2.1). Empirically, it typically has 10× more computation overhead and 15× more communication overhead than the plaintext FL.
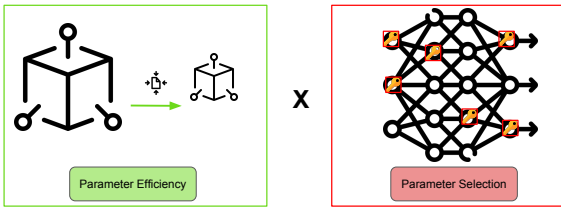
With this observation, we focus on reducing the input model sizes for HE functions as the primary optimization direction, which reduces the overheads of computation and communication simultaneously without loss of general usability and privacy. The input model sizes are determined in two stages: the resulting model parameters from the ML pipeline and the chosen model parameters used in HE functions.

Thus, we propose a universal overhead optimization scheme. As shown in Figure 8, input model sizes of HE can be optimized in two steps: (1) **Parameter Efficiency** to reduce the model sizes from an ML perspective and (2) **Parameter Selection** to selectively choose a certain portion of the model as inputs for HE functions while the rest of the model undergoes a plaintext aggregation at a certain privacy level. The optimization scheme can be defined as $r_{total} = r_{PE} \times r_{PS}$, where $r_{total}$ is the total optimization rate, $r_{PE}$ is the optimization rate from **Parameter Efficiency** and $r_{PS}$ is the optimization rate from **Parameter Selection**.

In the rest of §4, we will discuss in detail how we can optimize overheads in these two steps.

**Figure 7:** (w/o Optimization) Computational (up) and Computation (bottom) Overhead Comparison For Models of Different Sizes: FedML-HE vs. Nvidia FLARE vs. Plaintext Aggregation. Due to TenSeal's larger file sizes, FLARE did not finish the run on BERT on our 32GB memory machine (exceeding memory).



**Figure 8:** A Universal Overhead Optimization Scheme: Parameter Efficiency x Parameter Selection. Reducing the number of model parameters and selectively encrypting them significantly reduces FedML-HE's overhead while preserving privacy.

## 4.1 Parameter Efficiency

In general, there are two directions to downsize the communicated models from an ML perspective: **Model Compression** to minimize the parameters from a local model upon communication and **Parameter-Efficient Tuning** to find a small set of parameters to efficiently tune large models. Although we here introduce several efficiency optimization techniques, further efficiency techniques can also be easily integrated into our scheme. **Model Compression** is the typical method to reduce the number of communication parameters when training from scratch, which compresses the model to have a smaller size for communication [2, 3, 48, 51]. It may include sparsification (only communicate parameters with relatively higher weights), quantization (reduce the number of bits required to store the weights), and low-rank decomposition. The compression technique can largely reduce the number of parameters or the number of bits for each parameter to be encrypted and homomorphically aggregated.

**Parameter-Efficient Tuning** is widely used for fine-tuning large models. For example, LoRA [25] adjusts lightweight trainable parameters while keeping most pre-trained parameters frozen. Under this tuning scheme, only a small amount of trainable parameters will be shared in federated learning for HE.
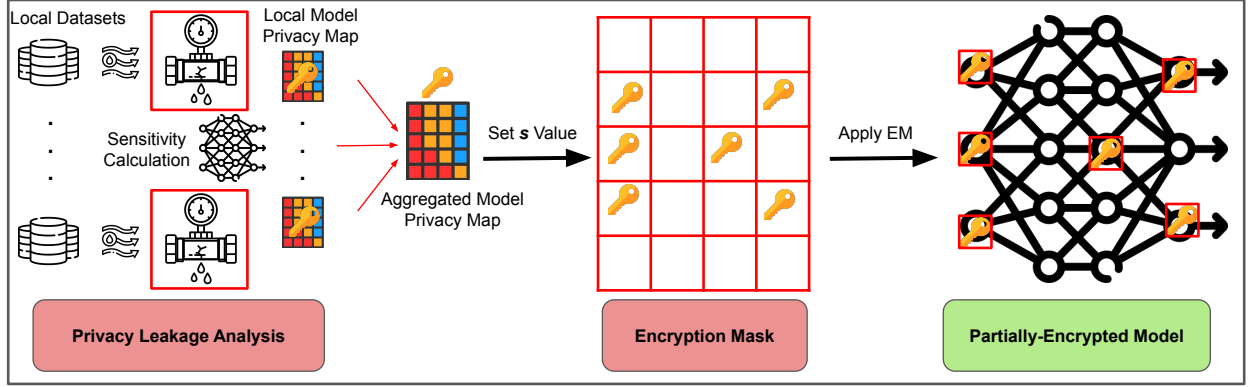
## 4.2 Parameter Selection

Fully encrypted models can guarantee no access to plaintext local models from the adversary. However, previous work on privacy leakage analysis shows that "partial transparency", e.g. hiding parts of the models [21, 34], can grant the adversary a limited chance to successfully perform attacks like gradient inversion attacks [31]. Although FLARE [43] code provides a configuration to only perform HE operations on certain layers, its interaction and tradeoffs regarding privacy leakage are unknown, leaving users a "Pandora's privacy box". On the contrary, to clarify, our solution can support parameter-wise selection for more fine-grained overhead control with privacy leakage analysis (Figure 9).

### 4.2.1 Parameter Privacy Sensitivity Map

Privacy leakage analysis can be done by directly performing gradient inversion attack [54] and evaluating the success rate of the attack, which can take much more time than the model training. Differently, Mutual Information (MI) [55] is proposed to calculate the usable information of model parameters while is practically similar to computing the area under the curve of attack models with even higher computation cost [34].

Sensitivity is then adopted for measuring the general privacy risk on gradients w.r.t. input, high-level features, or output [34, 37, 46]. Theoretically, if the privacy sensitivity of the plaintext parameters obtained by $\mathcal{A}$ is lower, the success rate of the attack can be also expected to be lower. Given model $\mathbf{W}$ and $K$ data samples with input matrix $\mathbf{X}$ and ground truth label vector $\mathbf{y}$, we compute

**Figure 9:** Parameter Selection: in the initialization stage, clients first calculate privacy sensitivities on the model using its own dataset and local sensitivities will be securely aggregated to a global model privacy map. The encryption mask will be then determined by the privacy map and a set selection value *s* per overhead requirements and privacy guarantee. Only the masked parameters will be aggregated in the encrypted form.

the sensitivity of each parameter $w_m$ by

$$\frac{1}{K} \sum_{k=1}^{K} \| J_m\left(y_k\right) \|,$$

where $J_m\left(y_k\right) = \frac{\partial}{\partial y_k}\left(\frac{\partial \ell(\mathbf{X},\mathbf{y},\mathbf{W})}{\partial w_m}\right) \in R$, $\ell(\cdot)$ is the loss function given $\mathbf{X}$, $\mathbf{y}$ and $\mathbf{W}$, and $\|\cdot\|$ calculates the absolute value. The intuition is to calculate how large the gradient of the parameter will change with the true output $y_k$ for each data point $k$.

### 4.2.2 Encryption Mask

As the privacy sensitivity analysis above as well as our results show in §5.3, different parts of a model contribute to attacks by revealing uneven amounts of information. Using this insight, we propose to only select and encrypt parts of the model that are more important and susceptible to attacks to reduce HE overheads while maintaining the privacy guarantee. We use encryption masks (EM) to effectively select encrypted parameters. EM works by first pre-calculating the model parameter privacy map via the privacy sensitivity and applying the mask with a set privacy threshold to pick parameters to encrypt that fit certain overhead expectations.

### 4.2.3 Mask Agreement

The encryption mask is dependent on the data it is processed on. With potentially different data distributions at each client, it is challenging to have all clients agree on a certain global encryption mask without revealing local datasets since the local privacy map carries information about the client data. To solve this problem, we encrypted aggregate local privacy maps at the server such that no further private information about local data is revealed. During the initialization, clients first calculate the local privacy sensitivities using their own

local data and then the server aggregates local sensitivity maps to a global privacy map (as shown in the first part of Figure 9). The encryption mask is configured using a privacy-overhead threshold *s* and the global privacy map is then shared among clients as part of the federated learning configuration.

## 5 Evaluation

In this section, we first evaluate the vanilla (without optimization) version of our framework's performance with an empathized focus on diagnosing the overheads (both in computation and communication) introduced by Homomorphic Encryption. Later in §5.3, we demonstrate the evaluation results to show how our proposed universal optimization scheme largely mitigates these overheads for real-world deployment.

### 5.1 Experiment Setup

**Models:** we tested our framework on models in different ML domains with different sizes (Table 3).

**Dataset:** MNIST dataset ($70k$ images) for experiments and the CIFAR-100 dataset ($50k$ images).

**HE Libraries:** we implement our HE core using both PALISADE and TenSEAL. Unless otherwise specified, our results show the evaluation of the PALISADE version (faster and with more HE functionalities).

**Default Crypto Parameters:** unless otherwise specified, we choose the multiplicative depth of 1, the scaling factor bit digit of 52, an HE packing batch size of 4096, and a securityLevel of 128 as our default HE cryptographic parameters during the evaluation.

**Machine:** we use an Intel 8-core 3.60GHz i7-7700 CPU with 32 GB memory and an NVIDIA Tesla T4 GPU on Ubuntu 18.04.6 LTS (Docker Env) for primarily microbenchmarking overheads of the HE operations.

| Model | Task | Model Size | HE Time (s) | Non-HE Time (s) | Comp Ratio | Ciphertext | Plaintext | Comm Ratio |
|---|---|---|---|---|---|---|---|---|
| Linear Model | Regression | 101 | 0.216 | 0.001 | 150.85 | 266.00 KB | 1.10 KB | 240.83 |
| TimeSeries Transformer | Time Series | 5,609 | 2.792 | 0.233 | 12.00 | 532.00 KB | 52.65 KB | 10.10 |
| MLP (2 FC) | Classification Regression | 79,510 | 0.586 | 0.010 | 60.46 | 5.20 MB | 311.98 KB | 17.05 |
| LeNet | CV | 88,648 | 0.619 | 0.011 | 57.95 | 5.97 MB | 349.52 KB | 17.50 |
| RNN(2 LSTM + 1 FC) | NLP | 822,570 | 1.195 | 0.013 | 91.82 | 52.47 MB | 3.14 MB | 16.70 |
| CNN (2 Conv + 2 FC) | CV | 1,663,370 | 2.456 | 0.058 | 42.23 | 103.15 MB | 6.35 MB | 16.66 |
| MobileNet | CV | 3,315,428 | 9.481 | 1.031 | 9.20 | 210.41 MB | 12.79 MB | 16.45 |
| ResNet-18 | CV | 12,556,426 | 19.950 | 1.100 | 18.14 | 796.70 MB | 47.98 MB | 16.61 |
| ResNet-34 | CV | 21,797,672 | 37.555 | 2.925 | 12.84 | 1.35 GB | 83.28 MB | 16.60 |
| ResNet-50 | CV | 25,557,032 | 46.672 | 5.379 | 8.68 | 1.58 GB | 97.79 MB | 16.58 |
| GroupViT | Multi Modal | 55,726,609 | 86.098 | 19.921 | 4.32 | 3.45 GB | 212.83 MB | 16.61 |
| Vision Transformer | CV | 86,389,248 | 112.504 | 17.739 | 6.34 | 5.35 GB | 329.62 MB | 16.62 |
| BERT | NLP | 109,482,240 | 136.914 | 19.674 | 6.96 | 6.78 GB | 417.72 MB | 16.62 |

**Table 3:** Vanilla Fully-Encrypted Models of Different Sizes: with 3 clients; Comp Ratio is calculated by time costs of HE over time costs of Non-HE; Comm Ratio is calculated by file sizes of HE over file sizes of Non-HE. CKKS is configured with default crypto parameters.

## 5.2 General HE Impacts On Overheads

| HE Batch Size | Scaling Bits | Comp (s) | Comm (MB) | Model Test Accuracy Δ (%) |
|---|---|---|---|---|
| 1024 | 14 | 8.834 | 407.47 | -0.28 |
| 1024 | 20 | 7.524 | 407.47 | -0.21 |
| 1024 | 33 | 7.536 | 407.47 | 0 |
| 1024 | 40 | 7.765 | 407.47 | 0 |
| 1024 | 52 | 7.827 | 407.47 | 0 |
| 2048 | 14 | 3.449 | 204.50 | -0.06 |
| 2048 | 20 | 3.414 | 204.50 | -0.13 |
| 2048 | 33 | 3.499 | 204.50 | 0 |
| 2048 | 40 | 3.621 | 204.50 | 0 |
| 2048 | 52 | 3.676 | 204.50 | 0 |
| 4096 | 14 | 1.837 | 103.15 | -1.85 |
| 4096 | 20 | 1.819 | 103.15 | 0.32 |
| 4096 | 33 | 1.886 | 103.15 | 0 |
| 4096 | 40 | 1.998 | 103.15 | 0 |
| 4096 | 52 | 1.926 | 103.15 | 0 |

**Table 4:** Computational & Communicational Overhead of Different Crypto Parameter Setups: tested with CNN (2 Conv+ 2 FC) and on 3 clients; model test accuracy Δs is the difference between the best plaintext global model and the best global encrypted global models.

We evaluate the HE-based training overheads (without our optimization in place) across various FL training scenarios and configurations. This analysis covers diverse model scales, HE cryptographic parameter configurations, client quantities involved in the task, and communication bandwidths. This helps us to identify bottlenecks in the HE process throughout the entire training cycle. We also benchmark our framework against other open-source HE solutions to demonstrate its advantages.

*5.2.1 Results on Different Scales of Models*
We evaluate our framework on models with different size scales and different domains, from small models like the linear model to large foundation models such as Vision Transformer [14] and BERT [13]. As Table 3, Figure 7 show, both computational and communicational overheads are generally proportional to model sizes.

Table 3 illustrates more clearly the overhead increase from the plaintext federated aggregation. The computation fold ratio is in general 5x ∼ 20x while the communication overhead can jump to a common 15x. Small models tend to have a higher computational overhead ratio increase. This is mainly due to the standard HE initialization process, which plays a more significant role when compared to the plaintext cost. The communication cost increase is significant for models with sizes smaller than 4096 (the packing batch size) numbers. Recall that the way our HE core packs encrypted numbers makes an array whose size is smaller than the packing batch size still requires a full ciphertext.

### 5.2.2 Results on Different Cryptographic Parameters

We evaluate the impacts of variously-configured cryptographic parameters. We primarily look into the packing batch size and the scaling bits. The packing batch size determines the number of slots packed in a single ciphertext while the scaling bit number affects the "accuracy" (i.e., how close the decrypted ciphertext result is to the plaintext result) of approximate numbers represented from integers.

From Table 4, the large packing batch sizes in general result in faster computation speeds and smaller overall ciphertext files attributed to the packing mechanism for more efficiency. However, the scaling factor number has an almost negligible impact on overheads.

Unsurprisingly, it aligns with the intuition that the higher bit scaling number results in higher "accuracy" of the decrypted ciphertext value, which generally means the encrypted aggregated model would have a close model test performance to the plaintext aggregated model. However, it is worth mentioning that since CKKS is an approximate scheme with noises, the decrypted aggregated model can yield either positive or negative model test accuracy Δs, but usually with a negative or nearly zero Δ.

### 5.2.3 Impact from Number of Clients

As real-world systems often experience a dynamic amount of participants within the FL system, we evaluate the overhead shift over the change in the number of clients. Figure 11 breaks down the cost distribution as the number of clients increases. With a growing number of clients, it also means proportionally-added ciphertexts as inputs to the secure aggregation function thus the major impact is cast on the server. When the server is overloaded, our system also supports client selection to remove certain clients without largely degrading model performance.

### 5.2.4 Communication Cost on Different Bandwidths

FL parties can be allocated in different geo-locations which might result in communication bottlenecks. Typically, there are two common scenarios: (inter) data centers and (intra) data centers. In this part, we evaluate the impact of the bandwidths on communication costs and how it affects the FL training cycle. We categorize communication bandwidths using 3 cases:

- Infiniband (IB): communication between intra-center parties. 5 GB/s as the test bandwidth.
- Single AWS Region (SAR): communication between inter-center parties but within the same geo-region (within US-WEST). 592 MB/s as the test bandwidth.
- Multiple AWS Region (MAR): communication between inter-center parties but across the different

geo-region (between US-WEST and EU-NORTH). 15.6 MB/s as the test bandwidth.

As shown in Figure 12, we deploy FedML-HE on 3 different geo-distributed environments, which are operated under different bandwidths. It is obvious that the secure HE functionality has an enormous impact on low-bandwidth environments while medium-to-high-bandwidth environments suffer limited impact from increased communication overhead during training cycles, compared to Non-HE settings.

### 5.2.5 Comparison with Other FL-HE Frameworks

We compare our framework to the other open-sourced FL frameworks with HE capability, namely NVIDIA FLARE (NVIDIA) and IBMFL.
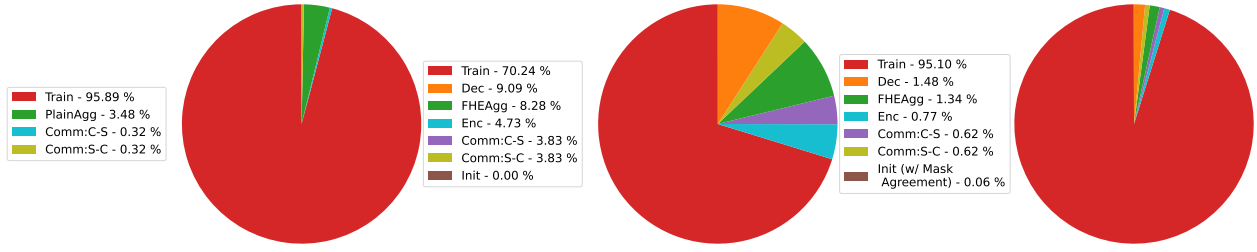
Both NVIDIA and IBMFL utilize Microsoft SEAL as the underlying HE core, with NVIDIA using Open-Minded's python tensor wrapper over SEAL and TenSEAL; IBMFL using IBM'spython wrapper over SEAL and HELayers (HELayers also has an HElib version). As §3.3 states that our HE core module can be replaced with different available HE cores, to give a more comprehensive comparison, we also implement a TenSEAL version of our framework for evaluation.

Table 5 demonstrates the performance summary of different FedML-HE frameworks using an example of a CNN model with 3 clients. Our PALISADE-powered framework has the smallest computational overhead due to the performance of the PALISADE library. In terms of communication cost, FedML-HE (PALISADE) comes second after IBMFL's smallest file serialization results due to the efficient packing of HELayers' Tile tensors [1].

Note that NVIDIA's TenSEAL-based realization is faster than the TenSEAL variant of our system. This is because NVIDIA scales each learner's local model parameters locally rather than weighing ciphertexts on the server. This approach reduces the need for the one multiplication operation usually performed during secure aggregation (recall that HE multiplications are expensive). However, such a setup would not suit the scenario where the central server does not want to reveal its weighing mechanism per each individual local model to learners as it reveals partial (even full in some cases) information about participants in the system.
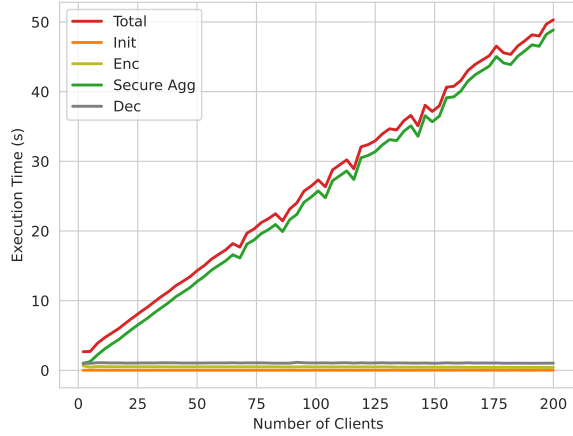
## 5.3 Optimizations

To mitigate the HE overhead surge, our universal optimization scheme works in two stages as described in §4: (1) **Parameter Efficiency** reduces the shared model parameters and (2) **Parameter Selection** selects certain portions of parameters for encrypted computation while

**Figure 10:** Time Distribution of A Non-HE Cycle on ResNet-50: assume an expected bandwidth of 200 MB/s and fixed local training time for plaintext FL (left), HE w/o optimization (middle), and HE w/ optimization (middle). Optimization setup uses *DoubleSqueeze* with $k = 1,000,000$ and encryption mask with $s = 30\%$.

| Frameworks | HE Core | Key Management | Deployment | Comp (s) | Comm (MB) | HE Multi-Party Functionalities |
|---|---|---|---|---|---|---|
| Ours | PALISADE | ✓ | ✓ | 2.456 | 105.72 | PRE, ThHE |
| Ours (w/ Opt) | PALISADE | ✓ | ✓ | 0.874 | 16.37 | PRE, ThHE |
| Ours | SEAL (TenSEAL) | ✓ | ✓ | 3.989 | 129.75 | — |
| Nvidia FLARE (9a1b226) | SEAL (TenSEAL) | ✓ | NVFLARE Dashboard | 2.826 | 129.75 | — |
| IBMFL (8c8ab11) | SEAL (HELayers) | ○ | Multi-Cloud OpenShift Orchestrator | 3.955 | 86.58 | — |
| Plaintext | — | — | — | 0.058 | 6.35 | — |

**Table 5:** Different Frameworks: tested with CNN (2 Conv + 2 FC) and on 3 clients; Github commit IDs are specified. For key management, our work uses a key authority server; FLARE uses a security content manager; IBMFL currently provides a local simulator.
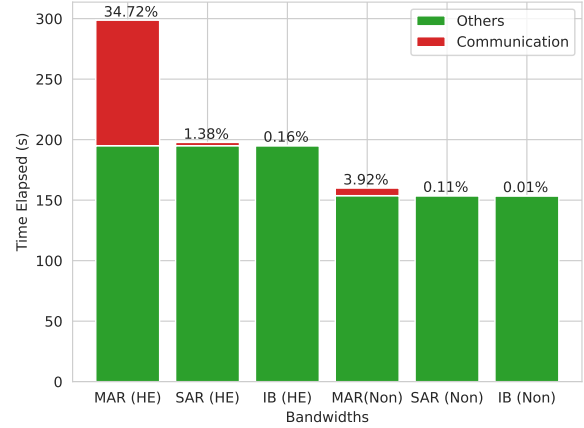


**Figure 11:** Step Breakdown of HE Computational Cost vs. Number of Clients (Up to 200): tested on fully-encrypted CNN



**Figure 12:** Impact of Different Bandwidths on Communication and Training Cycles on Fully-Encrypted ResNet-50: HE means HE-enabled training and Non means plaintext. Others include all other procedures except communication during training. Percentages represent the portion of communication cost in the entire training cycle.

leaving the rest in plaintext per desired overhead expectations and privacy promise. The overhead reductions from two stages stack as the final optimization.
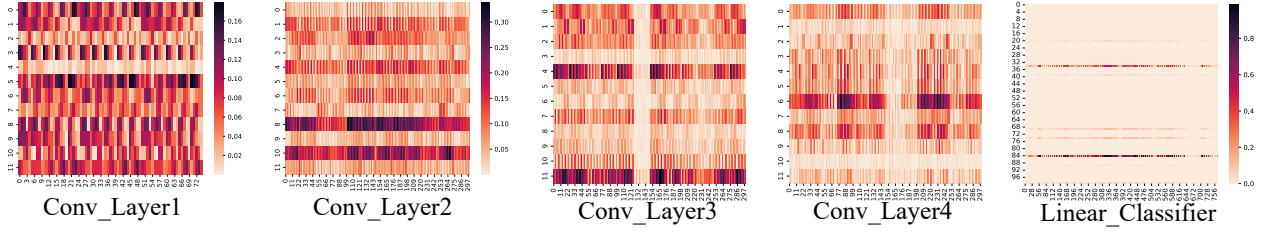
### 5.3.1 Optimized Overhead

In Table 6, we first case study different **Parameter Efficiency** ML techniques for both training-from-scratch and fine-tuning scenarios. Efficiently reducing the sizes of shared models directly helps with HE computation and communication efficiency.

| Conv_Layer1 | Conv_Layer2 | Conv_Layer3 | Conv_Layer4 | Linear_Classifier |

**Figure 13:** Model Privacy Map Calculated Using Parameter Privacy Sensitivity On LeNet: darker color indicates higher sensitivity. Each subfigure shows the sensitivity of parameters of the current layer. The sensitivity of parameters is imbalanced and many parameters have very little sensitivity (its gradient is hard to be affected by tuning the data input for attack).

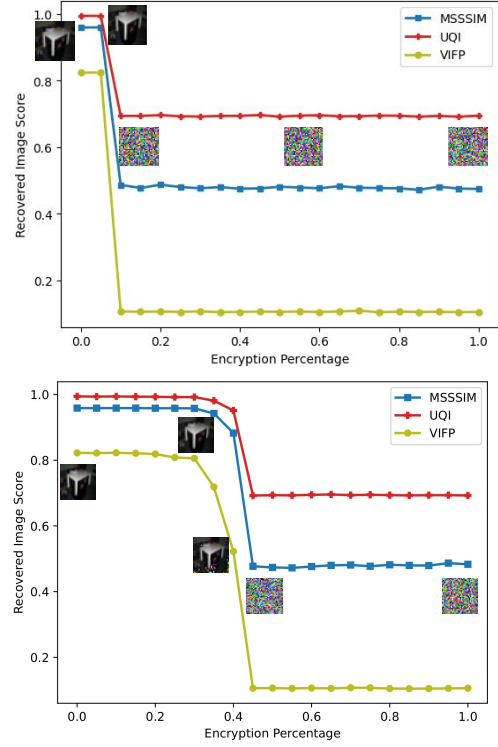| Efficiency Optimization | PT | CT | CT (Opt) |
|---|---|---|---|
| ResNet-18 (12 M) [48] | 47.98 MB | 796.70 MB | 19.03 MB |
| BERT (110 M) [25] | 417.72 MB | 6.78 GB | 16.66 MB |

**Table 6:** Parameter Efficiency Overhead: PT means plaintext and CT means ciphertext. Communication reductions are 0.60 and 0.96.

| Selection | Comp (s) | Comm | Comp Ratio | Comm Ratio |
|---|---|---|---|---|
| Enc w/ 0% | 17.739 | 329.62 MB | 1.00 | 1.00 |
| Enc w/ 10% | 30.874 | 844.49 MB | 1.74 | 2.56 |
| Enc w/ 30% | 50.284 | 1.83 GB | 2.83 | 5.69 |
| Enc w/ 50% | 70.167 | 2.83 GB | 3.96 | 8.81 |
| Enc w/ 70% | 88.904 | 3.84 GB | 5.01 | 11.93 |
| Enc w/ All | 112.504 | 5.35 GB | 6.34 | 16.62 |

**Table 7:** Overheads With Different Parameter Selection Configs Tested on Vision Transformer: "Enc w/ 10%" means performs encrypted computation only on 10% of the parameters; all computation and communication results include overheads from plaintext aggregation for the rest of the parameters.

We then examine the overhead optimization effects from **Parameter Selection**. We examine the overhead change when parameters with high privacy importance are selected and encrypted. Table 7 shows the overhead reduction from only encrypting certain parts of models, where both overheads are nearly proportional to the size of encrypted model parameters, which is coherent with the general relationship between HE overheads and input sizes.

Figure 10 dissects the training cycle composition for the HE framework (both with and without optimizations) and the plaintext framework respectively. For a medium-sized model, the overheads (both computation and communication) from HE shift some portion of the local training procedure to aggregation-related steps compared to Non-HE, but not with an infeasible margin relatively speaking. Though generally smaller models require shorter training time, the overheads of the HE-based aggregation also drop proportionally.



**Figure 14:** Selection Protection Against Gradient Inversion Attack On LeNet: attack results when protecting top-$s$ sensitive parameters (top) vs protecting random parameters (bottom). Each configuration is attacked 10 times and the best-recovered image is selected.

### 5.3.2 Effectiveness of Selection Defense

To evaluate the defense effectiveness of **Parameter Selection**, we first use privacy sensitivity to generate a privacy map (Figure 13) and then verify the effectiveness of selection by performing gradient inversion (DLG [58]).

DLG works by adjusting potential recovery data according to the loss feedback between the actual model gradient updates and the generated gradients.

We use image samples from CIFAR-100 to calculate the parameter sensitivities of the model. In the DLG attack experiments, we use Multi-scale Structural Similarity Index (MSSSIM), Visual Information Fidelity (VIF), and Universal Quality Image Index (UQI) as metrics to

measure the similarity between recovered images and original training images to measure the attack quality hence the privacy leakage[3]. In Figure 14, compared to random encryption selection where encrypting 42.5% of the parameters can start to protect against attacks, our top-10% encryption selection according to the model privacy map only alone can defend against the attacks, meaning lower overall overhead with the same amount of privacy protection.

**Empirical Selection Recipe** Our selection strategy works by first encrypting more important model parameters. Empirically, from our experimental investigation, model parameters in the beginning and end stages tend to be more important for information leakage [21] and attack defense, which can be used as a general guideline on top of model privacy maps.

## 6 Conclusion

In this paper, we propose an efficient Homomorphic-Encryption-based privacy-preserving federated learning system with a user/device-friendly deployment platform and a universal optimization scheme that can be practically deployed across distributed edge devices. To support realistic secure aggregation application deployment, our novel universal overhead optimization scheme is designed to largely reduce both computational and communicational overheads during deployment while providing a dynamic desired privacy guarantee. Future work includes quantitative and theoretical analysis of the trade-offs among privacy guarantee, system overheads and model performance compared to other approaches (including Difference Privacy and MPC approaches), and supporting decentralized serverless secure federated training using primitives like Threshold HE and Proxy Re-Encryption [5].

## References

[1] Ehud Aharoni, Allon Adir, Moran Baruch, Nir Drucker, Gilad Ezov, Ariel Farkash, Lev Greenberg, Ramy Masalha, Guy Moshkowich, Dov Murik, et al. 2011. HeLayers: A tile tensors framework for large neural networks on encrypted data.

[2] Alham Fikri Aji and Kenneth Heafield. 2017. Sparse communication for distributed gradient descent. *arXiv preprint arXiv:1704.05021* (2017).

[3] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. 2017. QSGD: Communication-efficient SGD via gradient quantization and encoding. *Advances in neural information processing systems* 30 (2017).

[4] Asma Aloufi, Peizhao Hu, Yongsoo Song, and Kristin Lauter. 2021. Computing blindfolded on data homomorphically encrypted under multiple keys: A survey. *ACM Computing Surveys (CSUR)* 54, 9 (2021), 1–37.

[5] Giuseppe Ateniese, Kevin Fu, Matthew Green, and Susan Hohenberger. 2006. Improved proxy re-encryption schemes with applications to secure distributed storage. *ACM Transactions on Information and System Security (TISSEC)* 9, 1 (2006), 1–30.

[6] Ayoub Benaissa, Bilal Retiat, Bogdan Cebere, and Alaa Eddine Belfedhal. 2021. TenSEAL: A Library for Encrypted Tensor Operations Using Homomorphic Encryption. arXiv:2104.03152 [cs.CR]

[7] Abhishek Bhowmick, John Duchi, Julien Freudiger, Gaurav Kapoor, and Ryan Rogers. 2018. Protection against reconstruction and its applications in private federated learning. *arXiv preprint arXiv:1812.00984* (2018).

[8] Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. 2017. Practical secure aggregation for privacy-preserving machine learning. In *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. 1175–1191.

[9] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. 2014. (Leveled) fully homomorphic encryption without bootstrapping. *ACM Transactions on Computation Theory (TOCT)* 6, 3 (2014), 1–36.

[10] David Byrd and Antigoni Polychroniadou. 2020. Differentially private secure multi-party computation for federated learning in financial applications. In *Proceedings of the First ACM International Conference on AI in Finance*. 1–9.

[11] Jung Hee Cheon, Andrey Kim, Miran Kim, and Yongsoo Song. 2017. Homomorphic encryption for arithmetic of approximate numbers. In *Advances in Cryptology–ASIACRYPT 2017: 23rd International Conference on the Theory and Applications of Cryptology and Information Security, Hong Kong, China, December 3-7, 2017, Proceedings, Part I 23*. Springer, 409–437.

[12] John Criswell, Nathan Dautenhahn, and Vikram Adve. 2014. KCoFI: Complete control-flow integrity for commodity operating system kernels. In *2014 IEEE symposium on security and privacy*. IEEE, 292–307.

[13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).

[15] Weidong Du, Min Li, Liqiang Wu, Yiliang Han, Tanping Zhou, and Xiaoyuan Yang. 2023. A efficient and robust privacy-preserving framework for cross-device federated learning. *Complex & Intelligent Systems* (2023), 1–15.

[16] Junfeng Fan and Frederik Vercauteren. 2012. Somewhat Practical Fully Homomorphic Encryption. Cryptology ePrint Archive, Paper 2012/144. https://eprint.iacr.org/2012/144 https://eprint.iacr.org/2012/144.

[17] Haokun Fang and Quan Qian. 2021. Privacy preserving machine learning with homomorphic encryption and federated learning. *Future Internet* 13, 4 (2021), 94.

[18] Craig Gentry. 2009. Fully homomorphic encryption using ideal lattices. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*. 169–178.

[19] Charles Gouert, Dimitris Mouris, and Nektarios Georgios Tsoutsos. 2022. New insights into fully homomorphic encryption

---

[3]The image similarity metric library used is at https://pypi.org/project/sewar/.

libraries via standardized benchmarks. *Cryptology ePrint Archive* (2022).

[20] Charles Gouert, Dimitris Mouris, and Nektarios Georgios Tsoutsos. 2022. New Insights into Fully Homomorphic Encryption Libraries via Standardized Benchmarks. Cryptology ePrint Archive, Paper 2022/425. https://eprint.iacr.org/2022/425 https://eprint.iacr.org/2022/425.

[21] Ali Hatamizadeh, Hongxu Yin, Holger R Roth, Wenqi Li, Jan Kautz, Daguang Xu, and Pavlo Molchanov. 2022. Gradvit: Gradient inversion of vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10021–10030.

[22] Chaoyang He, Emir Ceyani, Keshav Balasubramanian, Murali Annavaram, and Salman Avestimehr. 2022. Spreadgnn: Decentralized multi-task federated learning for graph neural networks on molecular data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 6865–6873.

[23] helib 2022. HELib. https://github.com/HomEnc/HElib. HELib.

[24] Briland Hitaj, Giuseppe Ateniese, and Fernando Perez-Cruz. 2017. Deep models under the GAN: information leakage from collaborative deep learning. In *Proceedings of the 2017 ACM SIGSAC conference on computer and communications security*. 603–618.

[25] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685* (2021).

[26] IBM. 2022. IBMFL Crypto. https://github.com/IBM/federated-learning-lib/blob/main/Notebooks/crypto_fhe_pytorch/pytorch_classifier_aggregator.ipynb. Accessed: 2023-1-25.

[27] Zhifeng Jiang, Wei Wang, and Yang Liu. 2021. Flashe: Additively symmetric homomorphic encryption for cross-silo federated learning. *arXiv preprint arXiv:2109.00675* (2021).

[28] Weizhao Jin, Bhaskar Krishnamachari, Muhammad Naveed, Srivatsan Ravi, Eduard Sanou, and Kwame-Lante Wright. 2022. Secure Publish-Process-Subscribe System for Dispersed Computing. In *2022 41st International Symposium on Reliable Distributed Systems (SRDS)*. IEEE, 58–68.

[29] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. 2020. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*. PMLR, 5132–5143.

[30] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. 2020. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems* 2 (2020), 429–450.

[31] Jiahao Lu, Xi Sheryl Zhang, Tianli Zhao, Xiangyu He, and Jian Cheng. 2022. APRIL: Finding the Achilles' Heel on Privacy for Vision Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10051–10060.

[32] Jing Ma, Si-Ahmed Naas, Stephan Sigg, and Xixiang Lyu. 2022. Privacy-preserving federated learning based on multi-key homomorphic encryption. *International Journal of Intelligent Systems* 37, 9 (2022), 5880–5901.

[33] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*. PMLR, 1273–1282.

[34] Fan Mo, Anastasia Borovykh, Mohammad Malekzadeh, Hamed Haddadi, and Soteris Demetriou. 2020. Layer-wise characterization of latent information leakage in federated learning. In *ICLR Distributed and Private Machine Learning workshop*.

[35] Viraaji Mothukuri, Reza M Parizi, Seyedamin Pouriyeh, Yan Huang, Ali Dehghantanha, and Gautam Srivastava. 2021. A survey on security and privacy of federated learning. *Future Generation Computer Systems* 115 (2021), 619–640.

[36] Milad Nasr, Reza Shokri, and Amir Houmansadr. 2019. Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In *2019 IEEE symposium on security and privacy (SP)*. IEEE, 739–753.

[37] Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. 2018. Sensitivity and Generalization in Neural Networks: an Empirical Study. In *International Conference on Learning Representations*.

[38] Nvidia. 2021. NVIDIA FLARE Federated Learning with Homomorphic Encryption. https://developer.nvidia.com/blog/federated-learning-with-homomorphic-encryption. Accessed: 2023-1-25.

[39] Pascal Paillier. 1999. Public-key cryptosystems based on composite degree residuosity classes. In *Advances in Cryptology—EUROCRYPT'99: International Conference on the Theory and Application of Cryptographic Techniques Prague, Czech Republic, May 2–6, 1999 Proceedings 18*. Springer, 223–238.

[40] palisade 2022. PALISADE Lattice Cryptography Library. https://gitlab.com/palisade/palisade-release. PALISADE Project.

[41] pybind 2022. pybind11: Seamless operability between C++11 and Python. https://github.com/pybind/pybind11.

[42] Sashank J Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and Hugh Brendan McMahan. 2021. Adaptive Federated Optimization. In *International Conference on Learning Representations*.

[43] Holger R Roth, Yan Cheng, Yuhong Wen, Isaac Yang, Ziyue Xu, Yuan-Ting Hsieh, Kristopher Kersten, Ahmed Harouni, Can Zhao, Kevin Lu, et al. 2022. NVIDIA FLARE: Federated Learning from Simulation to Real-World. *arXiv preprint arXiv:2210.13291* (2022).

[44] SEAL 2023. Microsoft SEAL (release 4.1). https://github.com/Microsoft/SEAL. Microsoft Research, Redmond, WA.

[45] Jinhyun So, Corey J Nolet, Chien-Sheng Yang, Songze Li, Qian Yu, Ramy E Ali, Basak Guler, and Salman Avestimehr. 2022. Lightsecagg: a lightweight and versatile design for secure aggregation in federated learning. *Proceedings of Machine Learning and Systems* 4 (2022), 694–720.

[46] Jure Sokolić, Raja Giryes, Guillermo Sapiro, and Miguel RD Rodrigues. 2017. Robust large margin deep neural networks. *IEEE Transactions on Signal Processing* 65, 16 (2017), 4265–4280.

[47] Dimitris Stripelis, Hamza Saleem, Tanmay Ghai, Nikhil Dhinagar, Umang Gupta, Chrysovalantis Anastasiou, Greg Ver Steeg, Srivatsan Ravi, Muhammad Naveed, Paul M Thompson, et al. 2021. Secure neuroimaging analysis using federated learning with homomorphic encryption. In *17th International Symposium on Medical Information Processing and Analysis*, Vol. 12088. SPIE, 351–359.

[48] Hanlin Tang, Chen Yu, Xiangru Lian, Tong Zhang, and Ji Liu. 2019. Doublesqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. In *International Conference on Machine Learning*. PMLR, 6155–6165.

[49] Stacey Truex, Nathalie Baracaldo, Ali Anwar, Thomas Steinke, Heiko Ludwig, Rui Zhang, and Yi Zhou. 2019. A hybrid approach to privacy-preserving federated learning. In *Proceedings of the 12th ACM workshop on artificial intelligence and security*. 1–11.

[50] Stacey Truex, Ling Liu, Mehmet Emre Gursoy, Lei Yu, and Wenqi Wei. 2019. Demystifying membership inference attacks in machine learning as a service. *IEEE Transactions on Services Computing* 14, 6 (2019), 2073–2089.

[51] Hongyi Wang, Scott Sievert, Shengchao Liu, Zachary Charles, Dimitris Papailiopoulos, and Stephen Wright. 2018. Atomo: Communication-efficient learning via atomic sparsification. *Advances in Neural Information Processing Systems* 31 (2018).

[52] Jianyu Wang, Rudrajit Das, Gauri Joshi, Satyen Kale, Zheng Xu, and Tong Zhang. 2022. On the unreasonable effectiveness of federated averaging with heterogeneous data. *arXiv preprint arXiv:2206.04723* (2022).

[53] Zhibo Wang, Mengkai Song, Zhifei Zhang, Yang Song, Qian Wang, and Hairong Qi. 2019. Beyond inferring class representatives: User-level privacy leakage from federated learning. In *IEEE INFOCOM 2019-IEEE conference on computer communications*. IEEE, 2512–2520.

[54] Wenqi Wei, Ling Liu, Margaret Loper, Ka-Ho Chow, Mehmet Emre Gursoy, Stacey Truex, and Yanzhao Wu. 2020. A framework for evaluating client privacy leakages in federated learning. In *Computer Security–ESORICS 2020: 25th European Symposium on Research in Computer Security, ESORICS 2020, Guildford, UK, September 14–18, 2020, Proceedings, Part I 25*. Springer, 545–566.

[55] Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stewart, and Stefano Ermon. 2020. A Theory of Usable Information under Computational Constraints. In *International Conference on Learning Representations*.

[56] Yuhang Yao and Carlee Joe-Wong. 2022. Fedgcn: Convergence and communication tradeoffs in federated training of graph convolutional networks. *arXiv preprint arXiv:2201.12433* (2022).

[57] Chengliang Zhang, Suyi Li, Junzhe Xia, Wei Wang, Feng Yan, and Yang Liu. 2020. Batchcrypt: Efficient homomorphic encryption for cross-silo federated learning. In *Proceedings of the 2020 USENIX Annual Technical Conference (USENIX ATC 2020)*.

[58] Ligeng Zhu, Zhijian Liu, and Song Han. 2019. Deep leakage from gradients. *Advances in neural information processing systems* 32 (2019).