
SECFWT: Efficient Privacy-Preserving Fine-Tuning of Large Language Models Using Forward-Only Passes

Jinglong Luo^{*1,2} Zhuo Zhang^{*1,2} Yehong Zhang^{†2} Shiyu Liu⁵
 Ye Dong⁶ Xun Zhou^{1,2} Hui Wang² Yue Yu² Zenglin Xu^{†2,3,4}
¹Harbin Institute of Technology, Shenzhen, ²Pengcheng Laboratory
³Fudan University, ⁴Shanghai Academy of AI for Science
⁵Southwestern University of Finance and Economics
⁶National University of Singapore
 {luojl, zhangyh02}@pcl.ac.cn, zenglin@gmail.com

Abstract

Large language models (LLMs) have transformed numerous fields, yet their adaptation to specialized tasks in privacy-sensitive domains, such as healthcare and finance, is constrained by the scarcity of accessible training data due to stringent privacy requirements. Secure multi-party computation (MPC)-based privacy-preserving machine learning offers a powerful approach to protect both model parameters and user data, but its application to LLMs has been largely limited to inference, as fine-tuning introduces significant computational challenges, particularly in privacy-preserving backward propagation and optimizer operations. This paper identifies two primary obstacles to MPC-based privacy-preserving fine-tuning of LLMs: (1) the substantial computational overhead of backward and optimizer processes, and (2) the inefficiency of softmax-based attention mechanisms in MPC settings. To address these challenges, we propose SECFWT, *the first MPC-based framework designed for efficient, privacy-preserving LLM fine-tuning*. SECFWT introduces a forward-only tuning paradigm to eliminate backward and optimizer computations and employs MPC-friendly Random Feature Attention to approximate softmax attention, significantly reducing costly non-linear operations and computational complexity. Experimental results demonstrate that SECFWT delivers substantial improvements in efficiency and privacy preservation, enabling scalable and secure fine-tuning of LLMs for privacy-critical applications. Code is available by clicking here.

1 Introduction

Fine-tuning pre-trained large language models (LLMs) Vaswani et al. [2017], Devlin et al. [2019], OpenAI [2023] using high-quality and domain-specific data presents significant potential for enhancing LLMs’ capabilities for real-world applications (e.g., hospitals, financial institutions). This technology is typically deployed in the cloud by commercial service providers, curating training data from terminal data custodians to fine-tune LLMs [Qu et al., 2021]. Given the increasing value of data, data custodians are often reluctant, or even legally constrained, to update it for training LLMs, stemming from concerns over data privacy, ownership, and regulatory compliance³. On the other hand, the updated model parameters are often not open-sourced due to commercial considerations

^{*} Authors contributed equally.

[†] Corresponding author

³Such as the EU’s GDPR or the US’s HIPAA.

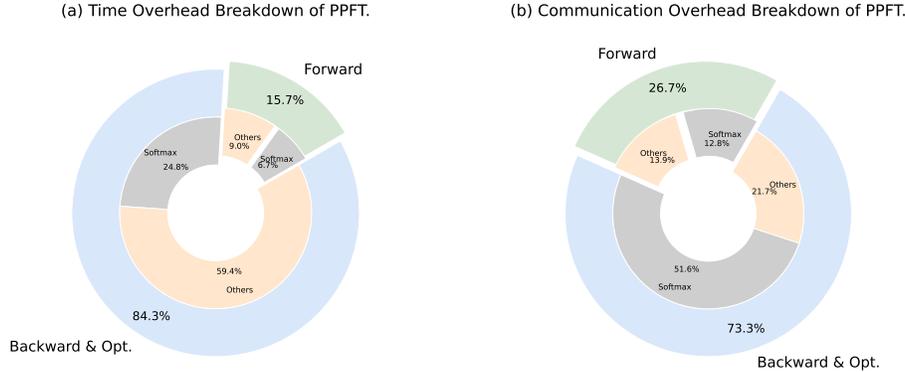


Figure 1: Breakdown of Time and Communication Overhead in Privacy-Preserving Fine-Tuning of RoBERTa_LARGE (24 layers, 1024 dimensions) based on MPC.

and the risk of privacy leaking⁴. Consequently, a critical bottleneck has emerged in the development of LLMs: the unavailability of private data and valuable proprietary model parameters significantly hampers LLMs’ improvement and applicability. This obstacle undermines collaborative research and innovation in the development of next-generation LLMs, particularly as high-quality public datasets are rapidly depleting [Villalobos et al., 2022].

Privacy-Preserving Machine Learning (PPML) [Mohassel and Zhang, 2017, Wagh et al., 2019] based on Multi-Party Computation (MPC) [Yao, 1986, Goldreich et al., 1987] offers a promising solution for secure collaborative model training. In this paradigm, model parameters and private data are firstly secret-shared among the involved parties. The parties then engage in a series of encrypted computations under a privacy-preserving environment by executing MPC protocols through multiple rounds of communication. All computations are performed over encrypted inputs, ensuring that each party learns only the outputs explicitly permitted by the protocol, without gaining access to private information such as the raw data or model parameters. Owing to its appealing privacy guarantee, MPC-based PPML has been successfully applied to the training and inference of various machine learning models, significantly advancing the deployment of privacy-preserving technologies in real-world applications [Hao et al., 2022, Luo et al., 2024, Pang et al., 2023, Lu et al., 2023].

Fine-tuning LLMs in sensitive domains like healthcare and finance demands robust privacy protections to safeguard user data, making the development of efficient MPC-based PPML solutions more essential. However, the substantial communication and computational overhead of MPC protocols (see Section 5.4) poses a significant barrier, severely impacting training efficiency compared to plaintext methods. For example, fine-tuning RoBERTa_LARGE [Liu et al., 2019] (with the batch size 1 on A100 GPUs) over a network with 3 Gbps bandwidth and 0.8 ms latency requires approximately 10 minutes and 838 GB of communication per batch update. These overheads, which scale with model size and architectural complexity, are particularly pronounced in transformer-based LLMs. As illustrated in Fig. 1, two primary factors dominate these overheads: a) *Backpropagation and optimization phases*, which account for up to 84.3% of communication time due to their intensive data exchanges, far exceeding the forward pass; and b) *Softmax operations in self-attention mechanisms*, which involve costly non-linear computations—exponentiation, maximum selection, and division—consuming up to 31.5% of total runtime under MPC. With increasing emphasis on protecting valuable model parameters and user data, there is an urgent need to devise efficient MPC-based methods to minimize these overheads and enable practical, privacy-preserving fine-tuning of LLMs.

These challenges necessitate tackling two pivotal research questions: (1) How can we minimize the significant computational and communication overhead of backpropagation in MPC-based LLM fine-tuning? (2) How can we design an efficient softmax approximation optimized for MPC frameworks? For (1), we propose parameter-efficient fine-tuning via Gradient-Free Optimization (GFO) [Rios and Sahinidis, 2013] within MPC settings, enabling parameter updates solely through forward passes, thus eliminating the costly backpropagation overhead. Unlike standard prompt tuning [Li and Liang, 2021, Liu et al., 2024], which requires caching all activation values for backpropagation, GFO streamlines

⁴This stems from model memorization during downstream training [Carlini et al., 2022].

the process, enhancing efficiency. For (2), we investigate linear attention mechanisms, specifically Random Feature Attention (RFA) [Peng et al., 2021, Choromanski et al., 2021], which circumvent the computational burden of softmax attention by offering linear time and space complexity, making them highly suitable for MPC environments.

This paper presents SECFWT, the *first privacy-preserving fine-tuning framework for LLMs* grounded in MPC environment. SECFWT efficiently updates model parameters just using forward passes, safeguarding the security of private data and trainable parameters. As depicted in Fig. 2, SECFWT comprises two core components: Forward-Only Tuning (FoT) and MPC-friendly Random Feature Attention. Based on GFO, FoT enhances training efficiency within MPC by employing an alternating plaintext–ciphertext computation scheme, eliminating the reliance on gradient-based optimization. To replace the Softmax operation in self-attention, SECFWT adopts RFA, which substantially reduces the computational overhead inherent in secure computation. Since the original RFA introduces nonlinear operations, particularly the cosine functions, we propose tailored optimizations that exploit trigonometric periodicity and sum-to-product formulas to enhance training efficiency.

Our Contributions: (1) We develop the first privacy-preserving fine-tuning framework SECFWT for LLMs training based on MPC, which efficiently updates model parameters just using forward passes. SECFWT ensures the confidentiality of both private data and trainable model parameters; (2) We introduce two novel components: MPC-friendly FoT and RFA. Extensive experiments demonstrate that SECFWT achieves a significant reduction in computation time (approximately **10.2×**) and communication cost (approximately **14.8×**), while incurring slight performance degradation (less than **3%**). Thus, this work pioneers the use of MPC for privacy-preserving LLMs fine-tuning, advancing the applicability of PPML in privacy-sensitive scenarios.

2 Related Work & Preliminaries

2.1 Privacy-preserving Machine Learning for LLMs

Existing privacy-preserving language models based on MPC and Homomorphic Encryption (HE) have primarily focused on improving inference efficiency by leveraging model architecture optimizations [Chen et al., 2022, Li et al., 2023, Zeng et al., 2022, Zhang et al., 2023, Liang et al., 2023] and privacy-preserving protocol enhancements [Hao et al., 2022, Zheng et al., 2023a, Gupta et al., 2023, Dong et al., 2023, Hou et al., 2023, Ding et al., 2023, Pang et al., 2023, Luo et al., 2024, Lu et al., 2023].

Compared to inference, fine-tuning large language models (LLMs) involves additional computational steps such as optimizer updates and backpropagation. These operations significantly increase the overall complexity and cost of privacy-preserving computation. As a result, only a limited number of studies in the PPML literature have explored secure fine-tuning of LLMs. Specifically, Lee et al. [2022] restricts fine-tuning to the classification head, while Li et al. [2024] focuses on reducing energy consumption during secure fine-tuning through hardware-level optimizations. A recent study [Rho et al., 2024a] proposes an efficient HE-based framework for privacy-preserving fine-tuning of language models, leveraging LoRA and Gaussian kernels to improve efficiency. While LoRA significantly reduces the number of tunable parameters, it does not eliminate the computational burden associated with backward and optimization, thereby limiting applicability to smaller models, such as two-layer BERT. Moreover, due to its HE implementation, non-linear operations within the language model must be approximated using polynomials, which inevitably degrades the fine-tuned model’s performance.

2.2 Softmax based Self-Attention & Random Feature Attention

2.2.1 Softmax based Self-Attention

Current large language models are constructed by stacking multiple Transformer layers. The core component of each Transformer layer is the self-attention mechanism. We omit a detailed discussion of the feed-forward network and other auxiliary components, as they remain unchanged in our work. Let n and d denote the sequence length and embedding dimension, respectively. The self-attention

mechanism is computed as follows:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} \right) \mathbf{V} \in \mathbb{R}^{n \times d}. \quad (1)$$

Here, the rows of Q , K , and V correspond to the query, key, and value vectors, respectively. The softmax function [Bridle, 1989] is applied row-wise, converting the similarity scores between each query and all key vectors into a probability distribution that weights the contribution of each value vector.

2.2.2 Random Feature Attention

To speed up the softmax operations in attention, recent work [Peng et al., 2021, Choromanski et al., 2021, Zheng et al., 2022] has employed random feature [Rahimi and Recht, 2007] methods to approximate the dot-then-exponentiate operation using kernel tricks. The main idea is to approximate the Gaussian kernel function via its Monte Carlo estimation:

$$\exp(-\|\mathbf{x} - \mathbf{x}'\|^2/\sigma^2) \approx \sum_{i=1}^M \varphi(\mathbf{x}, \omega_i) \varphi(\mathbf{x}', \omega_i), \quad (2)$$

where $\varphi(\mathbf{x}, \omega_i) = \sqrt{2/M} \cos(\omega_i^\top \mathbf{x} + b_i)$, with $\omega_i \sim \mathcal{N}(0, \sigma^2 I)$ and $b_i \sim U(0, 2\pi)$.

Let $\phi(\mathbf{x}) = \exp(\|\mathbf{x}\|^2/(2\sigma^2)) [\varphi(\mathbf{x}, \omega_1), \dots, \varphi(\mathbf{x}, \omega_M)]^\top$, the dot-then-exponentiate function can be approximated as:

$$\exp(\mathbf{x}^\top \mathbf{y}/\sigma^2) = \exp\left(\frac{1}{2\sigma^2} \|\mathbf{x}\|^2 + \frac{1}{2\sigma^2} \|\mathbf{y}\|^2\right) \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{y}\|^2\right) \approx \phi(\mathbf{x})^\top \phi(\mathbf{y}). \quad (3)$$

Substituting this approximation into the softmax attention, we obtain the RFA:

$$\begin{aligned} \text{softmax}(\mathbf{q}_t, \{\mathbf{k}_i\}_{i=1}^n, \{\mathbf{v}_i\}_{i=1}^n) &= \sum_i \frac{\exp(\mathbf{q}_t^\top \mathbf{k}_i/\sigma^2) \mathbf{v}_i^\top}{\sum_j \exp(\mathbf{q}_t^\top \mathbf{k}_j/\sigma^2)} \\ &\approx \sum_i \frac{\phi(\mathbf{q}_t)^\top \phi(\mathbf{k}_i) \mathbf{v}_i^\top}{\sum_j \phi(\mathbf{q}_t)^\top \phi(\mathbf{k}_j)} \\ &= \frac{\phi(\mathbf{q}_t)^\top \sum_i \phi(\mathbf{k}_i) \otimes \mathbf{v}_i}{\phi(\mathbf{q}_t)^\top \sum_j \phi(\mathbf{k}_j)} := \text{RFA}(\mathbf{q}_t, \{\mathbf{k}_i\}_{i=1}^n, \{\mathbf{v}_i\}_{i=1}^n), \end{aligned} \quad (4)$$

where $Q = \{\mathbf{q}_i\}_{i=1}^n$, $K = \{\mathbf{k}_i\}_{i=1}^n$, $V = \{\mathbf{v}_i\}_{i=1}^n$, and \otimes denotes the outer product between vectors. Leveraging this linearized formulation, RFA achieves linear time and memory complexity with respect to the sequence length.

2.3 Gradient-Free Optimization

Gradient-Free Optimization (GFO) [Rios and Sahinidis, 2013] algorithms focus on optimizing objective functions without relying on their gradients, instead utilizing only the function values (or fitness scores) of sampled solutions. Such methods are commonly referred to as black-box or zeroth-order optimization. The core framework of basic GFO algorithms follows a ‘‘sample-and-update’’ paradigm, iteratively searching for and estimating the optimum of the objective function. Since these algorithms do not require gradient information, they are particularly well-suited for problems where derivatives are difficult to obtain or computationally expensive to calculate.

Prior work Black-Box Tuning (BBT) Sun et al. [2022b] pioneered the application of GFO to LLMs prompt-tuning, achieving notable success. BBT aims to train a continuous prompt vector $p \in \mathbb{R}^D$ such that the prediction performance improves when the model is fed the optimal prompt vector p^* along with the input X . Mathematically,

$$p^* = \arg \min_{p \in \mathcal{P}} \mathcal{L}(f(p; X), Y), \quad (5)$$

where $f(\cdot)$ denotes the inference function of LLMs, $\mathcal{L}(\cdot)$ is the loss function, and \mathcal{P} represents the search space of the prompt vectors. However, the convergence speed of GFO algorithms tends

to degrade as the dimensionality of the search space increases. Motivated by the observation that LLMs exhibit low intrinsic dimensionality, BBT optimizes a lower-dimensional variable $z \in \mathbb{R}^d$ with $d \ll D$ in a smaller subspace, and projects it back to the original prompt space \mathcal{P} using a random projection matrix $A \in \mathbb{R}^{D \times d}$. The optimization objective becomes:

$$z^* = \arg \min_{z \in \mathcal{Z}} \mathcal{L}(f(Az; X), Y). \quad (6)$$

To efficiently search z , BBT uses the gradient-free optimizer Covariance Matrix Adaptation Evolution Strategy (CMA-ES) [Hansen, 2016]. Although BBT can protect the parameters of cloud-based LLMs by invoking their APIs as a black-box, it still requires users to upload their data to the cloud server, which is prohibited due to privacy and legal constraints. Directly integrating BBT into MPC appears to enable privacy-preserving fine-tuning of LLMs. Nevertheless, softmax operations present a challenge for practical BBT deployment in MPC settings due to their intrinsic nonlinearity.

2.4 2-out-of-2 Arithmetic Secret Sharing

For an integer ring $\mathbb{Z}_L = \{0, 1, \dots, L - 1\}$, a 2-out-of-2 arithmetic secret sharing scheme involves the following two algorithms:

- The sharing algorithm $Shr(x) \rightarrow ([x]_0, [x]_1)$ is used to generate the shares of x . Specifically, a value r is chosen uniformly at random from \mathbb{Z}_L , such that $[x]_0 = r$, and $[x]_1 = x - r \pmod{L}$ is computed.
- The reconstruction algorithm $Rec([x]_0, [x]_1) \rightarrow x$ is used to reconstruct x , i.e., $x = [x]_0 + [x]_1 \pmod{L}$.

Since the value r is sampled uniformly at random from \mathbb{Z}_L , neither share $[x]_0$ nor $[x]_1$ reveals any information about the secret x . We denote the arithmetic secret sharing of x as $[x] = ([x]_0, [x]_1)$.

In the field of secure MPC, numerous secure protocols have been developed for operating over secret shares $[x]$, including secure addition, multiplication, comparison, and various nonlinear activation functions. These cryptographic primitives are summarized in Section 5.4. In this work, we treat these primitives as black-box components and utilize them without requiring additional assumptions or modifications.

3 SECFWT

SECFWT adopts a three-party semi-honest (honest-but-curious) security model [Wagh et al., 2019, Knott et al., 2021], consisting of two non-colluding computation servers and a trusted dealer responsible for generating correlated randomness. During privacy-preserving fine-tuning, the servers are assumed to follow the protocol faithfully, but may attempt to infer unauthorized information from intermediate computation results. The privacy guarantees of SECFWT are inherited from the underlying secure computation framework.

3.1 Overview

As illustrated in Fig. 2, SECFWT *consists of the following six steps*: (1) Prior to privacy-preserving fine-tuning, the dataowner locally generates secret shares of the private training data, which are then distributed to the respective computation servers. (2) During fine-tuning, the two servers perform privacy-preserving inference using the secret-shared model parameters, which de and private inputs, obtaining secret shares of the inference results. (3) Each server sends its share of the inference results to the client. (4) Upon receiving both shares, the client reconstructs the correct inference result. (5) The client then computes the loss locally in plaintext using the corresponding ground-truth labels. (6) Finally, the client performs gradient-free optimization to update the prompt vector using the computed loss, iterating this process multiple times to obtain the final fine-tuned prompt.

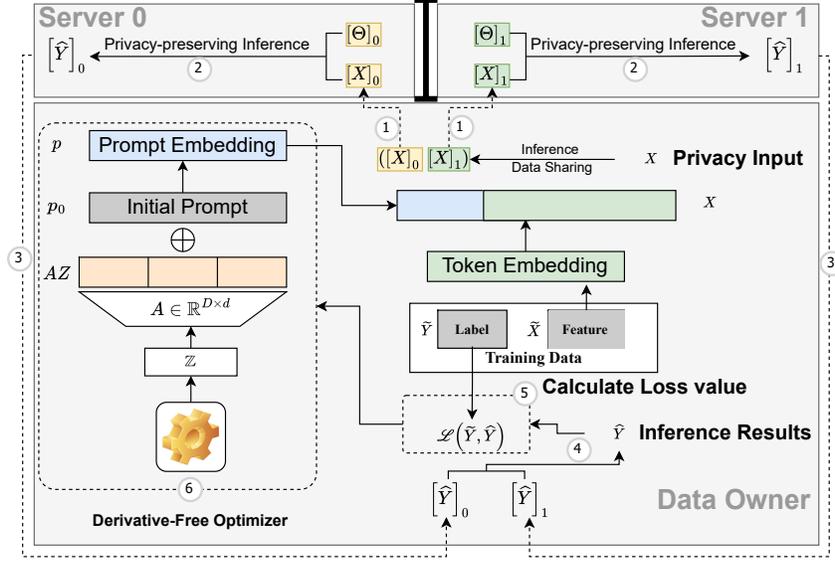


Figure 2: Overview of SECFWT. SECFWT leverages secure MPC to protect both training data and model parameters during fine-tuning. It addresses three key bottlenecks in privacy-preserving fine-tuning. First, it eliminates the computational overhead of backward and optimizer by adopting a FoT paradigm. Second, it mitigates the risk of privacy leakage caused by parameter updates by using derivative-free optimizer. Third, it improves the efficiency of privacy-preserving attention computation in LLMs by employing RFA.

3.2 Speedup With Forward-Only Tuning

3.2.1 Major Bottlenecks of Full Fine-Tuning

Full-parameter fine-tuning under MPC incurs prohibitive overhead. Specifically, conventional pre-trained language models (LLMs) typically involve three core stages: (1) forward propagation, (2) optimizer computation, and (3) backpropagation. In the optimizer computation stage, loss functions such as cross-entropy and maximum likelihood estimation involve numerous non-linear operations that are not well-suited for MPC. These non-linearities, coupled with the gradient update steps, lead to privacy-preserving optimizer computation costs that can even exceed the computation cost of a single Transformer layer Rho et al. [2024b].

Backpropagation, based on the chain rule, computes gradients by combining the derivative of activation functions with the weights of each layer. In deep neural networks with multiple hidden layers, this process starts from the output layer and proceeds backward to the input layer. Like the optimizer, it introduces a large number of MPC-unfriendly non-linear operations. Taken together, these factors result in significant computational and communication overhead, making full fine-tuning of LLMs under MPC highly inefficient and impractical.

3.2.2 Speedup Full Fine-tuning With Forward-Only Tuning

Forward-only tuning eliminates the need for gradient computation by executing only the forward pass, thereby fundamentally avoiding the high privacy-preserving computation costs associated with the complex non-linear operations in traditional full fine-tuning—particularly those arising from optimizer computation and backpropagation. Specifically, FoT computes a loss using the model’s inference results and the corresponding fine-tuning labels, and then feeds this loss into a derivative-free optimizer to update the prompt vector.

However, while derivative-free optimization avoids the overhead of gradient computation and backpropagation, it still incurs costs from loss computation and introduces additional privacy-preserving overhead due to the optimizer itself. To address this challenge, SECFWT introduces an innovative interaction paradigm in which encrypted communication between the servers and the data owner is leveraged. Under this paradigm, the data owner locally performs loss evaluation and derivative-free

optimization in plaintext, enabling high efficiency while maintaining privacy. By offloading these computations to the client side, SECFWT significantly reduces the privacy-preserving computational burden and provides a practical and efficient solution for fine-tuning large language models in privacy-sensitive settings.

3.3 Speedup With Random Feature Attention

3.3.1 Major Bottlenecks of Softmax-Based Attention

For a vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$, Softmax in Transformer converts it to an n -dimensional probability distribution with

$$\text{Softmax}(\mathbf{x})[i] = \frac{e^{x_i - \tau}}{\sum_{h=1}^n e^{x_h - \tau}}, \quad (7)$$

where $\tau = \max(\{x_h\}_{h=1}^n)$ is used to ensure stable numerical computations.

Bottleneck 1: Quadratic complexity with respect to sequence length. Given $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d}$, where n denotes the sequence length and d the embedding dimension, the complexity of Softmax-based attention scales as $O(n^2d)$. This quadratic dependence becomes prohibitively expensive for long input sequences.

Bottleneck 2: Numerous nonlinear operations incompatible with privacy-preserving computation. As shown in Eq. (7), computing the Softmax function involves three nonlinear operations—exponentiation, division, and maximization—all of which are costly to implement under MPC. These operations significantly inflate the overhead of privacy-preserving attention computation (see Appendix B for details), rendering Softmax-based attention inefficient in secure settings.

3.3.2 Speed Up Softmax-based Attention with Random Feature Attention

Since the softmax function in standard attention suffers from both high computational complexity and nonlinear operations, it is challenging to develop a privacy-preserving approximation. To address this issue, this paper employs RFA [Peng et al., 2021, Choromanski et al., 2021] as an alternative approximation to softmax. Compared to existing softmax approximation methods [Kitaev et al., 2020, Wang et al., 2020, Roy et al., 2021], RFA offers the following advantages:

- **Theoretical Guarantee on Approximation Error.** The approximation error is formally bounded, ensuring reliable accuracy.
- **Reduction in Computational Complexity of Softmax.** RFA reduces the complexity of softmax attention from $O(n^2d)$ to $O(ndr)$, where r represents the number of random features used.
- **Avoidance of Exponentiation and Maximum Operations in Softmax.** By bypassing these costly nonlinear operations, RFA significantly improves efficiency in privacy-preserving settings.

MPC-Friendly RFA Design. We present the specific design of the Privacy-preserving RFA based on MPC. According to Eq. (4), the computation of RFA involves multiplication, division, and cosine function operations. This implies that although RFA bypasses the exponential and maximum operations in softmax-based attention, it introduces cosine operations that are not friendly to MPC.

To address this challenge, we design an efficient MPC-based privacy-preserving cosine function protocol (Π_{cosine}) by leveraging the periodicity of trigonometric functions and the sum-to-product formulas. Detailed algorithmic descriptions are provided in Section 5.2. By executing Π_{cosine} , the privacy-preserving computation of the cosine function can be accomplished with only a single round of communication, transmitting 2ℓ -bit elements. Building upon this result, we further develop an efficient MPC-based privacy-preserving RFA protocol, which reduces the computational complexity of the Softmax-based attention mechanism while circumventing the need for expensive exponentiation and maximum operations.

4 Experiments

In this section, we present experimental results to demonstrate the effectiveness of SECFWT. Specifically, we provide a brief overview of the experimental setup in Section 4.1. Then, we report the efficiency and performance results of SECFWT in Section 4.5 and Section 4.2, respectively.

4.1 Experimental Setup

MPC-Backend & Testbeds. Our implementation is based on the MPC framework CrypTen⁵. We conduct our experimental evaluations on three servers, each equipped with an A100 GPU. To enable a comprehensive efficiency comparison, we utilize *Linux Traffic Control (TC)* to simulate various network conditions. Specifically: In the local area network (LAN) scenario, we set the bandwidth to 3 Gbps with a round-trip latency of 0.8 ms. For the wide area network (WAN) setting, we consider two different configurations: {100 Mbps, 80 ms} and {200 Mbps, 40 ms}, simulating diverse WAN environments.

Models and Datasets. To maintain consistency with prior work on FoT [Sun et al., 2022b,a], we adopt RoBERTa_{LARGE} [Liu et al., 2019] as our backbone model. It is worth noting that SECFWT is compatible with other Transformer-based large language models, such as GPT [Radford et al., 2019, Brown et al., 2020], T5 [Raffel et al., 2020], and BART [Lewis et al., 2020]. We leave privacy-preserving fine-tuning of these models to future work. For datasets, we follow standard practice and use the GLUE benchmark, which consists of representative tasks for evaluating natural language understanding.

Baseline. To demonstrate the effectiveness of SECFWT, we establish several baselines for comparison. Specifically, for performance evaluation, we consider the following four baselines:

- Full Fine-tuning with Softmax Attention (Full+SM): We fully fine-tune RoBERTa_{LARGE} on downstream task data using the standard softmax-based attention mechanism.
- Full Fine-tuning with Random Feature Attention (Full+RFA): We replace the softmax-based attention in RoBERTa_{LARGE} with RFA and then conduct full fine-tuning on downstream tasks.
- Forward-Only Tuning with Softmax Attention (FoT+SM): We perform FoT on RoBERTa_{LARGE} using downstream task data with the standard softmax-based attention.

For efficiency evaluation, we exclude Manual Prompting as it incurs no computational or communication overhead. Thus, we assess the full pipeline runtime and communication cost for the remaining three baselines.

4.2 Efficiency Comparison

We now evaluate the efficiency improvements achieved by our proposed scheme. Specifically, we first assess the performance gains brought by the privacy-preserving RFA mechanism. We then further evaluate the acceleration obtained when performing privacy-preserving fine-tuning using SECFWT.

4.3 Privacy-preserving RFA Efficiency

We compare Π_{PPRFA} with the conventional softmax-based privacy-preserving self-attention mechanism under varying sequence lengths, as illustrated in Section 4.2. The results clearly demonstrate that Π_{PPRFA} significantly reduces both computational and communication overhead compared to its softmax-based counterpart. Moreover, this advantage grows rapidly as the input sequence length increases, owing to the quadratic computational complexity of softmax-based attention, in contrast to the linear complexity of RFA.

4.4 End-to-end Efficiency

We evaluate the efficiency of SECFWT by measuring the time and communication cost required for fine-tuning on a single sample. The detailed results are presented in Table 1, with additional

⁵<https://github.com/facebookresearch/CrypTen>

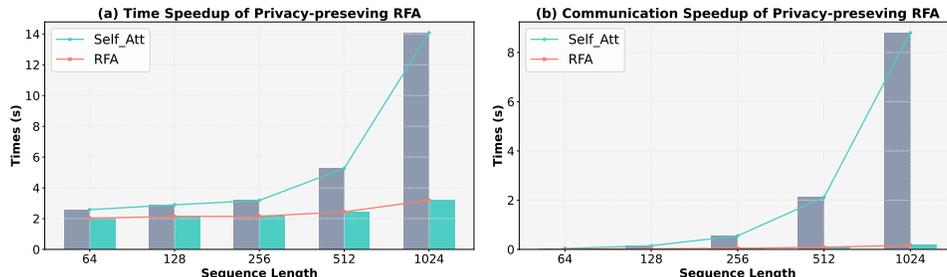


Figure 3: Comparison of Time and Communication Overhead Between Privacy-Preserving RFA and Softmax-Based Privacy-Preserving Attention.

Table 1: Efficiency Comparison of RoBERTa_{LARGE} in LAN Setting (3Gbps, 0.8ms). Bolded numbers indicate the best results. The results are the average of ten runs.

Methods	Forward		Backward		Optimizer		Total	
	Times(s)	Comm(GB)	Times(s)	Comm(GB)	Times(s)	Comm(GB)	Times(s)	Comm(GB)
Full + SM	186.966	259.091	364.228	559.721	13.058	18.964	564.253	837.776
FoT + SM	173.999	205.358	0.000	0.000	0.138	0.000	174.138	205.359
Full + RFA	92.693	60.426	12.234	59.475	12.595	18.964	84.859	138.865
SECFWT	54.17	56.545	0.000	0.000	0.103	0.000	55.172	56.545

evaluations under wide-area network (WAN) settings provided in Section 5.3. As shown in the table, SECFWT significantly reduces the time and communication overhead of privacy-preserving full fine-tuning. Specifically, by adopting FoT to eliminate the computational overhead associated with privacy-preserving optimizer operations and backpropagation, the fine-tuning time is improved by a factor of $3.24\times$, and the communication cost is reduced by $4.08\times$. On top of that, SECFWT further improves efficiency by leveraging RFA, ultimately achieving a $10.23\times$ speedup in time and a $14.82\times$ reduction in communication cost.

We further evaluate the time and communication improvements of SECFWT on downstream tasks, with detailed results presented in Table 2. Compared to full-parameter fine-tuning, SECFWT achieves comparable performance while delivering a $4.55\text{--}4.67\times$ speedup and a $12.55\times$ reduction in communication overhead for privacy-preserving fine-tuning across downstream tasks.

Table 2: We evaluate the feasibility of Fine-tuning-as-a-Service (FaaS) by comparing the end-to-end fine-tuning time, communication overhead, and the total amount of data uploaded/downloaded for completing private fine-tuning on the SST-2 and AG News datasets.

	FaaS	Training Time	Communication Volume	Upload (per query)	Download (per query)
SST-2 (number of classes: 2)					
Model Tuning	×	110.5 (h)	2838.81 (TB)	-	-
SECFWT	✓	24.3 (h)	226.18 (TB)	24 MB	0.5 KB
AG’s News (number of classes: 4)					
Model Tuning	×	212.6 (h)	5677.63 (TB)	-	-
SECFWT	✓	45.5 (h)	452.36 (TB)	54 MB	2 KB

4.5 Performance Comparison

We evaluate the performance of SECFWT on the GLUE benchmark. The detailed results are shown in Table Table 3. Experimental results demonstrate that SECFWT achieves comparable performance to full fine-tuning. Specifically, on binary classification tasks such as SST-2, MRPC, and Yelp P., SECFWT matches or even outperforms full fine-tuning—for instance, achieving better results on MRPC. However, for more challenging tasks such as the four-class classification in AG’s News, SECFWT suffers from a more noticeable performance drop. In practice, this gap can be mitigated by increasing the number of optimization iterations during tuning.

Table 3: Comprehensive performance comparison of SECFWT across various language understanding tasks. The results in the table report the mean and standard deviation over three runs. All experiments are conducted using the pretrained RoBERTa_{LARGE} model with 16 samples per class.

Method	SST-2 Acc	Yelp P. Acc	AG’s News F1	MRPC Acc	RTE Acc	Avg.
Full + SM	85.39 ± 2.84	91.82 ± 0.79	86.36 ± 1.85	77.25 ± 5.70	58.60 ± 6.21	79.88
Full + RFA	85.33 ± 3.85	90.63 ± 0.93	85.32 ± 2.55	75.22 ± 6.50	57.30 ± 6.58	78.76
FoT + SM	84.13 ± 0.16	89.70 ± 0.23	77.23 ± 0.23	79.47 ± 0.23	53.57 ± 0.23	76.82
SECFWT	83.83 ± 0.26	87.50 ± 1.52	76.13 ± 2.32	78.82 ± 2.32	52.22 ± 1.12	75.70

5 Conclusion

This paper proposes the first comprehensive MPC-based framework for privacy-preserving fine-tuning, named SECFWT. By adopting a FoT approach, it eliminates the need for backpropagation and optimizer computations, and leverages RFA to improve the efficiency of privacy-preserving attention computations. Experimental results demonstrate that SECFWT significantly enhances the efficiency of privacy-preserving fine-tuning, making scalable fine-tuning of LLMs feasible in privacy-sensitive applications.

References

- John Bridle. Training stochastic model recognition algorithms as networks can lead to maximum mutual information estimation of parameters. *Advances in neural information processing systems*, 2, 1989.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Ran Canetti. Universally composable security: A new paradigm for cryptographic protocols. In *Proceedings 42nd IEEE Symposium on Foundations of Computer Science*, pages 136–145. IEEE, 2001.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- Tianyu Chen, Hangbo Bao, Shaohan Huang, Li Dong, Binxing Jiao, Daxin Jiang, Haoyi Zhou, Jianxin Li, and Furu Wei. THE-X: Privacy-preserving transformer inference with homomorphic encryption. In *Findings of the Association for Computational Linguistics*, pages 3510–3520, 2022.
- Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. In *9th International Conference on Learning Representations, ICLR 2021*, 2021.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 4171–4186, 2019.
- Yuanchao Ding, Hua Guo, Yewei Guan, Weixin Liu, Jiarong Huo, Zhenyu Guan, and Xiyong Zhang. East: Efficient and accurate secure transformer framework for inference. *arXiv preprint arXiv:2308.09923*, 2023.
- Ye Dong, Wen-jie Lu, Yancheng Zheng, Haoqi Wu, Derun Zhao, Jin Tan, Zhicong Huang, Cheng Hong, Tao Wei, and Wenguang Cheng. PUMA: Secure inference of LLaMA-7B in five minutes. *arXiv preprint arXiv:2307.12533*, 2023.
- Oded Goldreich, Silvio Micali, and Avi Wigderson. How to play any mental game or A completeness theorem for protocols with honest majority. In *Proceedings of the 19th Annual ACM Symposium on Theory of Computing*, pages 218–229. ACM, 1987.

- Kanav Gupta, Neha Jawalkar, Ananta Mukherjee, Nishanth Chandran, Divya Gupta, Ashish Panwar, and Rahul Sharma. SIGMA: Secure GPT inference with function secret sharing. *Cryptology ePrint Archive, Paper 2023/1269*, 2023.
- Nikolaus Hansen. The cma evolution strategy: A tutorial. *arXiv preprint arXiv:1604.00772*, 2016.
- Meng Hao, Hongwei Li, Hanxiao Chen, Pengzhi Xing, Guowen Xu, and Tianwei Zhang. Iron: Private inference on transformers. *Advances in Neural Information Processing Systems*, 35:15718–15731, 2022.
- Xiaoyang Hou, Jian Liu, Jingyu Li, Yuhan Li, Wen jie Lu, Cheng Hong, and Kui Ren. CipherGPT: Secure two-party GPT inference. *Cryptology ePrint Archive, Paper 2023/1147*, 2023.
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *8th International Conference on Learning Representations, ICLR 2020*, 2020.
- Brian Knott, Shobha Venkataraman, Awni Hannun, Shubho Sengupta, Mark Ibrahim, and Laurens van der Maaten. CrypTen: Secure multi-party computation meets machine learning. *Advances in Neural Information Processing Systems*, 34:4961–4973, 2021.
- Garam Lee, Minsoo Kim, Jai Hyun Park, Seung-won Hwang, and Jung Hee Cheon. Privacy-preserving text classification on bert embeddings with homomorphic encryption. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3169–3175, 2022.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, 2020.
- Dacheng Li, Rulin Shao, Hongyi Wang, Han Guo, Eric P Xing, and Hao Zhang. MPCFormer: Fast, performant and private transformer inference with MPC. In *Proceedings of the Eleventh International Conference on Learning Representations, ICLR, 2023*.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- Yang Li, Wenhan Yu, and Jun Zhao. Privtuner with homomorphic encryption and lora: A p3eft scheme for privacy-preserving parameter-efficient fine-tuning of ai foundation models. *arXiv preprint arXiv:2410.00433*, 2024.
- Zi Liang, Pinghui Wang, Ruofei Zhang, Nuo Xu, and Shuo Zhang. MERGE: Fast private text generation. *arXiv preprint arXiv:2305.15769*, 2023.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt understands, too. *AI Open*, 5:208–215, 2024.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Wen-jie Lu, Zhicong Huang, Zhen Gu, Jingyu Li, Jian Liu, Cheng Hong, Kui Ren, Tao Wei, and WenGuang Chen. Bumblebee: Secure two-party inference framework for large transformers. *Cryptology ePrint Archive*, 2023.
- Jinglong Luo, Yehong Zhang, Jiaqi Zhang, Xin Mu, Hui Wang, Yue Yu, and Zenglin Xu. Secformer: Towards fast and accurate privacy-preserving inference for large language models. *arXiv preprint arXiv:2401.00793*, 2024.
- Payman Mohassel and Yupeng Zhang. SecureML: A system for scalable privacy-preserving machine learning. In *Proceedings of 2017 IEEE Symposium on Security and Privacy*, pages 19–38. IEEE, 2017.
- OpenAI. GPT-4 technical report. *ArXiv*, abs/2303.08774, 2023.

- Qi Pang, Jinhao Zhu, Helen Möllering, Wenting Zheng, and Thomas Schneider. BOLT: Privacy-preserving, accurate and efficient inference for transformers. *Cryptology ePrint Archive, Paper 2023/1893*, 2023.
- Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah A Smith, and Lingpeng Kong. Random feature attention. In *9th International Conference on Learning Representations, ICLR 2021*, 2021.
- Chen Qu, Weize Kong, Liu Yang, Mingyang Zhang, Michael Bendersky, and Marc Najork. Natural language understanding with privacy-preserving bert. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 1488–1497, 2021.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, 2007.
- Donghwan Rho, Taeseong Kim, Minje Park, Jung Woo Kim, Hyunsik Chae, Ernest K Ryu, and Jung Hee Cheon. Encryption-friendly llm architecture. *arXiv preprint arXiv:2410.02486*, 2024a.
- Donghwan Rho, Taeseong Kim, Minje Park, Jung Woo Kim, Hyunsik Chae, Ernest K Ryu, and Jung Hee Cheon. Encryption-friendly llm architecture. *arXiv preprint arXiv:2410.02486*, 2024b.
- Luis Miguel Rios and Nikolaos V Sahinidis. Derivative-free optimization: a review of algorithms and comparison of software implementations. *Journal of Global Optimization*, 56(3):1247–1293, 2013.
- Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse attention with routing transformers. *Transactions of the Association for Computational Linguistics*, 9:53–68, 2021.
- Tianxiang Sun, Zhengfu He, Hong Qian, Yunhua Zhou, Xuan-Jing Huang, and Xipeng Qiu. BBTv2: Towards a gradient-free future with large language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3916–3930, 2022a.
- Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. Black-box tuning for language-model-as-a-service. In *International Conference on Machine Learning*, pages 20841–20855. PMLR, 2022b.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- Pablo Villalobos, Jaime Sevilla, Lennart Heim, Tamay Besiroglu, Marius Hobbhahn, and Anson Ho. Will we run out of data? An analysis of the limits of scaling datasets in machine learning. *arXiv e-prints*, 2022.
- Sameer Wagh, Divya Gupta, and Nishanth Chandran. SecureNN: 3-Party secure computation for neural network training. *Proceedings on Privacy Enhancing Technologies*, pages 26–49, 2019.
- Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- Andrew Chi-Chih Yao. How to generate and exchange secrets. In *Annual Symposium on Foundations of Computer Science*, pages 162–167, 1986.
- Wenxuan Zeng, Meng Li, Wenjie Xiong, Wenjie Lu, Jin Tan, Runsheng Wang, and Ru Huang. MPCViT: Searching for MPC-friendly vision transformer with heterogeneous attention. *arXiv preprint arXiv:2211.13955*, 2022.

- Yuke Zhang, Dake Chen, Souvik Kundu, Chenghao Li, and Peter A Beereel. SAL-ViT: Towards latency efficient private inference on ViT using selective attention search with a learnable softmax approximation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5116–5125, 2023.
- Lin Zheng, Chong Wang, and Lingpeng Kong. Linear complexity randomized self-attention mechanism. In *International conference on machine learning*, pages 27011–27041, 2022.
- Mengxin Zheng, Qian Lou, and Lei Jiang. Primer: Fast private transformer inference on encrypted data. *arXiv preprint arXiv:2303.13679*, 2023a.
- Yu Zheng, Qizhi Zhang, Sherman SM Chow, Yuxiang Peng, Sijun Tan, Lichun Li, and Shan Yin. Secure softmax/sigmoid for machine-learning computation. In *Proceedings of the 39th Annual Computer Security Applications Conference*, pages 463–476, 2023b.

Appendices

The appendix is organized as follows. In Section 5.1, we discuss the limitations of SECFWT. In Section 5.2, we detail the privacy-preserving algorithms designed for SECFWT, including privacy-preserving cosine similarity and Random Feature Attention (RFA). In Section 5.3, we report the runtime overhead of executing SECFWT in a wide-area network (WAN) setting. In Section 5.4, we present the underlying MPC protocols upon which SECFWT is built. Finally, in Section 5.5, we provide a comprehensive security proof of SECFWT.

5.1 Limitation

As the first work to implement privacy-preserving fine-tuning of LLMs using MPC, SECFWT is currently evaluated only on the RoBERTa_{LARGE} model for natural language understanding tasks. We leave the extension of our framework to other Transformer-based LLMs—such as GPT (Brown et al., 2020), T5 (Raffel et al., 2020), and BART (Lewis et al., 2020)—for future work.

Moreover, we do not incorporate other optimization techniques that are orthogonal to SECFWT. In particular, we leave the secure handling of additional nonlinear operations in LLMs—such as GeLU and LayerNorm—as future work.

5.2 Privacy-preserving Protocols

Privacy-preserving Cosine. We propose an efficient privacy-preserving cosine protocol Π_{cosine} by exploiting the periodicity of the cosine function and trigonometric addition identity formulas. Here’s a detailed description of the algorithm steps: In the offline phase, the protocol initiates by generating pseudo-random values. Specifically, S_0 and the trusted third party T jointly produce $[t]_0, [u]_0, [v]_0$ by evaluating a pseudo-random function (PRF) with a specific key k_0 . Similarly, S_1 and T generate $[t]_1$ using a different key k_1 . Then, the trusted third party T recover the actual value $t = [t]_0 + [t]_1$, calculates $[u]_1 = \sin(t) - [u]_0$ and $[v]_1 = \cos(t) - [v]_0$. This phase is crucial for preparing necessary correlated randomness that will be used in the online phase.

In the online phase, the parties compute the $[\sin(x)]$ securely. First, each party S_j computes $[\delta]_j = [x]_j + [t]_j \pmod{20}$. Then, through one round of communication, the parties reconstruct δ by exchanging $[\delta]_0$ and $[\delta]_1$. Subsequently, we get $p = \sin(\delta)$ and $q = \cos(\delta)$. Finally, each party calculates $[y]_j = p[u]_j + q[v]_j$. This effectively leverages the precomputed correlated randomness with the current input $[x]$ to produce the $[\sin(x)]$ in a privacy-preserving manner. The Π_{cosine} requires only one round of communication during the online phase, with a communication cost of transmitting 2ℓ elements.

Algorithm 1: Protocol for Privacy-preserving Cosine Π_{cosine}

Input: For $j \in \{0, 1\}$, S_j holds the shares $[x]_j$; Pseudo-Random Function (PRF), and key k_j .

Output: For $j \in \{0, 1\}$, S_j returns the shares $[y]_j$, where $y = \sin(x)$.

```
/* Offline Phase */
1  $S_0, T : [t]_0, [u]_0, [v]_0 \leftarrow PRF(k_0)$ 
2  $S_1, T : [t]_1 \leftarrow PRF(k_1)$ 
3  $T : t = [t]_0 + [t]_1, [u]_1 = \sin(t) - [u]_0, [v]_1 = \cos(t) - [v]_0$ 
/* Online Phase */
4  $[\delta]_j = [x]_j + [t]_j \pmod{20}$ 
5  $\delta = [\delta]_0 + [\delta]_1$  // reconstruct  $\delta$  by 1 round of communication
6  $p = \sin(\delta), q = \cos(\delta)$ 
7  $[y]_j = p[u]_j + q[v]_j$ 
```

Privacy-preserving Feature Attention. The Privacy-preserving RFA Protocol (Π_{RFA}) is designed to enable computation of RFA with privacy preservation. The algorithm involves two parties, S_0 and S_1 , and a trusted third party T , to collaboratively compute the RFA while keeping the input data secure. In the offline phase, the protocol begins with the generation of pseudo-random values. Specifically, S_0 and the trusted third party T jointly produce $[t]_0, [u]_0, [v]_0$ by evaluating a PRF with a random seed r_0 , and also generate matrix W using another random seed r . On the other hand, S_1

and the trusted third party T generate $[t]_1$ by evaluating the PRF with a different random seed r_1 , and use the same matrix W generated earlier. Then, the trusted third party T recovers the actual value $t = [t]_0 + [t]_1$. Based on t , T computes $[u]_1 = \sin(t) - [u]_0$ and $[v]_1 = \cos(t) - [v]_0$. This offline phase essentially prepares some necessary random values and parameters, which will be used in the online phase. Although these values are related to trigonometric functions, they are computed in a way that preserves privacy as the actual values are hidden within the shares.

In the online phase, the algorithm focuses on computing the attention mechanism. First, for each query q_t at time step t and key k_i , the corresponding feature mappings are computed. This is done by taking the shares of q_t and k_i (i.e., $[q]_t$ and $[k]_i$) and applying a cosine-based transformation denoted as Π_{\cosine} , scaled by a factor of $\sqrt{2}/r$. The scaling factor is important to ensure proper normalization of the feature mappings.

Next, for each key-value pair (k_i, v_i) , the share $[z]_j$ is computed as the element-wise product (denoted by \otimes) between the feature - mapped key $[\phi(k_i)]_j$ and the value v_i . This effectively combines the key's feature representation with its associated value.

Then, the attention score $[s]_j$ is calculated as the dot product between the feature-mapped query $[\phi(q_t)]_j$ and the feature-mapped key $[\phi(k_i)]_j$. This dot product represents the similarity between the query and the key in the transformed feature space.

Finally, the output share $[y]_j$ is obtained by dividing $[z]_j$ by $[s]_j$. This step normalizes the combined key-value representation by the attention score, resulting in the weighted value that will be used as the output of the attention mechanism. The division here is crucial as it implements the attention-weighting process, where the value is scaled according to how relevant the corresponding key is to the query.

Algorithm 2: Privacy-preserving RFA Protocol (Π_{RFA})

Input: For $j \in \{0, 1\}$, S_j holds the shares $\{[q]_t, [k]_i, [v]_i\}$;

Output: For $j \in \{0, 1\}$, S_j returns the shares $[y]_j$, where $y = RFA([q]_t, [k]_i, [v]_i)$.

/* Offline Phase */

1 $S_0, T : [t]_0, [u]_0, [v]_0 \leftarrow PRF(r_0); W \leftarrow PRF(r)$

2 $S_1, T : [t]_1 \leftarrow PRF(r_1); W \leftarrow PRF(r)$

3 $T : t = [t]_0 + [t]_1, [u]_1 = \sin(t) - [u]_0, [v]_1 = \cos(t) - [v]_0$

/* Online Phase */

4 $[\phi(q_t)]_j = \sqrt{\frac{2}{r}} \Pi_{\cosine}(W[q_t]_j); [\phi(k_i)]_j = \sqrt{\frac{2}{r}} \Pi_{\cosine}(W[k_i]_j)$

5 $[z]_j = [\phi(k_i)]_j \otimes v_i$

6 $[s]_j = [\phi(q_t)]_j \cdot [\phi(k_i)]_j$

7 $[y]_j = [z]_j / [s]_j$

5.3 More Results

In this section, we present additional efficiency results of SECFWT. Specifically, reports the time overhead of SECFWT in a wide-area network (WAN) setting.

5.4 Underlying MPC Protocols

In this section, we provide a brief overview of the underlying protocols used and refer to the works of Knott et al. [2021] and Zheng et al. [2023b] for details. Let S_j with $j \in \{0, 1\}$ be two parties that are used to execute the MPC protocol. Each party S_j will be given one additive share $([u]_j, [v]_j) \in \mathcal{Z}_L$ of the operation inputs u and v for $j \in \{0, 1\}$.

5.4.1 Privacy-Preserving Linear Protocols

Privacy-preserving addition. It is implemented locally as $[u + v]_j = [u]_j + [v]_j$ for $j \in \{0, 1\}$.

Privacy-preserving multiplication. The multiplication is implemented with Beaver-triples: (a, b, c) where $a, b \in \mathcal{Z}_L$ are randomly sampled from \mathcal{Z}_L and $c = a \cdot b \pmod L$. Specifically, for each

Table 4: Efficiency Comparison of RoBERTa_{LARGE} in WAN Setting. Bolded numbers indicate the best results. The results are the average of ten runs.

Bandwidth & Latency	Methods	Forward (s)	Backward (s)	Optimizer (s)	Total (s)
200Mbps/40ms	Full + SM	190.815	409.07	13.217	613.103
	FoT + SM	180.902	0.000	0.249	181.151
	Full + RFA	110.512	13.935	17.105	141.553
	SECFWT	86.947	0.000	0.103	87.051
100Mbps/80ms	Full + SM	198.152	410.478	16.97	625.601
	FoT + SM	195.805	0.000	0.023	195.829
	Full + RFA	108.206	18.917	15.934	143.058
	SECFWT	89.074	0.000	0.125	89.179

$j \in \{0, 1\}$, S_j first calculates $[d]_j = [u]_j - [a]_j$ and $[e]_j = [v]_j - [b]_j$. Then, they send the $[d]_j$ and $[e]_j$ to each other and reconstruct $d = \text{Rec}([d]_0, [d]_1)$ and $e = \text{Rec}([e]_0, [e]_1)$. Finally, the additive share of $u \cdot v$ can be computed using $[u \cdot v]_j = -jd \cdot e + [u]_j \cdot e + d \cdot [v]_j + [c]_j$. To complete the SS-based multiplication, both parties need to spend 1 round of two-way communication and transmit 256 bits.

5.4.2 Privacy-Preserving Non-Linear Protocols

Privacy-preserving comparison. The comparison is implemented by the conversion between the additive shares and the binary shares. Specially, $[z] = [x - y]$ is converted to the binary shares $\langle z \rangle$ through additive circuit with $\log_2 \ell$ round of communication. Subsequently, the binary shares of z 's sign bit can be determined by $\langle b \rangle = \langle z \rangle \gg (\ell - 1)^6$. Finally, the additive shares of $x < y$ can be derived by converting $\langle b \rangle$ to $[b]$ with one round of communication. Thus, the implementation of privacy-preserving compare algorithm cost $\log_2 \ell + 1$ round of communication and transmit 3456 bits.

Privacy-preserving maximum. The maximum of the n -element vector x is implemented by calling $\log_2 n$ privacy-preserving comparisons using the tree reduction algorithm [Knott et al., 2021].

Privacy-preserving exponential. Function exponential is complex and usually implemented using the repeated-squaring approximation method

$$e^x = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{2^n}\right)^{2^n}, \quad (8)$$

which converts exponential calculations into addition and square operations. By fault, iterations is set $n = 8$ in [Knott et al., 2021].

Privacy-preserving reciprocal. Function reciprocal $\frac{1}{x}$ is implemented by Newton-Raphson method, which converts reciprocal calculations into addition and multiplication operations. The iterative formula is

$$y_{n+1} = y_n(2 - xy_n). \quad (9)$$

The initial value of the iteration is

$$y_0 = 3e^{\frac{1}{2}-x} + 0.003. \quad (10)$$

The number of iterations is set to 10 in [Knott et al., 2021] by default.

Privacy-preserving square root. \sqrt{x} is approximated by Newton-Raphson method in MPC, which converts exponential calculations into addition and multiplication operations. The iterative formula is

$$y_{n+1} = \frac{1}{2}y_n(3 - xy_n^2). \quad (11)$$

⁶ $\gg \ell$ denote shift ℓ bit to the right.

The initial value of the iteration is

$$y_0 = e^{-2.2(\frac{\pi}{2}+0.2)} + 0.198046875. \quad (12)$$

The number of iterations is set to 3 in [Knott et al., 2021] by default.

5.5 Security Analysis

SECFWT adheres to a semi-honest (also known as honest-but-curious) assumption similar to the works of Li et al. [2023] and Dong et al. [2023], where honest participants constitute the majority. Under this assumption, the security of SECFWT can be formally proved against static semi-honest adversaries denoted as \mathcal{A} , which can potentially corrupt no more than one of the servers in the hybrid model.

SECFWT is constructed from the well-established sub-protocols of Knott et al. [2021], Zheng et al. [2023b], and we invoke these protocols in a black-box manner. Leveraging the concept of composable security established by Canetti [2001], it is easy to see that the security of SECFWT is guaranteed in the sub-protocols hybrid model.