

# Whisper: Efficiently Preserving Both Your Prompt and Model Privacy in LLM Fine-Tuning

Anonymous ACL submission

## Abstract

We propose Whisper, an efficient Privacy-Preserving Fine-Tuning (PPFT) framework that provides Dual-privacy for cloud-based APIs by leveraging homomorphic encryption to protect both user inputs and fine-tuned model parameters. Whisper redesigns the encrypted fine-tuning pipeline at the architectural level by introducing a backpropagation-free encrypted fine-tuning paradigm, enabling practical deployments on top of existing Privacy-Preserving Inference (PPI) schemes. To further improve efficiency, Whisper optimizes the underlying cryptographic protocols with a group-based packing strategy and comparison with bootstrapping, substantially increasing throughput and reducing execution overhead for fine-tuning workloads. Extensive experiments demonstrate strong privacy guarantees and competitive model performance, achieving a  $7.43\text{-}38.39\times$  efficiency improvement over state-of-the-art cryptographic approaches.

## 1 Introduction

Fine-tuning has become a crucial stage for achieving customization and task adaptation in Large Language Models (LLMs) (Houlsby et al., 2019; Hu et al., 2022), serving as a key driver for their deployment in various domain-specific applications. By retraining pre-trained models on downstream task datasets, fine-tuning can significantly enhance model specialization and practical utility. However, when clients send data to the cloud via fine-tuning APIs provided by model providers, the transmission and storage of such data in plaintext may lead to severe privacy leakage risks.

To mitigate these concerns, several studies have explored Privacy-Preserving Fine-Tuning (PPFT) schemes based on Differential Privacy (DP) (Yu et al., 2022; Du et al., 2023), which aim to protect sensitive information during the cloud-based fine-tuning process. While these approaches can

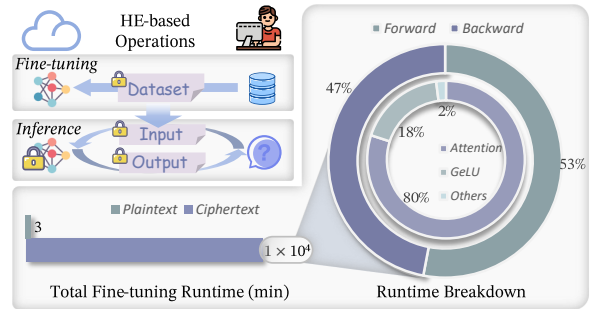


Figure 1: HE-Based Privacy-Preserving Fine-Tuning.

prevent attackers from extracting training data characteristics by launching inference attacks on the server side model (Kandpal et al., 2024), they remain ineffective against privacy threats arising from the plaintext storage and training of data on the cloud. According to a recent security report by Wiz Research<sup>1</sup>, sensitive information such as client prompts and log files could be directly retrieved from the ClickHouse cloud database used by DeepSeek. Such privacy leakage risks targeting client inputs are at least as harmful as model inference attacks and are often even more direct and damaging, potentially leading to the irreversible exposure of clients' private data.

From a cryptographic perspective, privacy-enhancing techniques, such as Homomorphic Encryption (HE) (Cheon et al., 2017) and Secure Multi-Party Computation (MPC) (Micali et al., 1987), can effectively protect the privacy of input data, thereby enabling Privacy-Preserving Machine Learning (PPML) (Cabrero-Holgueras and Pastrana, 2021). One direct application is that clients encrypt their local training data before uploading it, allowing the server to perform parameter updates entirely within the encrypted domain without ever accessing the raw inputs. However, executing forward and backward propagation entirely

<sup>1</sup><https://www.wiz.io/blog/wiz-research-uncovers-exposed-deepseek-database-leak>

069 over encrypted data incurs substantial runtime over- 116  
070 head. As shown in Fig. 1, the HE-based PPFT 117  
071 scheme (Rho et al., 2025) is up to  $10^4\times$  slower 118  
072 than plaintext execution, even when fine-tuning 119  
073 a 2-layer BERT<sub>base</sub> model. The optimized back- 120  
074 propagation still contributes nearly half of the cost, 121  
075 mainly from the Attention and GeLU. Moreover, 122  
076 since fine-tuned models remain on the cloud, all 123  
077 later inference must be performed under encryp- 124  
078 tion, leading to persistent computational overhead. 125

079 Consequently, existing PPFT architectures strug- 126  
080 gle to balance efficiency and privacy. To ad- 127  
081 dress this challenge, we propose **Whisper**—an 128  
082 encrypted PPFT framework built on a novel fine- 129  
083 tuning paradigm. Whisper achieves efficient Dual- 130  
084 privacy preservation, simultaneously protecting 131  
085 client inputs and the fine-tuned model while sub- 132  
086 stantially reducing computational overhead: 133

- 087 • **Novel Framework:** Whisper breaks the require- 134  
088 ment to perform backpropagation in the encrypted 135  
089 domain by introducing a PPFT framework that re- 136  
090 lies solely on forward computation. By generating 137  
091 fine-tuned dataset through Privacy-Preserving In- 138  
092 ference (PPI), Whisper fully localizes fine-tuning 139  
093 to the client. This ensures that client inputs remain 140  
094 encrypted and that the fine-tuned model parameters 141  
095 stay on the client side, enabling true Dual-privacy. 142
- 096 • **Efficient Protocol:** Whisper builds a non- 143  
097 interactive, high-throughput PPI protocol based on 144  
098 HE, and introduces targeted optimizations for en- 145  
099 crypted Attention and GeLU operations. These 146  
100 enhancements enable efficient processing of large- 147  
101 scale concurrent prompts directly in the encrypted 148  
102 domain, thereby allowing rapid generation of the 149  
103 fine-tuned dataset. 150
- 104 • **High Performance:** Whisper enables concurrent 151  
105 multi-prompt PPI for fine-tuned dataset generation, 152  
106 achieving a 7.43-38.39 $\times$  speedup over state-of-the- 153  
107 art (SOTA) PPI methods. By introducing no ran- 154  
108 dom noise, it improves model accuracy by 4–16% 155  
109 over DP-based approaches and inherently resists 156  
110 Data Reconstruction Attacks through ciphertext- 157  
111 only computation. 158

## 112 2 Preliminaries

### 113 2.1 Privacy-Preserving Inference

114 Privacy-Preserving Inference (PPI) is a central 165  
115 problem in PPML. Since inference involves only 166

forward computation, it is more amenable to en- 116  
cryptured execution, motivating extensive research 117  
on secure Transformer inference. 118

Existing PPI approaches can be broadly cate- 119  
gorized into MPC-based and HE-based methods. 120  
MPC protocols (Hao et al., 2022; Luo et al., 2024; 121  
Pang et al., 2024) rely on interactive execution and 122  
incur substantial communication overhead, leading 123  
to high latency in Wide-Area Network (WAN) set- 124  
tings. In contrast, HE-based schemes (Chen et al., 125  
2022; Zhang et al., 2025; Park et al., 2025) elim- 126  
inate interaction by operating directly on cipher- 127  
texts, but the computational cost of Transformers 128  
is significantly amplified under encryption, result- 129  
ing in prohibitive runtime overhead. As a result, 130  
both paradigms face scalability challenges under 131  
high-throughput and high-concurrency workloads. 132

Beyond cryptographic solutions, several works 133  
explore non-cryptographic defenses to improve ef- 134  
ficiency. CENTAUR (Luo et al., 2025a) applies 135  
vector permutation, while Morris et al. (2023) in- 136  
troduce Gaussian noise at the embedding layer. 137  
However, lacking formal security guarantees, these 138  
approaches remain vulnerable to Data Reconstruc- 139  
tion Attacks (DRA) (Pal et al., 2025), that recover 140  
sensitive inputs from intermediate results. 141

### 142 2.2 Privacy-Preserving Fine-Tuning

Privacy-Preserving Fine-Tuning (PPFT) is substan- 143  
tially more challenging than inference, as it re- 144  
quires parameter updates involving backpropaga- 145  
tion and gradient computation, which are particu- 146  
larly costly in the encrypted domain. Consequently, 147  
existing PPFT research remains limited. 148

P3EFT (Li et al., 2025) incorporates parameter- 149  
efficient techniques such as LoRA into HE-based 150  
frameworks to support encrypted fine-tuning on 151  
the cloud. However, since the fine-tuned model 152  
resides on the server, subsequent inference must 153  
still be executed over ciphertexts, incurring persis- 154  
tent computational overhead. Rho et al. (2025) re- 155  
duce training costs by replacing self-attention with 156  
Gaussian kernel attention, yet inference remains 157  
fundamentally encrypted. 158

SecP-Tuning (Luo et al., 2025b) proposes an 159  
MPC-based prompt tuning approach for efficient 160  
adaptation. However, it optimizes only client-side 161  
prompts, which are tightly coupled to a fixed back- 162  
bone model and become ineffective after model 163  
updates, thereby precluding the production of a 164  
reusable fine-tuned model for downstream use. 165

In parallel, DP has been integrated into fine- 166

tuning to mitigate inference attacks on trained models (Yu et al., 2022; Du et al., 2023). These methods primarily protect model parameters, while the privacy risks associated with the fine-tuning data itself are often insufficiently addressed.

## 2.3 Homomorphic Encryption

HE enables arbitrary polynomial computations to be performed directly over ciphertexts. In this work, we adopt the CKKS scheme (Cheon et al., 2017), which leverages SIMD encoding (Smart and Vercauteren, 2014) to encrypt a plaintext vector  $\mathbf{x} \in \mathbb{C}^{N/2}$  into a single ciphertext  $\tilde{\mathbf{x}} \in \mathcal{R}_q^2$ . The scheme supports a variety of algebraic operations in the ciphertext domain, including addition  $\widetilde{\mathbf{x} + \mathbf{y}} \leftarrow \tilde{\mathbf{x}} \boxplus \tilde{\mathbf{y}}$ , constant multiplication  $\widetilde{c\mathbf{x}} \leftarrow c \boxtimes \tilde{\mathbf{x}}$ , multiplication  $\widetilde{\mathbf{x}\mathbf{y}} \leftarrow \tilde{\mathbf{x}} \boxtimes \tilde{\mathbf{y}}$ , rotation  $\widetilde{\mathbf{x}^{(r)}} \leftarrow \text{Rot}(\tilde{\mathbf{x}}, r)$ , bootstrapping  $\tilde{\mathbf{x}}' \leftarrow \text{Bts}(\tilde{\mathbf{x}})$ , and expansion  $[\tilde{\mathbf{x}}_0, \tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{n-1}] \leftarrow \text{Expn}(\tilde{\mathbf{x}})$ . Formal definitions and implementation details of these operations are provided in the Appendix A.2.

## 2.4 Transformer

### 2.4.1 Attention

In Transformer models, Self-Attention (Vaswani et al., 2017) computes attention weights via  $\text{Softmax}(\mathbf{q}_n \mathbf{K}^\top)$  and applies them to the value matrix  $\mathbf{V}$  to obtain the output for each query vector:

$$\mathbf{o}_n = \lambda \mathbf{V}, \quad \lambda = \text{Softmax}(\mathbf{q}_n \mathbf{K}^\top / \sqrt{d}) \quad (1)$$

This computation requires  $\mathcal{O}(n)$  inner-product evaluations and an exponential nonlinearity. These operations are particularly costly under encrypted computation and further force HE-based PPML systems to process prompt requests strictly sequentially, thereby limiting overall system throughput.

To mitigate these limitations, recent work has proposed replacing Self-Attention with linear, recurrent-style formulations (Sun et al., 2023; Yang et al., 2024; Sun et al., 2025). These Linear-Attention methods maintain a cumulative state  $\mathbf{S}_n$  that is updated recursively and produce outputs:

$$\mathbf{o}_n = \mathbf{S}_n \mathbf{q}_n, \quad \mathbf{S}_n = D \mathbf{S}_{n-1} + \mathbf{v}_n \mathbf{k}_n^\top \quad (2)$$

Here,  $D$  denotes forgetting gate. This formulation removes both the exponential function and the  $\mathcal{O}(n)$  inner-product evaluations in Self-Attention, making the effective attention weights depend only on the per-token state  $\mathbf{S}_n$ . (see Appendix A.3 for details). We further find that this property can potentially enable HE-based SIMD techniques to pack

and process tokens from multiple prompts in parallel, thereby improving overall system throughput.

### 2.4.2 GeLU

In Transformer architectures, the GeLU activation introduces a smooth nonlinearity defined via the Gauss error function. Since non-linear functions cannot be directly evaluated in encrypted form, PPML systems typically approximate GeLU using piecewise and high-degree polynomial functions, as formalized in Appendix C. However, HE schemes do not support comparison (Cmp) over ciphertexts, forcing piecewise boundaries to be emulated by additional polynomial approximations. This substantially increases multiplicative depth and accelerates noise growth, often necessitating expensive bootstrapping (Bts) operations (Zhang et al., 2025; Park et al., 2025). As a result, evaluating GeLU becomes one of the dominant computational bottlenecks in HE-based PPML.

## 3 Whisper

### 3.1 Encryption-friendly Framework

*RQ 1: How can Whisper achieve Encryption-friendly and Dual-privacy PPFT?*

*Insight 1:* Directly implementing PPFT with cryptographic primitives often incurs significant runtime overhead. Section 1 (Fig. 1) shows that current approaches require expensive backpropagation for parameter updates, and the fine-tuned model remains on the server. Consequently, clients are still required to interact with the server in an encrypted manner during subsequent inference.

Inspired by the fine-tuning paradigm recently introduced by OpenAI<sup>2</sup>, we observe that a localized fine-tuning mechanism can be achieved through prompt-driven knowledge distillation (Chiang et al., 2023). As illustrated in Fig. 2 (a), the overall Whisper framework proceeds as follows: the client issues prompts  $\text{in}^{(i)}, i \in \{\text{prompts}\}$  to a cloud-based teacher model, receives the corresponding distributions  $\text{out}^{(i)}$ , and then constructs  $(\text{in}^{(i)}, \text{out}^{(i)})$  as fine-tuned dataset to locally train its student model.

The key advantage of this framework is that it decouples backpropagation from cloud-side computation, eliminating the high overhead of encrypted gradient updates and reducing PPFT to repeated instances of PPI, which is substantially more efficient under HE than encrypted backpropagation. As a result, Whisper enables direct parameter fine-tuning

<sup>2</sup><https://openai.com/index/api-model-distillation/>

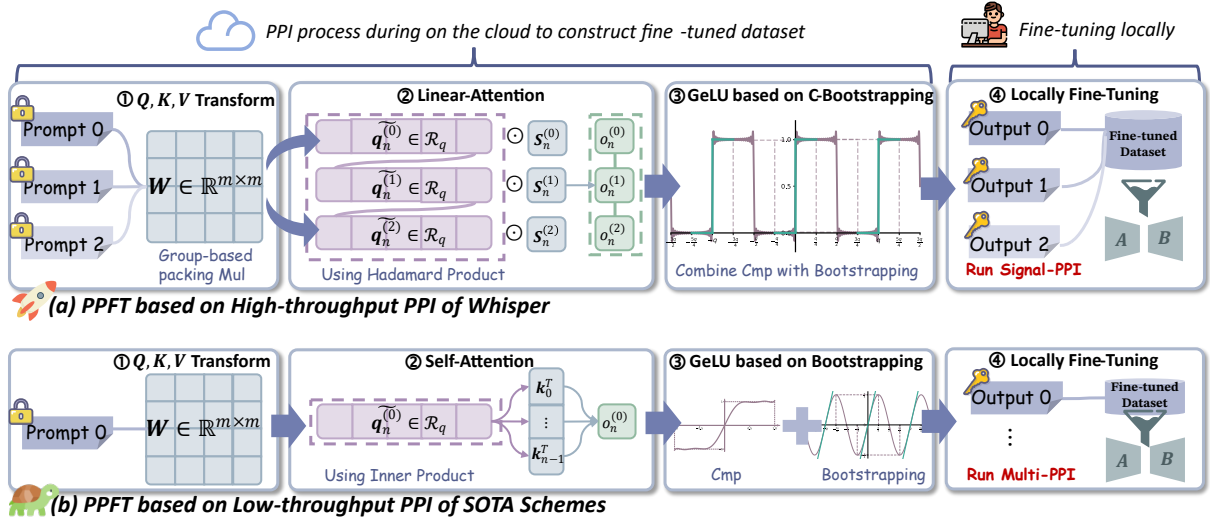


Figure 2: Whisper leverages cloud-side PPI to enable localized fine-tuning: (a) the high-throughput PPI design implemented in Whisper, and (b) the generalized PPI workflow used in current SOTA methods.

of a downstream model entirely on the client, rather than inference-time adaptation or prompt-based tuning. In contrast to vanilla knowledge distillation that relies on public or centrally curated data, Whisper derives task-specific supervision from encrypted prompts via PPI, without ever exposing plaintext data to the teacher.

**Threat Model.** Under a semi-honest security assumption, Whisper guarantees:

- i) *The server learns nothing about client prompts.*
- ii) *The client retains exclusive ownership of the fine-tuned model parameters.*

We exclude malicious adversaries and black-box model extraction (Liang et al., 2024), which are inherent to inference-as-a-service and orthogonal to our objectives, without enlarging the attack surface relative to plaintext inference and remaining compatible with standard defenses such as query limiting or model watermarking (Wang et al., 2024).

However, to construct the fine-tuned dataset, the client must issue a large number of prompts and obtain corresponding inference results from the cloud teacher model. This introduces a new technical challenge—how to efficiently support large-scale concurrent prompt inference under encryption.

### 3.2 Concurrent Processing

**RQ 2: How can Whisper efficiently support large-scale concurrent prompt inference under encryption?**

**Insight 2:** Existing PPI schemes based on MPC typically require multiple rounds of interactive communication between the client and the server. Under WAN conditions, such high communication rounds can lead to substantial latency overhead. To mitigate this issue, recent studies have shifted toward non-interactive approaches based on HE, enabling PPI without online interaction.

However, for LLMs, the core Self-Attention mechanism requires computing pairwise similarity scores between tokens within the same prompt. As shown in Step ② of Fig. 2 (b), conventional HE-based implementations pack all tokens of a prompt into a single SIMD ciphertext and rely on homomorphic rotations to compute the required inner products. The resulting attention weights are then obtained via the computationally expensive non-linear softmax operation,  $\lambda \leftarrow \text{Softmax}(q_n K^T)$ , and used to produce the attention output of the  $n$ -th token as  $o_n = \lambda V$ . This Self-Attention-oriented packing strategy restricts the system to processing only single-sample prompt (batch size = 1). As a result, the architecture is forced into a strictly sequential execution pattern, which makes it fundamentally incapable of scaling under high-concurrency inference workloads. Consequently, even SOTA HE-based PPI schemes become impractical when deployed within this PPFT framework.

To address this limitation, we draw inspiration from the Linear-Attention mechanism, which removes the dependency on pairwise token similarities in attention computation. Under this formulation, the attention weight for the  $n$ -th token in

prompt  $i$  depends only on its own accumulated state  $\mathcal{S}_n^{(i)}$ ,  $i \in \{\text{prompts}\}$  rather than on the similarities among the preceding  $n - 1$  tokens within the same prompt. As illustrated in Step ② of Fig. 2 (a), this property enables batching tokens from different prompts and packing them into a single SIMD ciphertext. As a result, attention computations can be executed in parallel across prompts, substantially improving inference throughput under high-concurrency workloads.

### 3.3 Protocol Optimization

**RQ 3: How can Whisper improve the overall efficiency of encrypted computation?**

**Insight 3:** Although HE eliminates interaction and reduces communication, ciphertext operations still incur substantial computational overhead—most notably from repeated rotation (key-switching) in linear layers and from the multiplicative depth consumption of polynomial-based non-linear approximations, which together lead to frequent bootstrapping and long execution time. To address these inefficiencies, as shown in Steps ① and ③ of Fig. 2(a), we introduce two targeted optimizations:

#### 3.3.1 Linear Protocol

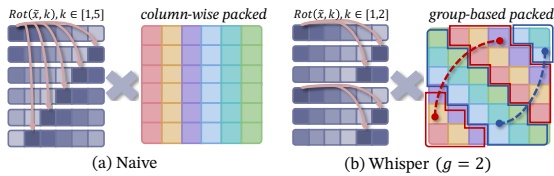


Figure 3: Column-wise packing (Naïve) vs. Group-based packing (Whisper).

In Step ① of our framework, the input ciphertext vector  $\tilde{\mathbf{x}} (\mathbf{x} \in \mathbb{C}^{1 \times m})$  is linearly mapped to  $(\tilde{\mathbf{g}}, \tilde{\mathbf{k}}, \tilde{\mathbf{v}})$  through homomorphic multiplications with plaintext weight matrices  $\mathbf{W}_{\{Q,K,V\}} \in \mathbb{C}^{m \times m}$ . This operation involves plaintext-ciphertext inner product (PCMul) computations, which require multiple homomorphic rotations—a major source of computational overhead. Although prior studies (Juvekar et al., 2018; Pang et al., 2024) have explored optimized packing strategies for the weight matrices to reduce the number of required rotations, the overall efficiency remains constrained by frequent and redundant rotation operations, which are inherently limited by the inefficiency of underlying key-switching procedures.

As illustrated in Fig. 3, we propose a group-based packing scheme built upon the diagonal

packing strategy (Halevi and Shoup, 2014). Unlike the naïve approach, which requires an independent rotation for each column vector in the matrix, our method performs diagonal-wise packing of the plaintext weight matrix and further partitions the packed diagonal vectors into  $g$  groups, each containing  $e = \frac{m}{g}$  vector elements, i.e.,  $\{\mathbf{w}_0^d, \dots, \mathbf{w}_{e-1}^d\}_0, \dots, \{\mathbf{w}_{e(g-1)}^d, \dots, \mathbf{w}_{m-1}^d\}_{g-1}$ . This grouping strategy allows different groups to reuse identical rotation operations on the ciphertext  $\tilde{\mathbf{x}}$ , thereby enabling more efficient computation of the inner product as follows:

$$\tilde{\mathbf{z}} = \sum_{i=0}^{g-1} \text{Rot} \left( \sum_{j=0}^{e-1} \text{Rot}(\tilde{\mathbf{x}}, j) \mathbf{w}_{ie+j}^d, ie \right) \quad (3)$$

This optimization reduces the rotation complexity from naïve  $\mathcal{O}(m^2)$  to  $\mathcal{O}(\frac{m}{g} + g - 2)$ , which is a convex function with respect to the group size  $g$ . Therefore, for a given matrix dimension  $m$ , there exists an optimal value of  $g$  that minimizes the overall number of rotations.

In practice, performing multiple rotations on the same ciphertext  $\tilde{\mathbf{x}}$  still incurs repeated key-switching operations, which remain expensive. To further reduce this cost, we adopt the hoisting technique (Halevi and Shoup, 2014), which reuses intermediate ciphertext decomposition results from key-switching to support batch rotation processing:

$$\text{Rot}(\tilde{\mathbf{x}}, 1), \dots, \text{Rot}(\tilde{\mathbf{x}}, m - 1) \leftarrow \text{Hoist}(\tilde{\mathbf{x}}) \quad (4)$$

This optimization allows roughly half of the rotation-related computations to be reused, and integrating it with the group-based packing strategy yields a further substantial performance gain. The detailed linear protocol and the corresponding convex optimization analysis are provided in the Appendix B.

#### 3.3.2 Non-linear Protocol

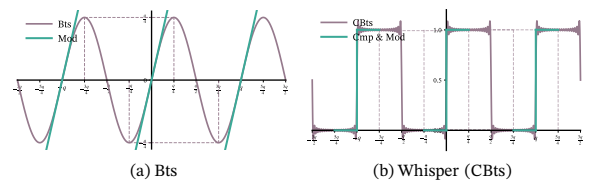


Figure 4: Bootstrapping (Naïve) vs. Comparison with Bootstrapping (Whisper).

Since HE schemes only support addition and multiplication over ciphertexts, non-linear functions  $f(x)$  are typically approximated by piecewise

† [UNK] represents out-of-vocabulary words

Ground Truth Prompt	The concept of a numbering machine is quite simple. It consists of various wheels inside and all wheels are stamped with numbers. But, what could be done with this?		
Reverse Example	The concept of a [UNK] machine is quite simple. It consists of various [UNK] inside and all [UNK] are [UNK] with numbers. But, what could be done with this?	Editor concept of a [UNK] machine is quite simple. It consists of various [UNK] inside and all [UNK] [unused146] [UNK] with numbers. But, what could be done ##n this?	[unused342] [unused231] ɾ n due full -  ʘ [unused619] [unused805] 5 ##head [unused908] bottom ɹ generation per ɐ australian pulling finger u [unused587] jane 1968 weren [unused578] dry e minute [unused704] parliament signed
ROUGLE-L	0.800	0.734	0.000
	(a) Shuffle-based Permutation	(b) Gaussian-noise Injection	(c) Whisper (based on cryptography)

Figure 5: Examples of how different PPI schemes defend against Hidden State Inversion attack.

polynomials, as described in Section 2. However, evaluating such functions often leads to rapid consumption of the modulus budget due to repeated ciphertext multiplications. When the modulus budget is exhausted, a bootstrapping (Bts) operation must be performed to refresh the ciphertext—raising its modulus and restoring its multiplicative depth so that subsequent encrypted computations remain correct. Specifically, given a ciphertext  $\tilde{x} \in \mathcal{R}_q^2$ , the process first applies ModRaise to lift its modulus to  $q' \gg q$ , and then performs modular reduction  $\tilde{x} + qI \pmod{q}$  within the interval  $[-\frac{q'}{2}, \frac{q'}{2}]$ , effectively raising the ciphertext’s modular and restoring its multiplicative depth.

However, modular reduction is inherently discontinuous, making it difficult to evaluate directly under HE. Prior work circumvents this by approximating the modulo function with smooth trigonometric forms (Cheon et al., 2018)—for example, replacing it with  $(\frac{q}{2\pi}) \sin((\frac{2\pi}{q}x))$ , as illustrated in Fig. 2 (a). In typical LLM inference pipelines, one Bts operation is already required after the Attention computation. Subsequently, evaluating the piecewise structure of the GeLU approximation introduces an additional non-linear comparison (Cmp) function. Since Cmp must also be approximated by a high-degree polynomial under HE, it consumes substantial multiplicative depth, thereby necessitating yet another Bts to refresh the modulus budget.

Inspired by functional Bts (Alexandru et al., 2025; Yang et al., 2025), we further observe that the approximation polynomial of Cmp can be reformulated with trigonometric components to achieve periodicity. This property enables the Cmp operation to be naturally integrated into the Bts process:

$$\text{Cmp}(x + qI \pmod{q}) = \frac{1 + \text{Sign}(\sin(\frac{2\pi}{q}x))}{2} \quad (5)$$

This extension enables both Cmp and Bts to be performed within a single homomorphic evaluation, effectively achieving functional fusion and significantly reducing the overall computational cost. We denote this integrated operation as Comparison with Bootstrapping (CBts). Furthermore, we can naturally replace the Sign function with its Fourier sine-series expansion, the expression can be equivalently written as:

$$\frac{1}{2} + \frac{2}{\pi} \sum_{k=1}^N \frac{\sin(\frac{2\pi(2k-1)}{q}x)}{2k-1} \quad (6)$$

As illustrated in Fig. 4, the approximation errors introduced by both the polynomial surrogate and the Cmp operation remain consistently small and well controlled, validating the effectiveness of our non-linear function construction. Detailed protocols and analyses for these non-linear components are provided in Appendix C.

## 4 Experiments

We evaluate our framework through a three-dimensional assessment that reflects the core properties required for PPFT: **Dimension I (Privacy):** Does Whisper prevent prompt leakage throughout encrypted pipeline? **Dimension II (Performance):** Does Whisper preserve task-level downstream performance? **Dimension III (Efficiency):** How much efficiency is gained through our HE optimizations under concurrent workloads?

**Baselines.** Since Whisper’s framework is compatible with arbitrary PPI schemes for constructing fine-tuned dataset, we include several SOTA cryptographic PPI methods (Pang et al., 2024; Zhang et al., 2025; Park et al., 2025) as baselines to comprehensively evaluate Whisper’s efficiency. In addition, we incorporate non-cryptographic

Table 1: Efficiency comparison of PCMul over varying matrix sizes (Runtime in s).

Framework	$\mathbb{C}^{128 \times 256} \times \mathbb{C}^{256 \times 256}$				$\mathbb{C}^{256 \times 256} \times \mathbb{C}^{256 \times 256}$				$\mathbb{C}^{128 \times 2048} \times \mathbb{C}^{2048 \times 2048}$			
	#KS	#Mul	Runtime	Factor	#KS	#Mul	Runtime	Factor	#KS	#Mul	Runtime	Factor
BOLT	136	128	0.19	1.00 ×	264	256	0.47	1.00 ×	3264	128	5.72	1.00 ×
NEXUS	514	32768	0.93	▼4.89 ×	1026	66536	2.16	▼4.60 ×	4112	262144	9.34	▼1.63 ×
Whisper	62	128	0.08	▲2.38 ×	124	256	0.23	▲2.04 ×	1552	128	2.23	▲2.57 ×

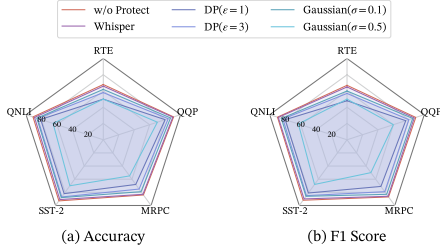


Figure 6: Accuracy and F1 performance on the GLUE.

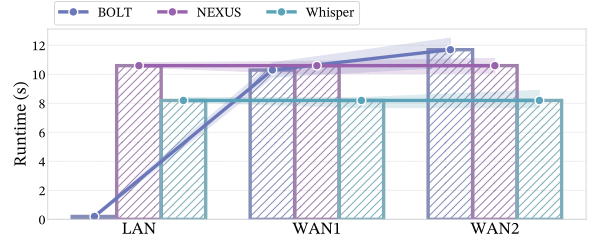


Figure 7: Efficiency Comparison of GeLU on  $\mathbb{C}^{2^{15}}$ .

PPI schemes, including Shuffle-based permutation method (Luo et al., 2025a) and Gaussian-noise injection at the embedding layer (Morris et al., 2023), as well as DP-based PPFT approaches (Du et al., 2023) that perturb model parameters directly. These baselines enable a thorough comparison of Whisper’s advantages in downstream performance and robustness.

**Models and Datasets.** To ensure comparability with prior work, we conduct experiments on both BERT and GPT-2, and evaluate performance across multiple standard datasets from the GLUE benchmark (Wang et al., 2018) (RTE, QQP, MRPC, SST-2, QNLI) and Fineweb-Edu dataset (Penedo et al., 2024). Additional experimental details are provided in the Appendix D.

#### 4.1 Privacy

We first note that the newly identified Hidden State Inversion (HSI) attack (Pal et al., 2025) provides significantly stronger reconstruction ability than prior DRA methods, enabling highly accurate recovery of token sequences. To evaluate Whisper’s prompt-privacy under this stronger threat model, we benchmark it against non-cryptographic defenses, including Gaussian-noise injection and Shuffle-based permutation.

**Attack Objective and Metric.** Given the intermediate hidden representation at layer  $l$ ,  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_N] \in \mathbb{C}^{N \times d}$ . The attacker aims to reconstruct the input sequence  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ . For each position  $i$  and candidate token from vo-

cabulary  $\mathbf{v} \in \mathcal{V}$ , a forward pass yields  $\hat{\mathbf{h}}_i(\mathbf{v}) = f_l(\mathbf{x}_{1:i-1}, \mathbf{v})$ , the attacker selects the sorted  $L_1$  distance as a fuzzy-matching criterion to iteratively recover the original sequence. Following (Pal et al., 2025), we adopt the ROUGE-L score to quantify the similarity between the reconstructed text and the ground-truth prompt. This metric captures both word-order and content-level overlap. Scores closer to 1 indicate more accurate reconstruction, and thus higher privacy leakage.

**Evaluation.** As shown in Fig. 5 on FineWeb-Edu dataset, Shuffle-based and Gaussian-noise defenses still yield high ROUGE-L scores (0.800 and 0.734), indicating that the reconstructed text remains closely aligned with the ground-truth prompt. This shows that plaintext-domain defenses based on noise injection or permutation fail to prevent reconstruction from intermediate representations. In contrast, Whisper achieves the lowest ROUGE-L score, producing semantically meaningless outputs. Without the client-side secret key, ciphertext-token reconstruction is computationally infeasible, highlighting the advantage of cryptographic protection for prompt privacy.

#### 4.2 Performance

We benchmark the downstream performance on GLUE after fine-tuning a BERT<sub>base</sub> model under four settings: (i) fine-tuning without any privacy protection, (ii) DP-based method, (iii) Gaussian-noise injection method, and (iv) our Whisper-based privacy-preserving approach.

As shown in Fig. 6, Whisper achieves down-

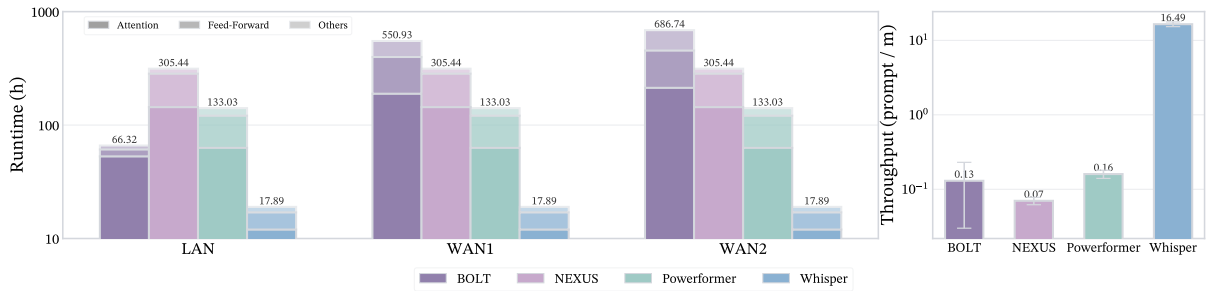


Figure 8: Runtime and throughput for generating a 128-sample fine-tuned dataset under different PPI schemes.

stream fine-tuning performance that is almost indistinguishable from the non-private baseline, with an average accuracy drop of less than 3% across all evaluated tasks. In contrast, training with DP leads to a noticeable degradation, reducing accuracy by 10-18% when privacy budget  $\epsilon$  is set between 1 and 3. Similarly, directly injecting Gaussian noise into the embedding layer causes consistent performance loss: when the noise standard deviation is  $\sigma = 0.1$ , accuracy drops by approximately 4-6%, and further degrades by more than 18% when  $\sigma$  increases to 0.5. Unlike noise-based methods that disrupt gradients or representations, Whisper avoids stochastic perturbations and preserves optimization stability through accurate non-linear approximations, enabling strong privacy protection without sacrificing model performance.

### 4.3 Efficiency

We conduct a comprehensive evaluation of Whisper’s efficiency improvements on both linear and non-linear layers, and further benchmark the processing time of several SOTA PPI schemes for constructing fine-tuned dataset under different network conditions (bandwidth, round-trip latency): LAN (3Gbps, 0.8ms), WAN1 (200Mbps, 40ms), and WAN2 (200Mbps, 80ms).

**Linear.** In the linear-layer evaluation (Table 1), Whisper substantially reduces the number of Rot operation, yielding runtime reductions of 2.04-2.57 $\times$  and 4.18-11.64 $\times$  compared to BOLT and NEXUS, respectively. This improvement primarily stems from our proposed group-optimization strategy, which effectively lowers the rotation complexity of inner-product computations. In addition, the incorporation of hoisting techniques enables reuse of key-switch primitives, further decreasing the number of rotations that must actually be executed.

**Non-linear.** In the non-linear-layer evaluation (Fig. 7), we focus on the efficiency of the GeLU operation. BOLT, which relies on MPC protocols, incurs multiple rounds of communication and thus suffers increasingly high latency under degraded network conditions. Although NEXUS avoids communication due to its fully homomorphic construction, the separation of Cmp and Bts evaluation causes its runtime to exceed that of BOLT in certain network settings. In contrast, Whisper integrates Cmp and Bts into a single operation via CBts, achieving significantly improved efficiency and reducing the runtime by 10.87-24.67% and 20.29% relative to BOLT and NEXUS, respectively.

**Dataset Generation.** Since fine-tuning is performed in plaintext, dataset generation via PPI becomes the dominant efficiency bottleneck in PPFT. Accordingly, in Fig. 8 we measure the time required by various SOTA PPI schemes to process 128 prompts on BERT<sub>base</sub> model and generate 128 outputs, each containing 10 tokens, to form a fine-tuned dataset. Whisper achieves substantial acceleration primarily due to its underlying Linear-Attention architecture: by decoupling the pairwise similarity computation among tokens from the accumulated state  $S$ , Whisper enables concurrent processing of multiple prompts, leading to a 103.06-235.57 $\times$  improvement in throughput and a 7.43-38.39 $\times$  reduction in runtime.

## 5 Conclusion

We present Whisper, an efficient HE-based PPFT framework that protects both user inputs and fine-tuned model privacy. To improve efficiency, Whisper relies on forward-only PPI, eliminating the need for costly encrypted backpropagation. We further design a suite of optimized protocols for both linear and nonlinear layers in language models. Extensive experiments demonstrate the effectiveness and practicality of Whisper.

## 614 Limitations

615 Following the standard security assumptions in HE-  
616 based PPML (Chen et al., 2022; Lee et al., 2023;  
617 Zhang et al., 2025; Park et al., 2025), our frame-  
618 work adopts a semi-honest security model, under  
619 which all parties are assumed to follow the pre-  
620 scribed protocol execution. While this assumption  
621 does not provide protection against fully malicious  
622 participants that arbitrarily deviate from the proto-  
623 col, it remains practical in real-world deployments.

624 **Minimal Server-Side Security Assumption.** In  
625 particular, although a malicious server could in  
626 principle deviate from the prescribed protocol,  
627 enforcing stronger guarantees against such be-  
628 havior would typically require interactive mech-  
629 anisms, such as MPC protocols augmented with  
630 message authentication codes (MACs) (Keller et al.,  
631 2016; Escudero et al., 2022) or zero-knowledge  
632 proofs (Atapoor et al., 2024; Zhou et al., 2025),  
633 to verify correct execution. In contrast, the threat  
634 model of Whisper imposes only a minimal secu-  
635 rity requirement, namely that the server must not  
636 be able to access any client plaintext information.  
637 Under our HE-based design, inference is entirely  
638 non-interactive and all computations are performed  
639 directly over encrypted data, which inherently pre-  
640 vents the server from observing or reconstructing  
641 any plaintext-sensitive information. As a result,  
642 this guarantee suffices for our threat model and  
643 remains consistent with the standard assumptions  
644 adopted in existing HE-based PPML systems.

645 **Client Incentive Compatibility.** Conversely, the  
646 client’s objective is to fine-tune a local model to  
647 improve performance on target downstream tasks.  
648 Any deviation from the prescribed protocol would  
649 be self-defeating, as it would prevent the client  
650 from acquiring useful knowledge from the teacher  
651 model and thus undermine this objective. Under  
652 this setting, Whisper does not aim to prevent black-  
653 box extraction of the teacher model via repeated  
654 queries (Tramèr et al., 2016; Liang et al., 2024),  
655 which is orthogonal to our goals and can be ad-  
656 dressed through standard complementary defenses  
657 such as query rate limiting or model watermarking  
658 (Jia et al., 2021; Wang et al., 2024).

659 Therefore, the adopted security model strikes a  
660 pragmatic balance between efficiency and protec-  
661 tion, providing sufficient guarantees for practical  
662 deployment, while not enlarging the attack surface  
663 compared to plaintext inference.

## References 664

- 665 Andreea Alexandru, Andrey Kim, and Yuriy Polyakov. 666  
2025. General functional bootstrapping using ckks. 667  
In *CRYPTO 2025*, pages 304–337. Springer.
- 668 Sebastian Angel, Hao Chen, Kim Laine, and Srinath 669  
Setty. 2018. Pir with compressed queries and amor- 670  
tized query processing. In *2018 IEEE symposium 671  
on security and privacy (IEEE S&P 2018)*, pages 672  
962–979. IEEE.
- 673 Shahla Atapoor, Karim Bagheri, Hilder VL Pereira, and 674  
Jannik Spiessens. 2024. Verifiable fhe via lattice- 675  
based snarks. *IACR Communications in Cryptology 676  
(CiC 2024)*, 1(1).
- 677 José Cabrero-Holgueras and Sergio Pastrana. 2021. Sok: 678  
Privacy-preserving computation techniques for deep 679  
learning. *Proceedings on Privacy Enhancing Tech- 680  
nologies (PETs 2021)*.
- 681 Tianyu Chen, Hangbo Bao, Shaohan Huang, Li Dong, 682  
Binxing Jiao, Daxin Jiang, Haoyi Zhou, Jianxin Li, 683  
and Furu Wei. 2022. The-x: Privacy-preserving trans- 684  
former inference with homomorphic encryption. In 685  
*Findings of the Association for Computational Lin- 686  
guistics (ACL 2022)*, pages 3510–3520.
- 687 Jung Hee Cheon, Kyoohyung Han, Andrey Kim, Miran 688  
Kim, and Yongsoo Song. 2018. Bootstrapping for 689  
approximate homomorphic encryption. In *EURO- 690  
CRYPT 2018*, pages 360–384. Springer.
- 691 Jung Hee Cheon, Andrey Kim, Miran Kim, and Yong- 692  
soo Song. 2017. Homomorphic encryption for arith- 693  
metic of approximate numbers. In *ASIACRYPT 2017*, 694  
pages 409–437. Springer.
- 695 Wei-Lin Chiang, Zhuohan Li, Ziqing Lin, Ying Sheng, 696  
Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan 697  
Zhuang, Yonghao Zhuang, Joseph E Gonzalez, and 698  
1 others. 2023. Vicuna: An open-source chatbot 699  
impressing gpt-4 with 90%\* chatgpt quality. See 700  
<https://vicuna.lmsys.org> (accessed 14 April 2023), 701  
2(3):6.
- 702 Minxin Du, Xiang Yue, Sherman SM Chow, Tianhao 703  
Wang, Chenyu Huang, and Huan Sun. 2023. Dp- 704  
forward: Fine-tuning and inference on language mod- 705  
els with differential privacy in forward pass. In *Pro- 706  
ceedings of the 2023 ACM SIGSAC Conference on 707  
Computer and Communications Security (CCS 2023)*, 708  
pages 2665–2679.
- 709 Daniel Escudero, Chaoping Xing, and Chen Yuan. 2022. 710  
More efficient dishonest majority secure computation 711  
over  $\mathbb{Z}_{2^k}$  via galois rings. In *CRYPTO 2022*, pages 712  
383–412. Springer.
- 713 Shai Halevi and Victor Shoup. 2014. Algorithms in 714  
helib. In *CRYPTO 2014*, pages 554–571. Springer.
- 715 Shai Halevi and Victor Shoup. 2021. Bootstrapping for 716  
helib. *Journal of Cryptology*, 34(1):7.

717	Meng Hao, Hongwei Li, Hanxiao Chen, Pengzhi Xing,	<i>Computer and Communications Security (AsiaCCS</i>	773
718	Guowen Xu, and Tianwei Zhang. 2022. Iron: Pri-	<i>2024)</i> , pages 1231–1245.	774
719	private inference on transformers. <i>Advances in Neu-</i>		
720	<i>ral Information Processing Systems (NeurIPS 2022)</i> ,	Jinglong Luo, Guanzhong Chen, Yehong Zhang, Shiyu	775
721	35:15718–15731.	Liu, Hui Wang, Yue Yu, Xun Zhou, Yuan Qi, and	776
		Zenglin Xu. 2025a. Centaur: Bridging the impos-	777
722	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	sible trinity of privacy, efficiency, and performance	778
723	Bruna Morrone, Quentin De Laroussilhe, Andrea	in privacy-preserving transformer inference. In <i>Pro-</i>	779
724	Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.	<i>ceedings of the 63rd Annual Meeting of the Asso-</i>	780
725	Parameter-efficient transfer learning for nlp. In <i>Inter-</i>	<i>ciation for Computational Linguistics (ACL 2025)</i> ,	781
726	<i>national Conference on Machine Learning (ICML</i>	<i>2019)</i> , pages 2790–2799. PMLR.	782
727			
728	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	Jinglong Luo, Yehong Zhang, Zhuo Zhang, Jiaqi Zhang,	783
729	Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,	Xin Mu, Hui Wang, Yue Yu, and Zenglin Xu. 2024.	784
730	Weizhu Chen, and 1 others. 2022. Lora: Low-	Secformer: fast and accurate privacy-preserving infer-	785
731	rank adaptation of large language models. <i>Inter-</i>	ence for transformer models via smpc. In <i>Findings of</i>	786
732	<i>national Conference on Learning Representations</i>	<i>the Association for Computational Linguistics (ACL</i>	787
733	<i>(ICLR 2022)</i> , 1(2):3.	<i>2024)</i> , pages 13333–13348.	788
734	Hengrui Jia, Christopher A Choquette-Choo, Varun	Jinglong Luo, Zhuo Zhang, Yehong Zhang, Shiyu	789
735	Chandrasekaran, and Nicolas Papernot. 2021. Entan-	Liu, Ye Dong, Hui Wang, Yue Yu, Xun Zhou, and	790
736	gled watermarks as a defense against model extrac-	Zenglin Xu. 2025b. Secp-tuning: Efficient privacy-	791
737	tion. In <i>30th USENIX Security Symposium (USENIX</i>	preserving prompt tuning for large language models	792
738	<i>Security 2021)</i> , pages 1937–1954.	via mpc. <i>arXiv preprint arXiv:2506.15307</i> .	793
739	Charanjit S Jutla and Nathan Manohar. 2022. Sine series	Vadim Lyubashevsky, Chris Peikert, and Oded Regev.	794
740	approximation of the mod function for bootstrapping	2010. On ideal lattices and learning with errors over	795
741	of approximate he. In <i>EUROCRYPT 2022</i> , pages	rings. In <i>EUROCRYPT 2010</i> , pages 1–23. Springer.	796
742	491–520. Springer.		
743	Chiraag Juvekar, Vinod Vaikuntanathan, and Anantha	Silvio Micali, Oded Goldreich, and Avi Wigderson.	797
744	Chandrakasan. 2018. {GAZELLE}: A low latency	1987. How to play any mental game. In <i>Proceedings</i>	798
745	framework for secure neural network inference. In	<i>of the Nineteenth ACM Symp. on Theory of Comput-</i>	799
746	<i>27th USENIX Security Symposium (USENIX Security</i>	<i>ing (STOC 1987)</i> , pages 218–229. ACM New York.	800
747	<i>2018)</i> , pages 1651–1669.		
748	Nikhil Kandpal, Krishna Pillutla, Alina Oprea, Peter	John Morris, Volodymyr Kuleshov, Vitaly Shmatikov,	801
749	Kairouz, Christopher A Choquette-Choo, and Zheng	and Alexander M Rush. 2023. Text embeddings	802
750	Xu. 2024. User inference attacks on large language	reveal (almost) as much as text. In <i>Proceedings of the</i>	803
751	models. In <i>Proceedings of the 2024 Conference on</i>	<i>2023 Conference on Empirical Methods in Natural</i>	804
752	<i>Empirical Methods in Natural Language Processing</i>	<i>Language Processing (EMNLP 2023)</i> , pages 12448–	805
753	<i>(EMNLP 2024)</i> , pages 18238–18265.	12460.	806
754	Marcel Keller, Emmanuela Orsini, and Peter Scholl.	Arka Pal, Rahul Krishna Thomas, Louai Zahran, Er-	807
755	2016. Mascot: faster malicious arithmetic secure	ica Choi, Akilesh Potti, and Micah Goldblum. 2025.	808
756	computation with oblivious transfer. In <i>Proceedings</i>	Hidden no more: Attacking and defending private	809
757	<i>of the 2016 ACM SIGSAC Conference on Computer</i>	third-party llm inference. In <i>ICLR 2025 Workshop on</i>	810
758	<i>and Communications Security (CCS 2016)</i> , pages	<i>Building Trust in Language Models and Applications</i> .	811
759	830–842.		
760	Seewoo Lee, Garam Lee, Jung Woo Kim, Junbum Shin,	Qi Pang, Jinhao Zhu, Helen Möllering, Wenting Zheng,	812
761	and Mun-Kyu Lee. 2023. Hetal: Efficient privacy-	and Thomas Schneider. 2024. Bolt: Privacy-	813
762	preserving transfer learning with homomorphic en-	preserving, accurate and efficient inference for trans-	814
763	ryption. In <i>International Conference on Machine</i>	formers. In <i>2024 IEEE Symposium on Security and</i>	815
764	<i>Learning (ICML 2023)</i> , pages 19010–19035. PMLR.	<i>Privacy (IEEE S&amp;P 2024)</i> , pages 4753–4771. IEEE.	816
765	Yang Li, Wenhan Yu, and Jun Zhao. 2025. Privtuner	Dongjin Park, Eunsang Lee, and Joon-Woo Lee. 2025.	817
766	with homomorphic encryption and lora: A p3eft	Powerformer: Efficient and high-accuracy privacy-	818
767	scheme for privacy-preserving parameter-efficient	preserving language model with homomorphic en-	819
768	fine-tuning of ai foundation models. <i>IEEE Trans-</i>	ryption. In <i>Proceedings of the 63rd Annual Meet-</i>	820
769	<i>actions on Wireless Communications</i> .	<i>ing of the Association for Computational Linguistics</i>	821
770	Jiacheng Liang, Ren Pang, Changjiang Li, and Ting	<i>(ACL 2025)</i> , pages 11090–11111.	822
771	Wang. 2024. Model extraction attacks revisited. In	Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov,	823
772	<i>Proceedings of the 19th ACM Asia Conference on</i>	Margaret Mitchell, Colin A Raffel, Leandro	824
		Von Werra, Thomas Wolf, and 1 others. 2024. The	825
		fineweb datasets: Decanting the web for the finest	826
		text data at scale. <i>Advances in Neural Informa-</i>	827
		<i>tion Processing Systems (NeurIPS 2024)</i> , 37:30811–	828
		30849.	829

830	Donghwan Rho, Taeseong Kim, Minje Park, Jung Woo Kim, Hyunsik Chae, Ernest K Ryu, and Jung Hee Cheon. 2025. Encryption-friendly llm architecture. In <i>International Conference on Learning Representations (ICLR 2025)</i> .	886
831		887
832		888
833		889
834		890
835	Nigel P Smart and Frederik Vercauteren. 2014. Fully homomorphic simd operations. <i>Designs, codes and cryptography</i> , 71(1):57–81.	891
836		892
837		893
838	Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram, Genghan Zhang, Yann Dubois, Xinlei Chen, Xiaolong Wang, Sanmi Koyejo, and 1 others. 2025. Learning to (learn at test time): Rnns with expressive hidden states. In <i>International Conference on Machine Learning (ICML 2025)</i> .	894
839		895
840		896
841		
842		
843		
844	Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. 2023. Retentive network: A successor to transformer for large language models. <i>arXiv preprint arXiv:2307.08621</i> .	
845		
846		
847		
848		
849	Florian Tramèr, Fan Zhang, Ari Juels, Michael K Reiter, and Thomas Ristenpart. 2016. Stealing machine learning models via prediction {APIs}. In <i>25th USENIX Security Symposium (USENIX Security 2016)</i> , pages 601–618.	
850		
851		
852		
853		
854	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in Neural Information Processing Systems (NeurIPS 2017)</i> , 30.	
855		
856		
857		
858		
859	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In <i>Proceedings of the 2018 EMNLP workshop BlackboxNLP: Analyzing and interpreting neural networks for NLP</i> , pages 353–355.	
860		
861		
862		
863		
864		
865		
866	Zhenyi Wang, Yihan Wu, and Heng Huang. 2024. Defense against model extraction attack by bayesian active watermarking. In <i>International Conference on Machine Learning (ICML 2024)</i> .	
867		
868		
869		
870	Songlin Yang, Bailin Wang, Yu Zhang, Yikang Shen, and Yoon Kim. 2024. Parallelizing linear transformers with the delta rule over sequence length. <i>Advances in Neural Information Processing Systems (NeurIPS 2024)</i> , 37:115491–115522.	
871		
872		
873		
874		
875	Zhaomin Yang, Chao Niu, Benqiang Wei, Zhicong Huang, Cheng Hong, and Tao Wei. 2025. Rboot: Accelerating homomorphic neural network inference by fusing relu within bootstrapping. <i>Cryptology ePrint Archive</i> .	
876		
877		
878		
879		
880	Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, and 1 others. 2022. Differentially private fine-tuning of language models. In <i>International Conference on Learning Representations (ICLR 2022)</i> , pages 1–19.	
881		
882		
883		
884		
885		
	Jiawen Zhang, Xinpeng Yang, Lipeng He, Kejia Chen, Wen-jie Lu, Yinghao Wang, Xiaoyang Hou, Jian Liu, Kui Ren, and Xiaohu Yang. 2025. Secure transformer inference made non-interactive. In <i>Network and Distributed System Security Symposium (NDSS 2025)</i> .	
	Zhelei Zhou, Yun Li, Yuchen Wang, Zhaomin Yang, Bingsheng Zhang, Cheng Hong, Tao Wei, and Wenguang Chen. 2025. Zhe: Efficient zero-knowledge proofs for he evaluations. In <i>2025 IEEE Symposium on Security and Privacy (IEEE S&amp;P 2025)</i> , pages 3328–3346. IEEE.	

## 897 A Details of Preliminaries

### 898 A.1 Notation

899 The notation used in this paper is summarized in Table 2.

Table 2: Summary of notations.

Notation	Description
$x, \mathbf{x}, X$	Scalar, vector, and matrix, respectively.
$\tilde{\mathbf{x}}$	HE ciphertext vector.
$\tilde{\mathbf{x}}$	Batched HE ciphertext vector.
$\chi$	Discrete gaussian distribution.
$\ell$	Multiplicative depth.
$N$	Polynomial degree in the HE scheme.

### 901 A.2 Homomorphic Encryption

902 Homomorphic Encryption (HE) supports arbitrary  
 903 polynomial operations over ciphertexts. In this  
 904 work, we adopt the CKKS scheme (Cheon et al.,  
 905 2017), which enables mapping plaintext complex  
 906 vectors in  $\mathbb{C}^{N/2}$  into ciphertext polynomials over  
 907 the ring  $\mathcal{R}_q \leftarrow \mathbb{Z}[X]/(X^N + 1)$  for approximate  
 908 homomorphic computation. The ciphertext mod-  
 909 ulus  $q$  is organized as a modulus chain  $q_\ell \leftarrow$   
 910  $q_0 \cdot \prod_{i=1}^{\ell} Q_i$ , allowing up to  $\ell$  multiplicative levels.  
 911 The homomorphic operations used in this paper are  
 912 summarized as follows.

- 913 • **KeyGen:** CKKS employs the RLWE assump-  
 914 tion (Lyubashevsky et al., 2010) to generate a key  
 915 pair (sk, pk). The secret key is sampled as  $s \leftarrow \chi$   
 916 in  $\mathcal{R}_q$ . The public key is defined as  $\text{pk} = (b, a) \in$   
 917  $\mathcal{R}_q^2$ , where  $b = -as + e, a \leftarrow \mathbb{Z}_q, e \leftarrow \chi$ .
- 918 • **SIMD Encoding:** A complex vector  $\mathbf{x} \in \mathbb{C}^{N/2}$  is  
 919 packed into the  $N/2$  SIMD slots of a polynomial  
 920 and scaled by a fixed-point factor  $\Delta$ , producing an  
 921 encoded plaintext  $\mathbf{m} \leftarrow \text{Encode}(\mathbf{x})$ , represented as  
 922 an integer polynomial  $\mathbf{m}_\Delta = \Delta \cdot \mathbf{m} \in \mathcal{R}$ .
- 923 • **Encryption:** Given a plaintext polynomial  $\mathbf{m}_\Delta$ , the  
 924 ciphertext is generated as  $\tilde{\mathbf{x}} = (c_0, c_1) = (bu + e_1 +$   
 925  $\mathbf{m}_\Delta, au + e_2)$ , where  $u, e_1, e_2 \leftarrow \chi$ .
- 926 • **Decryption:** Using the secret key, the plaintext is  
 927 recovered as  $\mathbf{m}_\Delta = c_0 + c_1s \pmod{q}$ , followed  
 928 by decoding and rescaling.
- 929 • **Addition:** For ciphertexts  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$ , homomorphic  
 930 addition is defined as  $\tilde{\mathbf{x}} + \tilde{\mathbf{y}} \leftarrow \tilde{\mathbf{x}} \boxplus \tilde{\mathbf{y}}$ .

- 931 • **Multiplication:** For ciphertexts  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$ , homomor-  
 932 phic multiplication is defined as  $\tilde{\mathbf{x}} \times \tilde{\mathbf{y}} \leftarrow \tilde{\mathbf{x}} \boxtimes \tilde{\mathbf{y}}$ , fol-  
 933 lowed by relinearization and rescaling as required  
 934 by CKKS.

- 935 • **Rotation:** A ciphertext  $\tilde{\mathbf{x}}$  can be rotated by  $r$  SIMD  
 936 slots using the rotation operation  $\tilde{\mathbf{x}}^{(r)} \leftarrow \text{Rot}(\tilde{\mathbf{x}}, r)$ .

- 937 • **Expansion:** Given a ciphertext  $\tilde{\mathbf{x}}$ , the expan-  
 938 sion operator produces a vector of ciphertexts  
 939  $[\tilde{\mathbf{x}}_0, \tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{m-1}] \leftarrow \text{Expn}(\tilde{\mathbf{x}})$ , with technical de-  
 940 tails provided in (Angel et al., 2018).

- 941 • **Bootstrapping:** Given a ciphertext  $\tilde{\mathbf{x}} \in \mathcal{R}_q^2$ , the  
 942 bootstrapping operation lifts it to a larger modu-  
 943 lus  $q' \gg q$ , thereby restoring noise budget and en-  
 944 abling additional multiplicative levels:  $\tilde{\mathbf{x}}' \in \mathcal{R}_{q'}^2 \leftarrow$   
 945  $\text{Bts}(\tilde{\mathbf{x}})$ . Further details are provided (Halevi and  
 946 Shoup, 2021).

### 947 A.3 Transformer

948 The Transformer, exemplified by models such as  
 949 BERT and GPT, is a deep neural architecture com-  
 950 posed of a stack of identical blocks, each con-  
 951 taining an Attention module followed by a Feed-  
 952 Forward Network (FFN). In this work, we focus  
 953 solely on the computations involved in the infer-  
 954 ence stage.

955 **Self-Attention.** Given an input matrix  $X$ , the  
 956 model first projects it into the corresponding query,  
 957 key, and value representations:

$$958 \mathbf{Q} = XW_Q, \mathbf{K} = XW_K, \mathbf{V} = XW_V \quad (7)$$

959 For query vector  $\mathbf{q}_n$  of token  $n$ , Self-Attention is  
 960 computed as:

$$961 \mathbf{o}_n = \lambda \mathbf{V}, \quad \lambda = \text{Softmax}(\mathbf{q}_n \mathbf{K}^\top / \sqrt{d_k}) \quad (8)$$

962 This mechanism captures contextual semantic rela-  
 963 tionships by measuring the similarity between  $\mathbf{q}_n$   
 964 and all key vectors through inner-product scores,  
 965 followed by Softmax normalization to obtain atten-  
 966 tion score. However, this process inevitably incurs  
 967  $\mathcal{O}(n)$  inner-product computations per token, along  
 968 with the additional cost of Softmax normalization.

969 **Linear-Attention.** To reduce the computational  
 970 complexity of Self-Attention, several works intro-  
 971 duce recurrent-style formulations into the attention  
 972 mechanism, leading to the class of Linear-Attention  
 973 methods:

$$974 \mathbf{o}_n = S_n \mathbf{q}_n, \quad S_n = DS_{n-1} + \mathbf{v}_n \mathbf{k}_n^\top \quad (9)$$

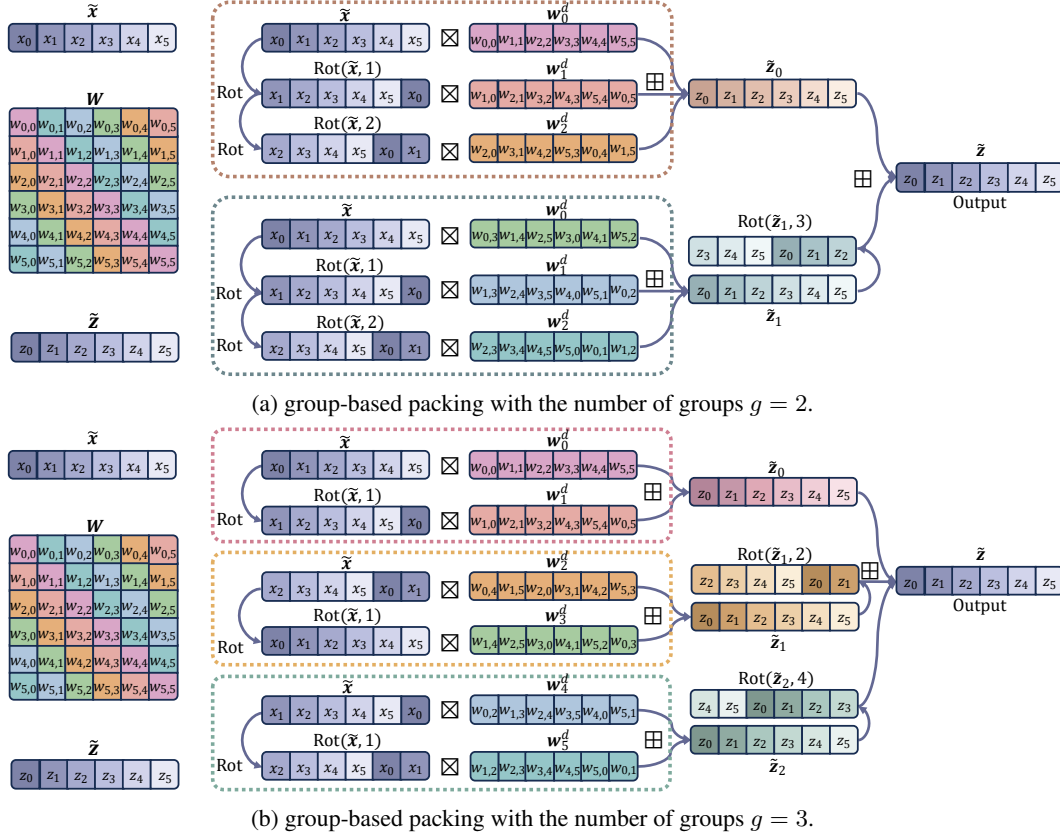


Figure 9: PCMul implemented with the group-based packing strategy, with 9a and 9b corresponding to different group sizes. The illustrated example uses  $z = x \times W$ , where  $x \in \mathbb{C}^{1 \times 6}$ ,  $W \in \mathbb{C}^{6 \times 6}$ , and  $z \in \mathbb{C}^{1 \times 6}$ .

By maintaining an accumulated state  $S$ , the model captures the semantic relationship between the current token  $n$  and all previous tokens, and uses this as the attention score. This formulation effectively removes the  $\mathcal{O}(n)$  inner-product operations and the Softmax normalization required in standard Self-Attention.

The state update depends on incorporating the current token information  $(k_n, v_n)$ , while forgetting outdated token contributions. Existing schemes differ in how they design the forgetting gate  $D$ . RetNet (Sun et al., 2023) applies a decay factor  $\gamma$  to gradually forget older token information. DeltaNet (Yang et al., 2024) updates the state by subtracting  $(S_{n-1}k_n)k_n^\top$ , thereby removing the component of  $S_{n-1}$  that is associated with  $k_n$ . TTT (Sun et al., 2025) instead constructs an optimizer over the current token information  $(v_n, k_n)$  to produce the updated state  $S_n$ .

**GeLU.** The Feed-Forward Network (FFN) introduces nonlinearity independently for each token:

$$\text{FFN}(x) = W_2 \text{GeLU}(W_1 x + b_1) + b_2 \quad (10)$$

Both BERT and GPT employ the Gaussian Error Linear Unit (GeLU) as the activation function, defined as:

$$\begin{aligned} \text{GeLU}(x) &= x \cdot \Phi(x), \\ &\approx 0.5x (1 + \tanh(\alpha(x + \beta x^3))) \quad (11) \end{aligned}$$

where  $\alpha = \sqrt{\frac{2}{\pi}}$ ,  $\beta = 0.044715$ .

where  $\Phi(\cdot)$  denotes the standard Gaussian cumulative distribution function. In practice, GeLU is often approximated using a smooth tanh-based formulation.

Since HE supports only polynomial operations, existing approaches (Zhang et al., 2025; Park et al., 2025) further approximate the GeLU function using piecewise polynomial representations.

## B Details of Linear Protocol

### B.1 Plaintext-Ciphertext Multiplication

**Group Optimization.** Inspired by diagonal packing (Halevi and Shoup, 2014) and the BSGS technique (Halevi and Shoup, 2021), we introduce a group-based packing strategy for plaintext-ciphertext multiplication (PCMul) that substan-

tially reduces the number of required homomorphic rotations.

As illustrated in Fig. 9, we begin by reorganizing the plaintext weight matrix  $\mathbf{W} \in \mathbb{C}^{m \times m}$  into its diagonal vectors and subsequently partitioning these diagonals into  $g$  groups, each containing  $e = \lceil \frac{m}{g} \rceil$  weight vectors:  $\{\mathbf{w}_0^d, \dots, \mathbf{w}_{e-1}^d\}_0, \dots, \{\mathbf{w}_{e(g-1)}^d, \dots, \mathbf{w}_{m-1}^d\}_{g-1}$ . For each group  $i \in [0, g)$ , the ciphertext vector  $\tilde{\mathbf{x}} (\mathbf{x} \in \mathbb{C}^{1 \times m})$  requires  $e - 1$  rotations, i.e.,  $\text{Rot}(\tilde{\mathbf{x}}, j)$  for  $j \in [0, e)$ . The corresponding plaintext weights  $\mathbf{w}_j^d$  can be shifted directly in the ciphertext and multiplied with the rotated  $\text{Rot}(\tilde{\mathbf{x}}, j)$  to obtain intermediate results  $\tilde{\mathbf{z}}_j$ . Since all groups reuse the same rotation outputs, the total number of rotations applied to  $\tilde{\mathbf{x}}$  remains exactly  $e - 1$ .

We then aggregate the intermediate results within each group:

$$\tilde{\mathbf{z}}_i = \sum_{j=1}^e \tilde{\mathbf{z}}_{ie+j} = \sum_{j=0}^{e-1} \text{Rot}(\tilde{\mathbf{x}}, j) \mathbf{w}_{ie+j} \quad (12)$$

and further rotate and sum the group-level results to produce the final output:

$$\tilde{\mathbf{z}} = \sum_{i=0}^{g-1} \text{Rot}(\tilde{\mathbf{z}}_i, ie) \quad (13)$$

In terms of complexity, the ciphertext vector  $\tilde{\mathbf{x}}$  incurs  $e - 1$  rotations, while the aggregation across all groups requires an additional  $g - 1$  rotations, yielding a total rotation cost of  $\mathcal{O}(e + g - 2)$  along with  $n$  plaintext-ciphertext multiplications. By further applying hoisting to reuse key-switch intermediates, the effective rotation cost can be reduced to  $\mathcal{O}((e + g - 2)/2)$ .

Crucially, this grouping mechanism turns the rotation complexity into a convex optimization problem: for any given  $m$ , there exists an optimal number of groups  $g$  that minimizes the overall cost. This property provides a principled basis for selecting parameters when designing efficient homomorphic linear transformations.

**Convexity of the Rotation Complexity.** The convex optimization analysis of the rotation complexity of PCMul, enabled by our group-based packing strategy, is presented as follows.

**Theorem 1 (Rotation Convexity)** *Let the rotation complexity:*

$$f(g) = \frac{m}{g} + g - 2 \quad (14)$$

where  $m > 0$  is a fixed constant and  $g > 0$ . Then  $f(g)$  is strictly convex over its domain  $(0, \infty)$ .

**Proof 1** *To establish convexity, we analyze the second-order derivative of  $f(g)$ .*

**Step 1: Compute first and second derivatives.** *The first derivative is*

$$f'(g) = 1 - \frac{m}{g^2} \quad (15)$$

*Differentiating again, the second derivative becomes*

$$f''(g) = \frac{2m}{g^3} \quad (16)$$

**Step 2: Verify positivity of the second derivative.** *Because the domain satisfies  $g > 0$ , we have  $g^3 > 0$ . Since the problem further assumes  $m > 0$ , it follows that*

$$f''(g) = \frac{2m}{g^3} > 0, \quad \text{for all } g > 0 \quad (17)$$

**Step 3: Conclude convexity.** *The strict positivity of the second derivative on the entire domain directly implies that  $f(g)$  is strictly convex on  $(0, \infty)$ .*

**Batched Group Optimization.** The example above illustrates the computation for a single prompt token. In general, the number of available SIMD slots  $n$  in the homomorphic polynomial ring is typically larger than the vector dimension  $m$ , naturally enabling ciphertext batching. Leveraging this property, our protocol—shown in Fig. 10—can pack token vectors from multiple prompts into a single ciphertext, thereby supporting cross-prompt batched linear transformations.

Formally, consider tokens from  $\tau = \lfloor n/m \rfloor$  prompts, denoted as  $[\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau-1)}] \in \mathbb{C}^{\tau \times m}$ . We first reorder the elements in column-major order and pack them into a single ciphertext vector:

$$\tilde{\mathbf{x}} \leftarrow \text{Enc}(x_0^{(0)}, \dots, x_0^{(\tau-1)} || \dots || x_{m-1}^{(0)}, \dots, x_{m-1}^{(\tau-1)}) \quad (18)$$

Following the same procedure as in the non-batched setting, the plaintext weight matrix  $\mathbf{W}$  is reorganized into a series of diagonal-wise vectors and partitioned into  $g$  groups. To match the packing

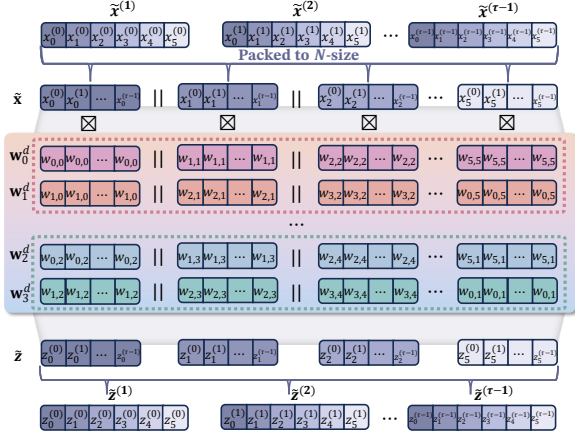


Figure 10: Batched PCMul protocol with group-based packing strategy.

layout of  $\tilde{\mathbf{x}}$ , each diagonal vector  $w_j^d$  in every group is expanded by repeating each element  $\tau$  times, ensuring that the resulting vector fully occupies the corresponding SIMD slots:

$$\mathbf{w}_j^d \leftarrow \underbrace{w_{j,1}^d, \dots, w_{j,1}^d}_{\tau} \parallel \dots \parallel \underbrace{w_{j,m}^d, \dots, w_{j,m}^d}_{\tau} \quad (19)$$

Within each group, we perform  $e - 1$  rotations on the packed vector  $\tilde{\mathbf{x}}$ , *i.e.*,  $\text{Rot}(\tilde{\mathbf{x}}, \tau j)$ ,  $j \in [0, e)$ , and multiply each rotated vector with its corresponding  $w_j^d$  to obtain intermediate results  $\tilde{\mathbf{z}}_j$ . Finally, the outputs of all groups are aggregated by applying rotations of stride  $\tau \cdot ie$  and summing them:

$$\tilde{\mathbf{z}} = \sum_{i=0}^{g-1} \text{Rot}(\tilde{\mathbf{z}}_i, \tau \cdot ie) \quad (20)$$

yielding the final batched linear-transformation result  $\tilde{\mathbf{z}}$ .

## B.2 Ciphertext-Ciphertext Multiplication

After obtaining the Attention inputs  $\mathbf{q}, \mathbf{k}, \mathbf{v} \leftarrow \mathbf{x} \times \mathbf{W}_{\{Q,K,V\}}$  via plaintext-ciphertext multiplication, Whisper replaces the standard Self-Attention with Linear-Attention to improve parallel processing capability. Unlike Self-Attention—which computes pairwise token interactions and normalizes them via Softmax—Linear-Attention maintains a continuously updated state vector  $\mathbf{S}$  to perform weighted aggregation. A key property is that each state update depends only on the current token  $n$ , eliminating the need for global normalization over all previous tokens:

$$\mathbf{o}_n = \mathbf{S}_n \mathbf{q}_n, \mathbf{S}_n = D \mathbf{S}_{n-1} + \mathbf{v}_n \mathbf{k}_n^\top \quad (21)$$

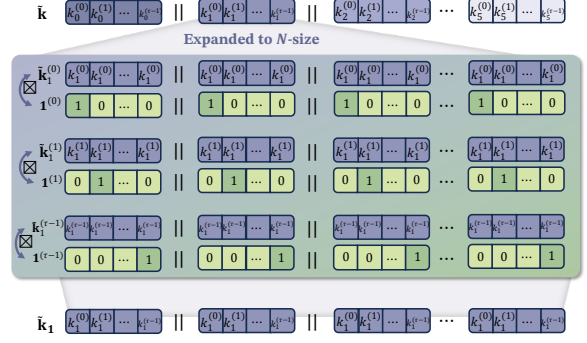


Figure 11: Batched CCMul protocol.

In the encrypted setting, implementing Linear-Attention requires computing the outer product between  $\tilde{\mathbf{v}}$  and  $\tilde{\mathbf{k}}$  under encryption (ciphertext-ciphertext multiplication, CCMul). This operation is efficiently supported by the homomorphic primitive `Expn`, which expands every element of  $\tilde{\mathbf{k}}$  into an equal-length vector, enabling direct execution of the required outer-product computation with  $\tilde{\mathbf{v}}$ .

Moreover, a key advantage of this computation pattern is its ability to update the Attention states  $\mathbf{S}$  for multiple prompts in parallel. When the packed vector  $\tilde{\mathbf{k}}$  contains tokens originating from different prompts:

$$\tilde{\mathbf{k}} \leftarrow \text{Enc}(k_0^{(0)}, \dots, k_0^{(\tau-1)} \parallel \dots \parallel k_{m-1}^{(0)}, \dots, k_{m-1}^{(\tau-1)}) \quad (22)$$

we can reorganize the expanded vectors produced by `Expn` using block-wise one-hot masking, as illustrated in Fig. 11. This reorganization aligns the expanded structure with the multi-prompt packing layout.

Multiplying the reconstructed vectors with the correspondingly packed  $\tilde{\mathbf{v}}$  then yields all prompt-specific state updates  $\mathbf{S}$  in a single operation. This mechanism enables the system to perform Attention updates for multiple prompts simultaneously, substantially improving throughput in highly concurrent settings.

## B.3 Secure Attention

Building on the implementations of the linear protocols PCMul and CCMul, along with their corresponding batched variants, we construct a Secure Attention protocol that supports concurrent processing of multiple prompts. To further improve efficiency, we incorporate an offline–online execution paradigm into the protocol design. The complete protocol is presented in Protocol 1.

---

### Protocol 1: Secure Attention

---

**Input:**  $\mathcal{C}$  holds input  $\mathbf{x} \in \mathbb{C}^{1 \times \tau m}$  batched  $\tau$  prompts;  
 $\mathcal{S}$  holds weight  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{C}^{\tau m \times \tau m}$ .  
**Output:**  $\mathcal{S}$  holds SIMD packed ciphertext vector  $\tilde{\mathbf{z}}$ ,  
where  $\mathbf{z} \leftarrow \text{Attention}(\mathbf{x})$ .

§ **Offline**

1 **for each matrix**  $\mathbf{W} \in \{\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V\}$  **do**  
2      $\mathcal{C}$  samples random vector  $\mathbf{a}$  for  $\mathbf{x}$  and different  
   weight matrix.  
3      $\mathcal{C}$  encrypts each random vector after SIMD  
   encoding to obtain  $\tilde{\mathbf{a}} \leftarrow \text{Enc}(\mathbf{a})$ , and then  
   transmits them to  $\mathcal{S}$ .  
4      $\mathcal{S}$  computes  $\tilde{\mathbf{b}} \leftarrow \text{PCMul}(\tilde{\mathbf{a}}, \mathbf{W})$  to obtain the  
   ciphertext mask  $\tilde{\mathbf{b}}$ .

5 **end**

§ **Online**

6 **for each matrix**  $\mathbf{W} \in \{\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V\}$  **do**  
7      $\mathcal{C}$  computes  $\mathbf{c} \leftarrow \mathbf{x} - \mathbf{a}$  and sends to  $\mathcal{S}$ .  
8 **end**  
9  $\mathcal{S}$  computes linear transformation:  
    $\tilde{\mathbf{q}} \leftarrow (\mathbf{c} \times \mathbf{W}_Q) \boxplus \tilde{\mathbf{b}}$ ;  
    $\tilde{\mathbf{k}} \leftarrow (\mathbf{c} \times \mathbf{W}_K) \boxplus \tilde{\mathbf{b}}$ ;  
    $\tilde{\mathbf{v}} \leftarrow (\mathbf{c} \times \mathbf{W}_V) \boxplus \tilde{\mathbf{b}}$ ;  
10  $\mathcal{S}$  updates  $\tilde{\mathbf{S}} := D \boxtimes \tilde{\mathbf{S}} \boxplus \text{CCMul}(\tilde{\mathbf{k}}^\top, \tilde{\mathbf{v}})$ .  
11  $\mathcal{S}$  executes  $\{\tilde{\mathbf{q}}_0, \tilde{\mathbf{q}}_1, \dots, \tilde{\mathbf{q}}_{k-1}\} \leftarrow \text{Expn}(\tilde{\mathbf{q}})$ .  
12 **for each row-wise vector**  $\tilde{\mathbf{s}}_j, j \in \{0, k-1\}$  **of**  $\tilde{\mathbf{S}}$  **do**  
13      $\mathcal{S}$  executes  $\tilde{\mathbf{z}} \leftarrow \boxplus_{j=0}^{k-1} (\tilde{\mathbf{q}}_j \boxtimes \tilde{\mathbf{s}}_j)$ .  
14 **end**  
15 **return**  $\tilde{\mathbf{z}}$ .

---

## C Details of Non-linear Protocol

### C.1 Comparison

In fixed-point HE schemes, each ciphertext multiplication increases the underlying plaintext scale by a factor of  $\Delta$ . To compensate for this scale growth, HE systems employ a modulus chain:

$$q_\ell \leftarrow q_0 \cdot (Q_1 \cdot Q_2 \cdots Q_{\ell-1}) \quad (23)$$

where each rescale step (after a multiplication) removes the highest modulus level, effectively reducing the ciphertext modulus and restoring the scale. Once the chain is fully consumed and only the base modulus  $q_0$  remains, no further ciphertext multiplications are possible. At this point, bootstrapping (Bts) must be invoked to replenish the modulus budget and restore the available multiplicative depth.

In standard Bts of CKKS, the system first performs ModRaise, lifting the ciphertext  $\tilde{x} \in \mathcal{R}_q^2$  into a larger modulus  $q' \gg q$ . This operation embeds  $\tilde{x}$  into the extended modulus domain—often expressed as a centered lifting such as  $\tilde{x} + qI$ , thereby expanding the representable interval from  $[-\frac{1}{2}q, \frac{1}{2}q]$  to  $[-\frac{1}{2}q', \frac{1}{2}q']$ . The Bts procedure then applies a modulo reduction (mod  $q$ ) to ensure the

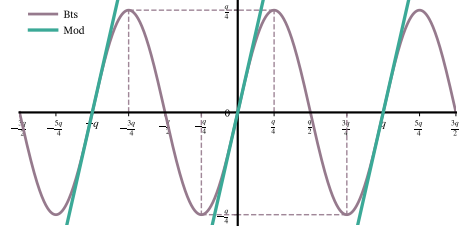


Figure 12: Implementing the Mod function via the periodicity of the Sin function.

decrypted value falls back into the correct range. However, modular reduction is inherently periodic and discontinuous, whereas HE natively supports only addition and multiplication—*i.e.*, continuous polynomial operations. Consequently, existing schemes (Cheon et al., 2018; Jutla and Manohar, 2022) exploit the periodic structure of trigonometric functions to rewrite the modulo operation as a composition of smooth periodic functions, which can then be approximated via Taylor series expansions and evaluated homomorphically (see Fig. 12).

In LLMs inference under HE, Bts commonly appears adjacent to non-linear operations such as Cmp for GeLU, which themselves are implemented through polynomial approximations.

$$\text{GeLU}(x) = \begin{cases} 0 & \text{if } x < -\phi_0 \\ P_0(x) & \text{if } -\phi_0 < x \leq -\phi_1 \\ P_1(x) & \text{if } -\phi_1 < x \leq \phi_2 \\ x & \text{if } x > \phi_2 \end{cases} \quad (24)$$

Therefore, if both Bts and these non-linear functions can be expressed within the same trigonometric-based approximation framework, they can be fused into a single homomorphic evaluation. This functional integration reduces computational overhead and significantly improves overall efficiency.

To enable such fusion, we first observe that the decision boundary of Cmp depends solely on the sign of its input. Meanwhile, HE bootstrapping naturally involves evaluating functions in the modular domain. Leveraging this observation, we exploit the periodicity of the sine function to express the sign of  $x$  modulo  $q$ , thereby embedding Cmp into the same functional space used for Bts, as CBts:

$$\text{Cmp}(x + qI \pmod{q}) = \frac{1 + \text{Sign}(\sin(\frac{2\pi}{q}x))}{2} \quad (25)$$

1215 This formulation leverages the fact that:

- 1216 •  $\sin\left(\frac{2\pi}{q}x\right)$  has the same sign as  $x$  within any
- 1217 principal interval of length  $q$ .
  
- 1218 • enabling Cmp to be computed directly on
- 1219 wrapped ciphertext values during Bts.

1220 Thus, both operations can be executed in a single  
 1221 homomorphic evaluation pipeline, reducing mul-  
 1222 tiplicative depth and eliminating a full additional  
 1223 Bts cycle.

1224 To homomorphically evaluate Eq. 25, we ap-  
 1225 proximate the Sign function using its well-known  
 1226 Fourier sine-series. For any odd, periodic function  
 1227 with period  $2\pi$ , the Fourier series contains only odd  
 1228 sine terms:

$$1229 \text{Sign}(\sin \theta) = \frac{4}{\pi} \sum_{k=1}^{\infty} \frac{\sin((2k-1)\theta)}{2k-1}, \quad \theta = \frac{2\pi}{q}x \quad (26)$$

1230 Substituting Eq. 26 into Eq. 25 yields:

$$1231 \text{Cmp}(x+qI \pmod{q}) = \frac{1}{2} + \frac{2}{\pi} \sum_{k=1}^{\infty} \frac{\sin\left(\frac{2\pi(2k-1)x}{q}\right)}{2k-1} \quad (27)$$

1232 To make the computation HE-compatible, the in-  
 1233 finite series must be truncated. The degree- $N$  ap-  
 1234 proximation is:

$$1235 \text{Cmp}_N(x+qI \pmod{q}) = \frac{1}{2} + \frac{2}{\pi} \sum_{k=1}^N \frac{\sin\left(\frac{2\pi(2k-1)x}{q}\right)}{2k-1} \quad (28)$$

1236 where each sine term can be further approximated  
 1237 using its Taylor expansion.

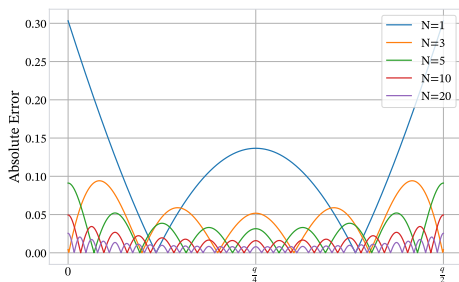


Figure 13: Absolute error of the Cmp function under different approximation degrees.

## 1238 C.2 Truncation Error Analysis

### Theorem 2 (Truncated Cmp Error Bound)

Let

$$E_N(x) = \text{Cmp}(x) - \text{Cmp}_N(x) \quad (29)$$

denote the truncation error of the Cmp function approximated by the first  $N$  terms of its Fourier sine expansion. Then the absolute error satisfies:

$$|E_N(x)| \leq \frac{4}{\pi} \sum_{k=N+1}^{\infty} \frac{1}{2k-1} \quad (30)$$

$$\approx \frac{4}{\pi} (\log N + \gamma + O(1/N))$$

where  $\gamma$  is the Euler–Mascheroni constant. Moreover, over any compact interval  $[a, b]$  that excludes the discontinuity points, the convergence is uniform and improves to:

$$|E_N(x)| = O\left(\frac{1}{N}\right), \quad x \in [a, b], b < a + q. \quad (31)$$

**Proof 2** The Cmp function is derived from a square-wave-like Sign function. Its truncated Fourier sine expansion takes the form

$$1240 \text{Cmp}_N(x) = \frac{1}{2} + \frac{2}{\pi} \sum_{k=1}^N \frac{\sin\left(\frac{2\pi(2k-1)x}{q}\right)}{2k-1} \quad (32)$$

By standard Fourier analysis of odd periodic square-wave functions, the truncation error of the series satisfies

$$1244 |E_N(x)| = \left| \frac{2}{\pi} \sum_{k=N+1}^{\infty} \frac{\sin((2k-1)x)}{2k-1} \right| \quad (33)$$

$$1245 \leq \frac{4}{\pi} \sum_{k=N+1}^{\infty} \frac{1}{2k-1}. \quad 1246$$

Using the asymptotic expansion of the partial sum of the harmonic series over odd indices yields

$$1248 \sum_{k=N+1}^{\infty} \frac{1}{2k-1} = \log N + \gamma + O(1/N) \quad (34) \quad 1249$$

which establishes the global error bound. On any compact interval excluding discontinuity points, classical Fourier convergence theory ensures uniform convergence of the sine series with rate  $O(1/N)$ , proving the local bound.

As illustrated in Fig. 13, within the periodic interval  $[0, \frac{1}{2}q]$ , the minimum approximation error between our CBts operation and the plaintext Cmp operation remains bounded within  $[1.4 \times 10^{-5}, 1.8 \times 10^{-4}]$  across different expansion orders  $N$ , and this approximation error does not accumulate across layers. Near the discontinuity ( $x = 0 \pmod{q}$ ), the series exhibits the classical Gibbs phenomenon; however, these points correspond to a measure-zero region and do not affect correctness in HE inference, where inputs are approximate floating-point ciphertext encodings.

---

### Protocol 2: Secure GeLU

---

**Input:**  $\mathcal{S}$  holds SIMD packed ciphertext vector input  $\tilde{x}$ .

**Output:**  $\mathcal{S}$  holds SIMD packed ciphertext vector  $\tilde{z}$ , where  $z \leftarrow \text{GeLU}(x)$ .

- 1 Compare  $x$  with the breakpoints:
    - $\tilde{b}_0 \leftarrow \text{CBts}(\tilde{x} \boxplus \phi_0) // b_0 = 1\{x > -\phi_0\};$
    - $\tilde{b}_1 \leftarrow \text{CBts}(\tilde{x} \boxplus \phi_1) // b_1 = 1\{x > -\phi_1\};$
    - $\tilde{b}_2 \leftarrow \text{CBts}(\tilde{x} \boxplus \phi_2) // b_2 = 1\{x > \phi_2\};$
  - 2 Compute segment selection:
    - $\tilde{s}_0 \leftarrow 1 \boxminus \tilde{b}_0 // s_0 = 1\{x < -\phi_0\};$
    - $\tilde{s}_1 \leftarrow \tilde{b}_0 \boxminus \tilde{b}_1 // s_1 = 1\{-\phi_0 < x < -\phi_1\};$
    - $\tilde{s}_2 \leftarrow \tilde{b}_1 \boxminus \tilde{b}_2 // s_2 = 1\{-\phi_1 < x < \phi_2\};$
    - $\tilde{s}_3 \leftarrow \tilde{b}_2 // s_3 = 1\{x > \phi_2\};$
  - 3 Compute GeLU:
    - $\tilde{z} \leftarrow (\tilde{s}_0 \boxtimes 0) \boxplus (\tilde{s}_1 \boxtimes P_0(\tilde{x})) \boxplus (\tilde{s}_2 \boxtimes P_1(\tilde{x})) \boxplus (\tilde{s}_3 \boxtimes \tilde{x});$
  - 4 return  $\tilde{z}$ .
- 

### C.3 Secure GeLU

Following most existing works that implement GeLU via the piecewise polynomial approximation shown in Eq. 24, we design a Secure GeLU protocol using CBts, which integrates both the Bts and Cmp operations. The full protocol is presented in Protocol 2.

## D Details of Experiments

### D.1 Evaluation Setting

All experiments were conducted on a machine equipped with a 64-core CPU and 128GB RAM. Each experiment was repeated 10 times, and the reported results are averages. The detailed experimental settings are described below.

- For the HE hyperparameters, we set the slot size to  $2^{15}$ , enabling 21 levels of multiplicative depth and providing 128-bit security.
- For performance evaluation, we fine-tune a pre-trained BERT<sub>base</sub> model and assess its downstream

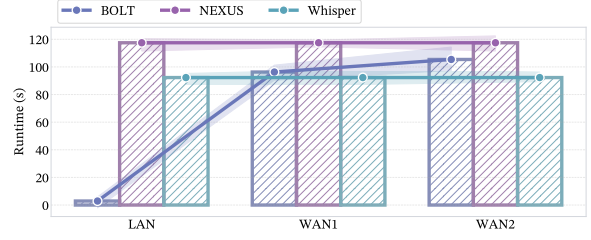


Figure 14: Efficiency Comparison of GeLU on  $\mathbb{C}^{2^{20}}$ .

performance on the GLUE benchmark. In the Whisper setting, to simulate a cloud-based teacher model, we assume that the teacher has already learned the target dataset and generates a fine-tuned dataset via PPI to enable PPFT. In contrast, noise-injection-based baselines assume direct access to the fine-tuning samples and perform model fine-tuning on the cloud.

- In the efficiency evaluation, we conduct fine-tuned dataset generation experiments using BERT<sub>base</sub> and GPT-2<sub>medium</sub>. The inputs consist of 128 prompts, each containing 128 tokens; correspondingly, the outputs are the inference results for 128 prompts, with each output comprising 10 (BERT<sub>base</sub>) or 5 (GPT-2<sub>medium</sub>) tokens.

### D.2 Additional Efficiency Evaluation

Table 3: Efficiency comparison of PCMul on  $\mathbb{C}^{k \times m} \times \mathbb{C}^{m \times m}$  ( $k = 256, m = 2048$ , and Runtime in s).

Framework	#KS	#Mul	Runtime	Factor
BOLT	6656	256	11.14	1.00 ×
NEXUS	8208	524288	20.69	▼1.85 ×
Whisper	3040	256	4.56	▲2.44 ×

**Linear.** We further evaluate the efficiency of the linear protocol PCMul under large-scale matrix computation scenarios in Table 3. For the NEXUS scheme, additional  $\mathcal{O}(\frac{km}{N})$  Expand operations with complexity  $\mathcal{O}(\log_2 N)$  are required to unpack the packed ciphertexts, which further introduces a multiplication overhead of  $\mathcal{O}(km)$ . Although BOLT reduces the theoretical rotation complexity by adopting the BSGS technique, its overall computational cost remains high in large-scale matrix operations, leading to a substantial increase in the number of key-switching and multiplication calls during execution. In contrast, our proposed group-based packing strategy reduces the rotation complexity to  $\mathcal{O}(k(\frac{m}{g} + g - 2)/2)$  and, when combined with

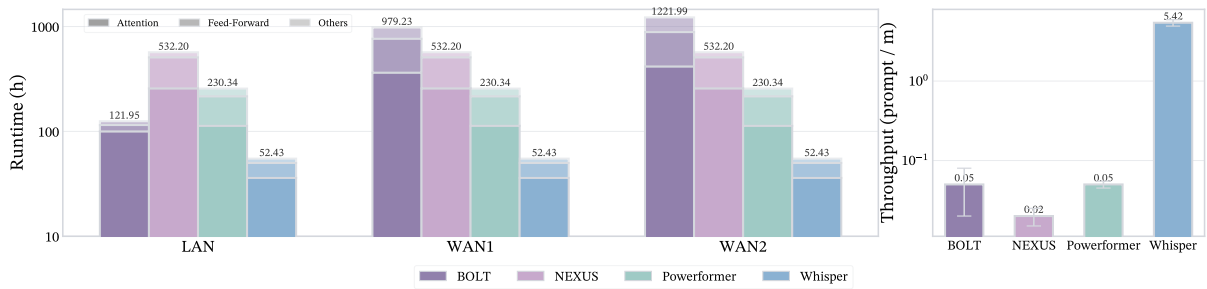


Figure 15: Runtime and throughput for generating a 128-sample fine-tuned dataset under different PPI schemes on GPT-2 model.

1318 hoisting, effectively decreases the number of key-  
 1319 switch operations in practice. As a result, our ap-  
 1320 proach maintains efficiency gains comparable to  
 1321 those observed in small-scale settings even under  
 1322 large-scale matrix computation workloads.

1323 **Non-linear.** We further evaluate the computa-  
 1324 tional efficiency of the non-linear GeLU protocol  
 1325 under large input lengths in Fig. 14. The BOLT  
 1326 scheme, which implements non-linear operations  
 1327 based on MPC, involves multiple rounds of in-  
 1328 teraction and is therefore highly sensitive to net-  
 1329 work communication conditions. In contrast, in  
 1330 the HE-based NEXUS scheme, the Bts and Cmp  
 1331 operations are executed independently, introducing  
 1332 additional computational overhead during polyno-  
 1333 mial evaluation. By comparison, Whisper adopts  
 1334 a Fourier sine-series approximation and exploits  
 1335 the periodicity of trigonometric functions to fuse  
 1336 the Bts operation with the non-linear approxima-  
 1337 tion Cmp, resulting in the proposed CBts protocol.  
 1338 This design eliminates the need to evaluate two sep-  
 1339 arate polynomials. As a result, Whisper achieves  
 1340 an approximately 27% improvement in runtime  
 1341 efficiency, even in large-input regimes.

1342 **Dataset Generation.** We further evaluate the  
 1343 end-to-end efficiency of fine-tuned dataset gener-  
 1344 ation on the GPT-2<sub>medium</sub> model in Fig. 15. Under  
 1345 the same setting where SIMD is used for cipher-  
 1346 text packing, Whisper benefits from the inherent  
 1347 parallelism of the Linear-Attention architecture,  
 1348 enabling the simultaneous processing of multiple  
 1349 prompt requests. In contrast, approaches based  
 1350 on Self-Attention can only execute prompt infer-  
 1351 ence sequentially. In large-scale fine-tuned dataset  
 1352 generation tasks, Whisper integrates protocol op-  
 1353 timizations across both linear and nonlinear lay-  
 1354 ers, achieving an 108.40-271.01 $\times$  improvement in  
 1355 throughput and reducing runtime by 4.39-23.30 $\times$

1356 compared to the SOTA PPI scheme. These results  
 1357 demonstrate that Whisper can generate fine-tuned  
 1358 datasets more efficiently, thereby effectively sup-  
 1359 porting the practical deployment of PPFT.