

THOR: Secure Transformer Inference with Homomorphic Encryption

Jungho Moon
Hanyang University
email: moonjunggho@hanyang.ac.kr

Xiaoqian Jiang
University of Texas,
Health Science Center at Houston
email: Xiaoqian.Jiang@uth.tmc.edu

Dongwoo Yoo
Yonsei University
email: aydw0507@yonsei.ac.kr

Miran Kim
Hanyang University
email: miran@hanyang.ac.kr

Abstract—As language models are increasingly deployed in cloud environments, privacy concerns have become a significant issue. To address this, we design THOR, a secure inference framework for transformer models on encrypted data. Specifically, we first propose new fast matrix multiplication algorithms based on diagonal-major order encoding and extend them to parallel matrix computation through the compact ciphertext packing technique. Second, we design efficient protocols for secure computations of four non-linear functions such as softmax, LayerNorm, GELU, and Tanh, by integrating advanced underlying approximation methods with tailored optimizations. Our matrix multiplication algorithms reduce the number of key-switching operations in the linear layers of the attention block in the BERT-base model by up to 14.5x, compared to the state-of-the-art HE-based secure inference protocol (Park et al., Preprint). Combined with cryptographic optimizations, our experimental results demonstrate that THOR provides secure inference for the BERT-base model with a latency of 10.43 minutes on a single GPU, while maintaining comparable inference accuracy on the MRPC dataset.

1. Introduction

The rapid advancement of *transformer models* has changed the way machines understand and generate sequence data, such as human language. These models are mainly based on *attention mechanism* [1], which allows the model to focus dynamically on different entities in the input sequence. This capability enables transformers to weigh the importance of various entities in a sequence, regardless of their position, thereby capturing long-range dependencies more effectively than traditional recurrent neural networks (RNNs) and long short-term memory networks (LSTMs). The introduction of BERT (Bidirectional Encoder Representations from Transformers) [2] exemplified the power of transformer models, leading to significant advancements in various natural language processing tasks, including question answering, named entity recognition, and translation.

However, substantial computational demands of deep learning pose challenges for practical applications. As a result, deploying transformer models in a cloud environment is a natural choice.

However, this approach has significant concerns when dealing with sensitive data, such as protected health information. The main issue is the potential exposure of such data during inference. Traditional methods of model deployment on cloud computing platforms, such as OpenAI’s GPT-4o [3] or Meta’s LLaMA [4], do not adequately address these privacy concerns, highlighting the need for new approaches to ensure data confidentiality.

To address privacy concerns, several cryptographic techniques such as secure multi-party computation (MPC), homomorphic encryption (HE), or their combination, have been adopted to secure transformer inference. In their scenario, the service provider holds a fine-tuned transformer model and provides inference services to clients whose input data is sensitive. Recent studies [5, 6, 7, 8, 9, 10] utilize oblivious transfer-based or secret sharing-based MPC techniques to protect private information across all layers, including the clients’ input. Although these interactive approaches achieve faster execution times, they still incur high communication costs due to the inherent overhead of MPC.

Several studies have explored secure transformer inference under HE. Although HE-based methods come with a substantial performance overhead, they allow the service provider (or server) to perform computations on encrypted data without decryption. THE-X [11] is the first HE-based transformer inference system; however, it sacrifices accuracy by employing HE-compatible replacements for non-linear functions. Most notably, its protocol delegates the computations of ReLU to the client, thereby revealing the comparison results involved in softmax and GELU approximations. Two recent works [12, 13] are HE-based non-interactive systems that perfectly protect data privacy during transformer inference. However, their works focus on the BERT-tiny model (a distilled student model of BERT with 4.4 million parameters), which limits the scalability and efficiency of matrix multiplication algorithms when

extended to larger models. In addition, they do not fully utilize all plaintext slots during matrix multiplications, leading to substantial performance overhead. Furthermore, their non-linear approximation methods are restricted to a small domain, reducing effectiveness for larger models. On the other hand, NEXUS [14] evaluated the BERT-base model (110M parameters [2]). Although NEXUS is capable of making inferences with high throughput, it requires a substantial number of operations even for a single prediction, resulting in high latency. This work employs various packing methods, such as component-wise, row/column-wise, and diagonal-wise packings; however, the details of changes in the packing structure during secure inference are not fully explained. In particular, the output of one encoder layer is used as the input for the plaintext-ciphertext matrix multiplication in the subsequent attention layer, requiring substantial computational costs to convert it into component-wise format.

Our contributions. The main challenges to achieving accurate and low-latency inference for large transformer models under HE are efficient homomorphic matrix multiplications required for attention mechanisms and secure computation of non-linear operations. In this study, we present THOR, a secure Transformer inference framework that leverages Homomorphic encryption to ensure data privacy. To the best of our knowledge, this is the first fully non-interactive secure transformer inference framework for the BERT-base model. Our contributions are three-fold:

Homomorphic matrix multiplication. We propose computation-efficient homomorphic matrix multiplication algorithms optimized for transformer architectures, aimed at reducing the computational overhead associated with encrypted operations. In a nutshell, we formulate matrix multiplication as entrywise additions and multiplications along matrix diagonals. In the context of HE, we leverage all plaintext slots by compactly packing multiple diagonal vectors and processing computations over packed ciphertext. Based on this insight, we introduce a plaintext-ciphertext matrix multiplication algorithm that performs parallel matrix multiplications between submatrices and aggregates the intermediate results to produce the final output. We also design a ciphertext-ciphertext matrix multiplication along with a lazy rotation optimization that reduces the number of rotations on input ciphertexts. Instead of performing rotations directly on the inputs, we defer these operations to the intermediate results by applying rotations of the same amounts. As a result, we can reduce the number of ciphertext rotations by roughly half, thus enhancing the computational efficiency.

Accurate and efficient non-linear homomorphic computations. Homomorphic evaluation of non-linear functions is an important factor in achieving accurate inference. We first propose an efficient method for evaluating the softmax function

over large intervals using a *two-phase method*: approximating the exponential function within a small domain and then performing repeated squaring-and-normalization operations. Specifically, we adopt Goldschmidt’s iterative method to approximate the inverse and adaptively adjust the relaxation coefficient to achieve faster convergence. This structured approach enables efficient computation while preserving high precision, even on large intervals. We also apply the adaptive iterative method to efficiently compute the inverse square root for the LayerNorm evaluation. Another contribution is the evaluation of the GELU function. We approximate it using a composition of two low-degree polynomials and leverage the Paterson-Stockmeyer algorithm for efficient polynomial evaluation.

Enc-to-end transformer inference. We design an end-to-end HE-based framework for secure transformer inference. The experimental results show that THOR provides secure inference on $n = 128$ tokens with the BERT-base model with a latency of 10.43 minutes on a single GPU, while maintaining a comparable inference accuracy.

2. Preliminaries

Notation. The binary logarithm will be denoted by $\log(\cdot)$. We denote by $[n]$ the set $\{0; 1; \dots; n - 1\}$, where $n \geq \mathbb{N}$. $a \pmod n$ or $[a]_n$ denotes the unique integer r such that $r \in [n]$ and $r = a \pmod n$. The entry of the i -th row and j -th column of an $m \times n$ matrix A is denoted by $A[i; j]$ or A_{ij} . A_i and A^j denote the i -th row and j -th column vectors. $hu; vi$ denotes the inner product.

2.1. Homomorphic Encryption

Homomorphic Encryption (HE) is a cryptographic primitive that allows computations to be performed directly on encrypted data, without the need for decryption. Since Gentry’s breakthrough on HE [15], numerous HE schemes with various properties have been proposed. Among them, the Cheon-Kim-Kim-Song (CKKS) scheme [16] has received significant attention due to its support for approximate arithmetic on encrypted complex numbers. The CKKS scheme also allows operations to be executed in a single instruction multiple data (SIMD) manner, enabling efficient parallel computation of large-scale datasets. Specifically, a plaintext vector $\mathbf{m} \in \mathbb{C}^s$ is first transformed into a plaintext polynomial, which is then encrypted as a ciphertext, where $s = N/2$ represents the number of plaintext slots for a ring dimension N .

The CKKS scheme supports the following basic arithmetic and advanced operations:

Vectorized arithmetic operations: Element-wise addition and multiplication can be performed between two ciphertexts (Add and Mult), or between a ciphertext and a plaintext (SAdd and SMult)

Rotation: Given a ciphertext ct representing $\mathbf{m} \in \mathbb{C}^S$, it returns a ciphertext representing the left rotation $\text{rot}^r(\mathbf{m})$ by r amounts. We abuse the same notation for the operation on ciphertexts.

Complex conjugation: This operation applies element-wise complex conjugation to a ciphertext.

Key-switching: One can transform a ciphertext with a secret key s into a ciphertext that encrypts the (approximately) same message but with a different key s^ℓ . For example, ciphertext-ciphertext multiplication involves a key switching operation to obtain a normal ciphertext with respect to the original secret key after performing the tensor product of the inputs.

Bootstrapping: The bootstrapping operations [17] refreshes a ciphertext by increasing its coefficient modulus, allowing further operations without significant loss of precision.

2.2. Transformer

The transformer architecture [1] is designed to handle sequence-to-sequence tasks effectively, particularly in natural language understanding (NLU) and generation. The transformer is based on an encoder-decoder architecture, with both parts having a similar structure. Therefore, we mainly focus on the encoder below. The transformer takes as input $X \in \mathbb{R}^{n \times d}$, where n indicates the number of tokens and d denotes the model dimension.

Encoder. The encoder is composed of a stack of L layers, each comprising *multi-head attention* followed by two normalizations, with feed-forward layers in between. Multi-head attention operates as follows: (i) The input sequence $X \in \mathbb{R}^{n \times d}$ is linearly projected to H heads so that $Q_h = XW_{Q,h}$; $K_h = XW_{K,h}$ and $V_h = XW_{V,h}$ with the parameter matrices $W_{Q,h}, W_{K,h} \in \mathbb{R}^{d \times d_k}$, and $W_{V,h} \in \mathbb{R}^{d \times d_v}$ for $h \in [H]$. (ii) The scaled dot-product attention is computed in each head:

$$\text{Att}_h = \text{Softmax} \left(\frac{Q_h K_h^\top}{d_k} \right) V_h;$$

(iii) The results of multiple heads are concatenated and linearly transformed by $W^O \in \mathbb{R}^{Hd_v \times d}$ to produce the final result. Following the parameters from [1], we use $d_k = d_v = d/H$. A feed-forward layer consists of two linear transformations with a GELU (Gaussian Error Linear Unit function) activation in between.

Pooler and classifier. The first token from the L encoder layers (corresponding to the first row of the result matrix) is fed into a pooler, which is a dense layer followed by a Tanh activation. The output is then fed into a classifier, which is another dense layer.

3. Basic Matrix Multiplication

In this section, we introduce HE-friendly matrix multiplication algorithms, which serve as the foundation of the linear layers in secure transformer inference.

3.1. Definitions & Lemmas

Definition 1 (Upper diagonal vector). For a matrix $A = (a_{ij}) \in \mathbb{R}^{a \times b}$ and $0 \leq k < \min\{a, b\}$, the k -th upper diagonal vector $U_k(A) \in \mathbb{R}^{\max\{a, b\}}$ is defined by

$$U_k(A)[t] = A[t \pmod{a}; k + t \pmod{b}]$$

for $0 \leq t < \max\{a, b\}$.

Definition 2 (Lower diagonal vector). For a matrix $B = (b_{ij}) \in \mathbb{R}^{b \times c}$ and $0 \leq \ell < \min\{b, c\}$, the ℓ -th lower diagonal vector $L_\ell(B) \in \mathbb{R}^{\max\{b, c\}}$ is defined by

$$L_\ell(B)[t] = B[\ell + t \pmod{b}; t \pmod{c}]$$

for $0 \leq t < \max\{b, c\}$.

We start with useful lemmas that characterize the upper and lower diagonal vectors, which are later used for matrix multiplications in private transformer inference.

Lemma 1. For any $a, b \in \mathbb{R}^n$ and $r \in \mathbb{Z}$, we have $h \cdot a; b \cdot i = h \cdot \text{rot}^r(a); \text{rot}^r(b) \cdot i$. That is, $\text{rot}^r(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is an isometry on the real inner product space \mathbb{R}^n .

Lemma 2. For any $A \in \mathbb{R}^{a \times b}$ and $k \in \mathbb{Z} \setminus [0; \min\{a, b\}]$, we have $U_k(A) = L_k(A^\top)$.

Lemma 3. For $m \leq n$, we have

$$L_\ell(A) = \begin{cases} m + \ell \cdot (U_{m-\ell}(A)) & \text{if } A \in \mathbb{R}^{m \times n}; \\ n + \ell \cdot (U_{m-\ell}(A)) & \text{if } A \in \mathbb{R}^{n \times m}; \end{cases}$$

Lemma 4. For a square matrix $A \in \mathbb{R}^{n \times n}$ and any $t \in [n]$, we have $(U_k(A)[t])_{k \in [n]} = A_t \in \mathbb{R}^n$.

3.2. Upper-lower Matrix Multiplication

We first propose matrix multiplication algorithms for the case when the upper and lower diagonal vectors of two input matrices are provided. In a nutshell, it follows from Lemma 1 that given two square matrices of A and B of size n , for any $r; t \in [n]$, we have $h \cdot \text{rot}^t(A_{r+t}); \text{rot}^t(B^t) \cdot i = h \cdot A_{r+t}; B^t \cdot i = C[r+t; t]$, which corresponds to the t -th entry of $L_r(C)$. In a SIMD fashion, $\text{rot}^t(A_{r+t})$ is obtained from the t -th entries of $\text{rot}^r(U_{\cdot-r}(A))$ for $\cdot \in [n]$. Similarly, $\text{rot}^t(B^t)$ is obtained from the t -th entries of $L_\ell(B)$ for $\ell \in [n]$. Thus, this can be formulated as the following propositions.

Proposition 1. For $n \leq m$, suppose that n is divisible by m . Given two matrices $A \in \mathbb{R}^{n \times m}$ and $B \in \mathbb{R}^{m \times n}$, for $r \in [n]$, the r -th lower diagonal vector of the matrix AB can be computed as follows:

$$L_r(AB) = \sum_{\ell \in [m]} \text{rot}^r(U_{[\cdot-r]_m}(A)) \cdot L_\ell(B); \quad (1)$$

Proof. See Appendix A.1. \square

Proposition 2. For $n \leq m$, suppose that n is divisible by m . Given two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times n}$, for

$r \geq [m]$, the r -th lower diagonal vector of the matrix AB can be computed as follows:

$$L_r(AB) = \sum_{i \geq [n]} (U_{i-r}^0(A) \cdot L_i(B));$$

where the upper diagonal vector $U_k^0(A) \in \mathbb{R}^n$ is defined by

$$U_k^0(A) = \begin{cases} U_k(A) & \text{if } k < m; \\ m \cdot \lfloor k/m \rfloor (U_{(k \bmod m)}(A)) & \text{if } m \leq k < n; \end{cases}$$

Proof. See Appendix A.2. \square

3.3. Lower-lower Matrix Multiplication

We now propose matrix multiplication algorithms for the case when the lower diagonal vectors of two input matrices are provided.

Definition 3. For a matrix $A \in \mathbb{R}^{n \times m}$ with $n \geq m$, the lower diagonal vector $L_k^0(A) \in \mathbb{R}^n$ is defined by

$$L_k^0(A) = \begin{cases} L_k(A) & \text{if } k < m, \\ m \cdot \lfloor k/m \rfloor (L_{(k \bmod m)}(A)) & \text{if } m \leq k < n; \end{cases}$$

Corollary 1. For $n \geq m$, suppose that n is divisible by m . Given two matrices $A \in \mathbb{R}^{n \times m}$ and $B \in \mathbb{R}^{m \times n}$, for $r \geq [n]$, the r -th lower diagonal vector of the matrix AB can be computed as follows:

$$L_r(AB) = \sum_{i \geq [m]} (L_{[n \cdot i + r]}^0(A) \cdot L_i(B)); \quad (2)$$

Proof. See Appendix A.3. \square

Corollary 2. For $n \geq m$, suppose that n is divisible by m . Given two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times n}$, for $r \geq [m]$, the r -th lower diagonal vector of the matrix AB can be computed as follows:

$$L_r(AB) = \sum_{i \geq [n]} (L_{[m \cdot i + r]}(A) \cdot L_i(B));$$

Proof. See Appendix A.4. \square

4. Homomorphic Matrix Multiplication

In this section, we explain how to evaluate the proposed matrix multiplication algorithms using a combination of compact ciphertext packing and SIMD techniques. The main goal of our algorithms is to minimize the number of key-switching operations, particularly rotations, which is the dominant process in homomorphic operations.

4.1. Matrix Multiplication on Packed Ciphertexts

We introduce a compact packing method that packs multiple diagonal vectors of an input matrix into a single ciphertext, ensuring that the output diagonals of the matrix product are efficiently packed within the resulting ciphertext. Furthermore, multiple diagonals from different matrices fit into a single plaintext to maximize the slot usage, enabling parallel matrix computations in a SIMD manner.

4.1.1. Compact Matrix Packing. For simplicity of presentation, we assume that the matrix $A \in \mathbb{R}^{n \times m}$ is provided in plaintext, while the matrix $B \in \mathbb{R}^{m \times n}$ is provided in encrypted form, with $n \geq m$. If the number of plaintext slots s is equal to n , each n -dimensional diagonal vector of A and B can be encoded in a single ciphertext. Since it requires m scalar multiplications to compute one diagonal vector of the output matrix by evaluating Eq. (1), the total computation requires mn scalar multiplications in total.

We now assume that the slot parameter s is larger than n . A straightforward approach is to pad zeros after the n slots and process with the algorithm as described above. However, this method is space-inefficient since the slots after n remain unused. To address this inefficiency, we aim to produce multiple diagonal vectors of AB within a single ciphertext by utilizing all available slots. This can be achieved by compactly packing multiple diagonal input vectors and performing SIMD computations over the packed ciphertexts. Specifically, given that $s = cn$ for some integer c , we assume that c many lower diagonal vectors of B are packed into a single ciphertext. Let $\text{ct} \cdot B_j$ be an encryption of the stacked diagonals $(L_{cj}(B); L_{cj+1}(B); \dots; L_{c(j+1)-1}(B))$ for $j \geq [m/c]$. Then, for $r \geq [n/c]$, an encryption of $(L_{cr}(C); L_{cr+1}(C); \dots; L_{c(r+1)-1}(C))$ can be obtained by evaluating

$$\sum_{i \geq [c] \cdot j \geq [m/c]} \text{SMult}(\text{ct} \cdot B_j; \text{pt} \cdot A_{j \cdot c + i}); \quad (3)$$

where $\text{pt} \cdot A_{j \cdot c + i}$ denotes a plaintext of the internally rotated $(i + cj)$ -th upper diagonal $(c^i(U_{i+cj}(A)); c^{i+1}(U_{i+cj}(A)); \dots; c^{(r+1)-1}(U_{i+cj}(A)))$. Although the number of rotations is increased, the total number of scalar multiplications is reduced to $mn=c$.

4.1.2. Parallel Matrix Multiplications. Assume that s is divisible by n and let $s = s \cdot n$. Given that $m \geq n$, we aim to perform matrix multiplications of $A^{(z)} \in \mathbb{R}^{n \times m}$ and $B^{(z)} \in \mathbb{R}^{m \times n}$ for $z \geq [s]$. Let $A = fA^{(z)}g_{z \geq [s]}$ and $B = fB^{(z)}g_{z \geq [s]}$. For $k \geq [m]$, we define the k -th interlaced upper diagonal vector $\tilde{U}_k(A) \in \mathbb{R}^s$ as

$$\tilde{U}_k(A)[t] = U_k(A^{(t \bmod s)})[bt = sc] \quad (4)$$

for $t \geq [s]$. In other words, the entries of the upper diagonal vectors $U_k(A^{(z)})$ are interlaced within the slots. The t -th interlaced lower diagonal vector $\tilde{L}_t(B) \in \mathbb{R}^s$ can be defined in a similar way. We then perform s -many matrix multiplications to obtain the matrix $C^{(z)} = A^{(z)}B^{(z)}$ in one shot by slightly modifying the single matrix multiplication algorithm in Eq. (1). Specifically, the r -th interlaced lower diagonal vector of the output matrices $C^{(z)}$ can be computed as follows: $\tilde{L}_r(C) = \sum_{i \geq [m]} r^s (U_{i-r}(A)) \cdot \tilde{L}_i(B)$.

In practice, the BERT-base model has $H = 12$ heads, so H many matrix multiplications are executed in parallel. Given that each diagonal vector of the intermediate matrix in transformer has a size of $n = 128$, we get a power of two integer c that is less than $s = (nH) = 2^{15} = (2^7 \cdot 12)$, resulting in $c = 16$.

4.2. Plaintext-ciphertext Matrix Multiplication

We observe that all the plaintext-ciphertext matrix multiplications in the linear layers of the BERT model can be expressed as $C = AB$, where $A \in \mathbb{R}^{d \times d}$ is provided in plaintext (e.g., a large weight matrix) and $B \in \mathbb{R}^{d \times n}$ is provided in encrypted form (e.g., an encrypted output from the previous layer), given that $n \mid d$ (see Section 6.2 for details). Without loss of generality, we assume that d is divisible by n and denote $d = d \cdot n$; otherwise, zeros can be padded. Additionally, we assume $s = cd$, where c is an integer that divides n , and we denote $\# = n/c$.

4.2.1. Extension of Basic Matrix Multiplication. One approach is to modify the basic upper-lower matrix multiplication algorithms in Section 3.2 to leverage the d -dimensional diagonal vectors of A and B . For $r \in [n]$, the r -th lower diagonal of C can be computed as $L_r(C) = \sum_{j \in [d]} L_r(U_{\cdot, r})(A) \cdot L^{\ell}(B)$, where $L^{\ell}(B) \in \mathbb{R}^d$ is defined by

$$L^{\ell}(B) = \begin{cases} L_{\cdot}(B) & \text{if } \cdot < n; \\ n \cdot b^{\cdot} = nc(L_{\cdot \bmod n}(B)) & \text{if } n \leq \cdot < d; \end{cases}$$

As described in Section 4.1.1, we assume that c -many d -dimensional diagonal vectors of B are fully packed into a ciphertext. Let $\text{ct}:B_j$ be an encryption of $(L_{c_j}(B); L_{c_{j+1}}(B); \dots; L_{c_{(j+1)-1}}(B))$ for $j \in [\#]$. We first need rotations within diagonal vectors to obtain encryptions of $L^{\ell}(B)$. Specifically, we generate a ciphertext $\text{ct}:B_{nt+j}$ of

$$\begin{aligned} & (L_{c_j}^{\ell}(B); L_{c_{j+1}}^{\ell}(B); \dots; L_{c_{(j+1)-1}}^{\ell}(B)) \\ & = (L_{-nt+c_j}^{\ell}(B); L_{-nt+c_{j+1}}^{\ell}(B); \dots; L_{-nt+c_{(j+1)-1}}^{\ell}(B)) \end{aligned}$$

for $t \in [1; d]$. That is, we need to rotate within a subslot, but the CKKS scheme only supports 1-dimensional rotations. To achieve this, we perform an *internal rotation*¹, which requires two homomorphic rotations and two scalar multiplications. Since we repeat this operation for each rotation amount $t \in [1; d]$ on $\#$ many input ciphertexts $\text{ct}:B_j$, this step requires $2(d-1)\# = 2d=c$ rotations and $2(d-1)\# = 2d=c$ scalar multiplications.

Similar to Eq. (3), we obtain ciphertexts for each set of c -many output diagonals by evaluating

$$\sum_{j \in [d]} \text{SMult}(L^{\ell}(\text{ct}:B_j); \text{pt}:A_{j, \cdot, r});$$

for $r \in [\#]$, using appropriate plaintexts $\text{pt}:A_{j, \cdot, r}$ from the matrix A . For each $j \in [d=c]$, the ciphertext $\text{ct}:B_j$ is rotated $(c-1)$ times, requiring $(c-1)(d-c) = d$ rotations. In conclusion, this algorithm takes approximately $2d=c + d = O(d)$ rotations and $2d=c + c(d-c)\# = 2d=c + dn=c$ scalar multiplications.

1. For example, given that a ciphertext representing the vector $(u_0; u_1; u_2; u_3; v_0; v_1; v_2; v_3)$, it can be rotated to the left by one and right by three. Each ciphertext is applied by the mask to extract $(u_1; u_2; u_3; 0; v_1; v_2; v_3; 0)$ and $(0; 0; 0; u_0; 0; 0; 0; v_0)$. At the end, we get a ciphertext of $(u_1; u_2; u_3; u_0; v_1; v_2; v_3; v_0)$ by summing the two intermediate ciphertexts.

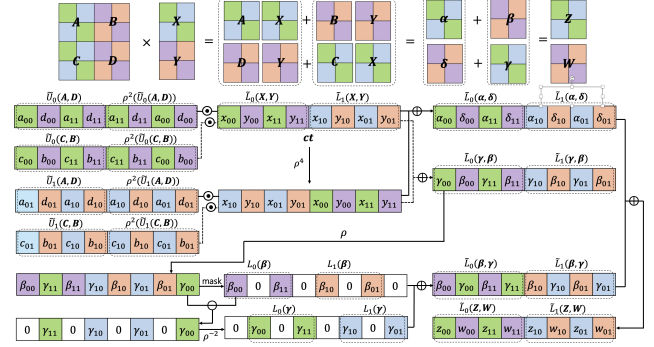


Figure 1. Diagonal block plaintext-ciphertext matrix multiplication. The ciphertext ct represents a stacked vector of $L_0(X; Y)$ and $L_1(X; Y)$. For simplicity, we do not include the complexification.

Algorithm 1 PLAINTEXT-CIPHERTEXT MATMUL

Input: Plaintexts $\text{pt}:A_{i,j}; \cdot, r$ of the internally rotated upper diagonal vectors of $A \in \mathbb{R}^{d \times d}$; ciphertexts $\text{ct}:B_j$ of the lower-diagonally stacked vectors of $B \in \mathbb{R}^{d \times n}$ for $i \in [d=n]; j \in [n=c]; \cdot \in [c]$

Output: Ciphertexts $\text{ct}:C_r$ of the lower-diagonally stacked vectors of C for $r \in [n=c]$

- 1: $d = d \cdot n, \# = n/c$
- 2: **for** $i = 0$ to $d-1$ **do**
- 3: **for** $r = 0$ to $\#-1$ **do**
- 4: **for** $j = 0$ to $\#-1, \cdot = 0$ to $c-1$ **do**
- 5: $\text{ct}:i_{j, \cdot, r} \leftarrow \text{SMult}(L^{\ell}(\text{ct}:B_j); \text{pt}:A_{i, j, \cdot, r})$
- 6: $\text{ct}:i_{r, \cdot} \leftarrow \sum_{j \in [d], j \in [n]} \text{ct}:i_{j, \cdot, r}$
- 7: **for** $r = 0$ to $\#-1$ **do**
- 8: $\text{ct}:C_r \leftarrow \text{SMult}(\text{ct}:i_{r, \cdot}; \text{pt}:0)$
- 9: **for** $i = 1$ to $d-1$ **do**
- 10: $\text{ct}:i_{r, R} \leftarrow \text{SMult}(\text{ct}:i_{r, \cdot}; \text{pt}:i)$
- 11: $\text{ct}:i_{r, L} \leftarrow \text{ct}:i_{r, R} \cdot \text{ct}:i_{r, R}$
- 12: $\text{ct}:i_{r, R}^{\ell} \leftarrow L^{\ell}(d \cdot n)(\text{ct}:i_{r, R}) + L^{\ell}(i)(\text{ct}:i_{r, L})$
- 13: $\text{ct}:C_r \leftarrow \text{ct}:C_r + \text{ct}:i_{r, R}^{\ell}$
- 14: **return** $\text{ct}:C_r$

4.2.2. Diagonal Block Matrix Computation. To further enhance efficiency, we partition the two matrices into $n \times n$ blocks and perform the matrix product using the n -dimensional diagonal vectors from the submatrices. Specifically, we denote the $(i; k)$ -th submatrix of A by $A_{ik} \in \mathbb{R}^{n \times n}$ and the k -th submatrix of B by $B_k \in \mathbb{R}^{n \times n}$ for $i; k \in [d]$. Then, the submatrices in the resulting matrix C are calculated by blockwise multiplication, where $C_k = \sum_i A_{ik} B_k$ for $k \in [d]$. This blockwise computation is efficiently implemented using HE through the following three steps: (i) perform parallel matrix multiplications between submatrices in the same diagonal direction, (ii) rotate the intermediate results by an appropriate amount to align the partial results for final submatrix output into the same slots, and (iii) sum the aligned values to obtain the final blockwise product.

A concrete example is provided in Fig. 1, and the detailed procedure is described in Algorithm 1. First, for each

TABLE 1. COMPARISON OF PLAINTEXT-CIPHERTEXT MATRIX MULTIPLICATION METHODS OF AB WHERE $A \in \mathbb{R}^{d \times d}$ AND $B \in \mathbb{R}^{d \times n}$ WITH $n = d$ AND $s = cd$.

Method	#Smult	#Rot
Basic	$\frac{2d}{c} + \frac{dn}{c}$	$d + \frac{2d}{c}$
DiagBlock	$\frac{d}{c} + \frac{dn}{c}$	$n + \frac{2d}{c}$

$\ell \in [n]$, the ℓ -th lower diagonal vectors from the matrices $B_0; \dots; B_{d-1}$ are interlaced to form the vector $\tilde{L}_\ell(B) \in \mathbb{R}^d$, as shown in Eq. (4). Next, c such vectors are concatenated to be fully-packed into a single ciphertext. Specifically, for $j \in [m]$, we define

$$\tilde{L}_{[c(j+1)]}(B) = (\tilde{L}_{c(j)}(B); \dots; \tilde{L}_{c(j+1)-1}(B)) \in \mathbb{R}^s$$

and denote the corresponding ciphertext as $\text{ct}:B_j$. Similarly, we define $A^{(i)} = fA_{i,0}; A_{i+1,1}; \dots; A_{i+d-1,d-1}g$ as the set of submatrices in the i -th diagonal order for $i \in [d]$. The plaintext of their interlaced compactly packed upper diagonal vectors for $A^{(i)}$ is denoted as $\text{pt}:A_{i,j}; r$ for $\ell \in [d]$, and $j; r \in [m]$, as explained in Section 4.1.1.

Let $C_{i,k} = A_{i,k}B_k$ and $C^{(i)} = fC_{i,0}; C_{i+1,1}; \dots; C_{i+d-1,d-1}g$. For each $i \in [d]$ and $r \in [m]$, we begin by computing the concatenated lower diagonal vectors $L_{[c(r+1)]}(C^{(i)})$ (lines 5 and 6 in Algorithm 1). The parallel submatrix multiplication requires $(c-1)m$ rotations and $n \times d \times m = dn = c$ scalar multiplications. Similar to an internal rotation operation, these results are internally rotated within the output diagonals of the same order, yielding the product result $fC_{0,d-1}; \dots; C_{i,0}; \dots; C_{d-1,d-1}g$ (lines 10 to 12). This process involves $2(d-1)m = 2d = c$ rotations and $(d-1)m = d = c$ scalar multiplications to compute all the aligned results of the submatrix products. Finally, each lower diagonal vector of the intermediate results is added to produce the final output (line 13). We note that this ciphertext $\text{ct}:C_r$ represents the plaintext vector of $\tilde{L}_{[c(r+1)]}(C)$ where $C = fC_0; C_1; \dots; C_{d-1}g$.

Although additional computations are required to rotate plaintext slots for interactions between values of different submatrices, intermediate submatrix products can be parallelized in a SIMD manner. Consequently, the number of rotations is reduced to approximately $n + 2d = c = O(n)$, enabling the block matrix multiplication algorithm to outperform the naive method. We present the results in Table 1.

4.3. Ciphertext-ciphertext Matrix Multiplication

When two input matrices are provided in encrypted form, the upper-lower matrix multiplication algorithms can be directly applied in a similar manner to the plaintext-ciphertext matrix multiplication method. However, it is very challenging to generate encryptions of the plaintexts $\text{pt}:A_{j,r}$ as explain in Section 4.1.1 from the stacked upper diagonal vectors of A . For example, consider the product AB with $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times n}$. Suppose that

we are given the diagonally stacked ciphertext representing $(U_{c(j)}(A); U_{c(j+1)}(A); \dots; U_{c(j+1)-1}(A))$ for $j \in [m]$, where $m = m/c$. We apply the the internal rotation method to generate ciphertexts that are internally rotated by $r \in [1;n]$ positions for each input ciphertext. This procedure requires $2(n-1)m$ rotations and scalar multiplications. After that, it takes roughly mn scalar multiplications to generate the packed ciphertexts in the form of $(c^{(j)}(U_k(A)); c^{(j+1)}(U_k(A)); \dots; c^{(j+1)-1}(U_k(A)))$ for each $k \in [m]$ and $\ell \in [n=c]$. Instead of this naive extension, we leverage both the lower diagonals of the input matrices and perform the matrix computation in parallel to enhance performance.

4.3.1. Packed Lower-lower Matrix Multiplication. Suppose that we are given the ciphertext $\text{ct}:A_j$ representing the lower-diagonally stacked vectors $(L_{c(j)}(A); L_{c(j+1)}(A); \dots; L_{c(j+1)-1}(A))$ for $j \in [m]$. First, two rotated ciphertexts are generated as follows:

$$\begin{aligned} \text{ct}_{j,r}^1 &= \text{rot}^{(r)}(\text{ct}:A_j) \\ \text{ct}_{j,r}^0 &= \text{rot}^{(r+c)}(\text{ct}:A_j) \end{aligned}$$

Then $\text{ct}_{j,r}^1$ is applied by two masks γ_0 and γ_1 to extract the first $(n-\ell)$ entries of $(L_{c(j+r)}(A))$ for each $\ell \in [c-r < c]$ and $0 \leq r < [c]$, respectively. Similarly, $\text{ct}_{j,r}^0$ is masked to extract the last ℓ entries of each $(L_{c(j+r)}(A))$ using γ_2 and γ_3 . This procedure requires approximately $2mn$ rotations and $4mn$ scalar multiplications. Next, these ciphertexts are appropriately added together to obtain the internally rotated ciphertext $\text{ct}:A_{j,r}$ representing $(L_{c(j-\ell)}(A); L_{c(j-\ell+1)}(A); \dots; L_{c(j+1)-\ell}(A))$, i.e.,

$$\begin{aligned} \text{ct}:A_{j,r} &= (\text{ct}_{j,r-1}^1 \gamma_0 + \text{ct}_{j,r}^1 \gamma_1) \\ &+ (\text{ct}_{j,r-1}^0 \gamma_2 + \text{ct}_{j,r}^0 \gamma_3); \end{aligned} \quad (5)$$

We now assume that the ciphertext $\text{ct}:B_r$ contains c copies of the ℓ -th lower diagonal vector of B for $\ell \in [n]$. Corollary 2 implies that the diagonally stacked ciphertext of $(L_{c(j)}(C); L_{c(j+1)}(C); \dots; L_{c(j+1)-1}(C))$ can be obtained by evaluating

$$\text{ct}:C_j = \sum_{\ell \in [n]} \text{Mult}(\text{ct}:A_{j,r}; \text{ct}:B_r); \quad (6)$$

4.3.2. Lazy Rotation Technique. We introduce the lazy rotation optimization, which reduces the required number of rotations by roughly a factor of two. The main observation is that the same rotation amount ℓ is used for the ciphertexts $\text{ct}:A_{j,r}$ across different j . This consistency in rotation amounts simplifies the computation and enables efficient parallel processing.

We observe from Eq. (5) that the ciphertext $\text{ct}:A_{j,r}$ can be computed as

$$\begin{aligned} \text{ct}:A_{j,r} &= \text{ct}_{j,r-1}^1 \gamma_0 + \text{ct}_{j,r}^1 \gamma_1 \\ &+ \text{rot}^{(r)}(\text{ct}_{j,r-1}^0) \gamma_2 + \text{rot}^{(r)}(\text{ct}_{j,r}^0) \gamma_3 \\ &= \text{ct}_{j,r-1}^1 \gamma_0 + \text{ct}_{j,r}^1 \gamma_1 \\ &+ \text{rot}^{(r)}(\text{ct}_{j,r-1}^0 \gamma_2) + \text{ct}_{j,r}^0 \gamma_3; \end{aligned}$$

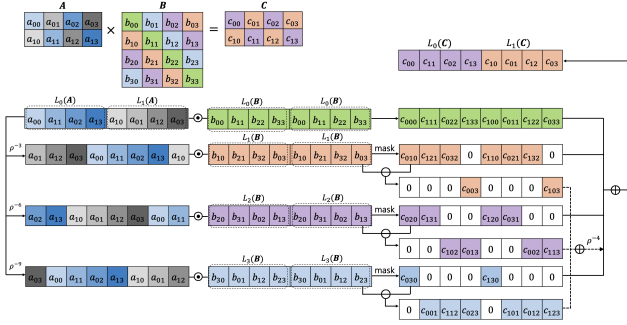


Figure 2. Ciphertext-ciphertext matrix multiplication. We denote $c_{ijk} = a_{ij}b_{jk}$. For simplicity, we do not include the complexification.

The first equality follows from the fact that $ct_{j,\cdot}^{\rho} = n(ct_{j,\cdot})$. By substituting this into Eq. (6), the ciphertext $ct:C_j$ can be expressed as follows:

$$\sum_{\ell \in [2[n]]} \text{Mult}(ct_{j,\cdot}^{\rho-\ell}; \rho_0 + ct_{j,\cdot}^{\rho-\ell}; ct:B_{\cdot}) + \sum_{\ell \in [2[n]]} \text{Mult}(n(ct_{j,\cdot}^{\rho-\ell}); n(\rho_2) + ct_{j,\cdot}^{\rho-\ell}; n(\rho_3)); ct:B_{\cdot})$$

Let $ct:C_{j,\cdot} = \text{Mult}(ct_{j,\cdot}; ct:B_{\cdot})$. Then the first term can be expressed as follows:

$$\sum_{\ell \in [2[n]]} (\text{Mult}(ct_{j,\cdot}^{\rho-\ell}; ct:B_{\cdot}) \rho_0 + \text{Mult}(ct_{j,\cdot}^{\rho-\ell}; ct:B_{\cdot}) \rho_1) = \sum_{\ell \in [2[n]]} (ct:C_{j,\cdot}^{\rho-\ell} \rho_0 + ct:C_{j,\cdot}^{\rho-\ell} \rho_1)$$

The second term can be expressed as follows:

$$n\left(\sum_{\ell \in [2[n]]} \text{Mult}(ct_{j,\cdot}^{\rho-\ell}; n(\rho_2) + ct_{j,\cdot}^{\rho-\ell}; n(\rho_3)); ct:B_{\cdot}\right) = n\left(\sum_{\ell \in [2[n]]} ct:C_{j,\cdot}^{\rho-\ell} n(\rho_2) + ct:C_{j,\cdot}^{\rho-\ell} n(\rho_3)\right)$$

The first equation follows from the fact that $ct:B_{\cdot} = n(ct:B_{\cdot})$.

We visualize the high-level idea of the lazy rotation technique in Fig. 2, and the detailed procedure is provided in Algorithm 2. The cost of the lazy rotation method is as follows. First, we apply one rotation to compute $ct_{j,\cdot}^{\rho}$ for each $j \in [m]$ and $1 \leq \rho < n$ (line 5). Next, we perform mn homomorphic multiplications between $ct_{j,\cdot}^{\rho}$ and $ct:B_{\cdot}$ (line 6). After that, we carry out $4mn$ scalar multiplications of the resulting ciphertexts with the masks $\rho_0, \rho_1, n(\rho_2)$ and $n(\rho_3)$ (lines 9 to 11). We note that the number of scalar multiplications can be reduced by $3mn$ using the fact that $ct:C_{j,\cdot}^{\rho} n(\rho_3) = ct:C_{j,\cdot}^{\rho} ct:C_{j,\cdot}^{\rho_0} ct:C_{j,\cdot}^{\rho_1} ct:C_{j,\cdot}^{\rho_2} n(\rho_2)$ (line 12). Finally, the output ciphertexts are appropriately summed with rotations by n positions (lines 14 to 16).

The main advantage of the lazy rotation technique is the reduction in the number of rotations from $2mn$ to mn .

Algorithm 2 CIPHERTEXT-CIPHERTEXT MATMUL

Input: Ciphertexts $ct:A_j$ of the lower-diagonally stacked vectors of $A \in \mathbb{Z}^{m \times n}$ for $j \in [m=c]$; ciphertexts $ct:B_{\cdot}$ of the replicated lower diagonal vectors of $B \in \mathbb{Z}^{n \times n}$ for $\cdot \in [n]$

Output: Ciphertexts $ct:C_j$ of the lower-diagonally stacked vectors of C

- 1: $m = c$
- 2: **for** $j = 0$ to $m - 1$ **do**
- 3: $ct:C_j = \text{Mult}(ct:A_j; ct:B_0)$
- 4: **for** $\rho = 1$ to $n - 1$ **do**
- 5: $ct_{j,\cdot}^{\rho} = n[\rho]_{c+n}(ct:A_j)$
- 6: $ct:C_{j,\cdot} = \text{Mult}(ct_{j,\cdot}^{\rho}; ct:B_{\cdot})$
- 7: **for** $\rho = 1$ to $n - 1$ **do**
- 8: **for** $j = 0$ to $m - 1$ **do**
- 9: $ct_{j,\cdot}^{\rho,0} = \text{SMult}(ct:C_{j,\cdot}^{\rho}; \rho_0)$
- 10: $ct_{j,\cdot}^{\rho,1} = \text{SMult}(ct:C_{j,\cdot}^{\rho}; \rho_1)$
- 11: $ct_{j,\cdot}^{\rho,2} = \text{SMult}(ct:C_{j,\cdot}^{\rho}; n(\rho_2))$
- 12: $ct_{j,\cdot}^{\rho,3} = ct:C_{j,\cdot}^{\rho} ct_{j,\cdot}^{\rho,0} ct_{j,\cdot}^{\rho,1} ct_{j,\cdot}^{\rho,2}$
- 13: **for** $j = 0$ to $m - 1$ **do**
- 14: $ct:C_j^{\rho} = \sum_{\ell=1}^{\rho < n} (ct_{j,\cdot}^{\rho-\ell,0} + ct_{j,\cdot}^{\rho-\ell,1})$
- 15: $ct:C_j^{\rho,0} = \sum_{\ell=1}^{\rho < n} (ct_{j,\cdot}^{\rho-\ell,2} + ct_{j,\cdot}^{\rho-\ell,3})$
- 16: $ct:C_j = ct:C_j + ct:C_j^{\rho} + n(ct:C_j^{\rho,0})$
- 17: **return** $ct:C_j$

TABLE 2. COMPARISON OF CIPHERTEXT-CIPHERTEXT MATRIX MULTIPLICATION METHODS OF AB WHERE $A \in \mathbb{Z}^{m \times n}$ AND $B \in \mathbb{Z}^{n \times n}$ WITH $n = m$ AND $s = cn$.

Method	Cplx	#Smult	#Mult		#Rot
			#TProd	#KS	
Basic	\times	$\frac{4mn}{c}$	$\frac{mn}{c}$	$\frac{m}{c}$	$\frac{2mn}{c}$
Lazy rotations	\times	$\frac{3mn}{c}$	$\frac{mn}{c}$	$\frac{2m}{c}$	$\frac{mn}{c}$
	\checkmark	$\frac{3mn}{2c}$	$\frac{mn}{2c}$	$\frac{2m}{c}$	$\frac{m(n+2)}{2c}$

When combined with the lazy relinearization technique (will be explained in Section 4.5), the number of key-switching operations during homomorphic multiplication is doubled to $2mn$ compared to the naive approach. However, this increase is relatively small compared to the reduction in rotations, so it does not result in significant efficiency degradation. We also remark that the ciphertext-ciphertext matrix multiplication can be parallelized in a SIMD manner without incurring additional costs, as explained in Section 4.1.2.

4.4. Row-wise Computation over Diagonal Representation

Assume that the upper diagonal vectors of A are given in encrypted form, where $ct:A_k$ denotes the encryption of the k -th upper diagonal vector $U_k(A)$. It follows from Lemma 4 that the vector formed by the t -th entries of the upper diagonal vectors corresponds to the t -th row of the matrix

TABLE 3. COMPARISON OF HOMOMORPHIC MATRIX MULTIPLICATIONS. THE CONCRETE NUMBERS ARE BASED ON THE BERT-BASE PARAMETERS: $n = 128$, $d = 768$, $d_k = 64$, $H = 12$, AND $s = 2^{15}$.

Operation	Method	#KeySwitch	
		Equation	Number
$fXW_{Q,h}; XW_{K,h}; XW_{V,h}g_{h2[H]}$ ($X \in \mathbb{R}^{n \times d}; W \in \mathbb{R}^{d \times d_k}$)	BOLT	$\frac{63nd}{s} + (\frac{s}{64n} - 1) \frac{nd_k}{s} - 3H$	297
	NEXUS	$6dd_kH$	3538944
	Powerformer	$(\frac{P}{3n + d_k} + (\frac{P}{3n + d_k} + \log(\frac{n}{d_k})) \frac{3H}{2}) \frac{d}{n}$	2612
	THOR	$\frac{n}{2} + 6(\frac{d}{2d_k} - 1) \frac{d_k}{c} + 3 \frac{d_k}{c}$	180
$fQ_hK_h^Tg_{h2[H]}$ ($Q_h; K_h \in \mathbb{R}^{n \times n=2}$)	BOLT	$(2n - 2 + n \log(\frac{n}{2})) H$	12264
	NEXUS	n^2H	196608
	Powerformer	$(\frac{5n}{2} + 4 \frac{P}{n}) \frac{H}{2}$	2196
	THOR	$(\frac{6n}{c} + 2) + \frac{n}{8}(\log c + 3) + \frac{2n}{c} + 2(n - 1) + \frac{3n}{2c}$	444
$f_hV_hg_{h2[H]}$ ($h \in \mathbb{R}^{n \times n}; V_h \in \mathbb{R}^{n \times n=2}$)	BOLT	$(2n - 1) \log n - H$	21420
	NEXUS		
	Powerformer	$(4n + 3 \frac{P}{14n} + 8 \frac{P}{n} + 2 \frac{P}{2n}) \frac{H}{2}$	4614
	THOR	$\frac{n}{4}(\log c + 3) + \frac{n}{c}(1 + \frac{n+2}{4})$	492
Concat($Att_0; \dots; Att_{H-1}$) W^O ($Att_h \in \mathbb{R}^{n \times n=2}; W^O \in \mathbb{R}^{d \times d}$)	BOLT	$\frac{31nd}{s} + (\frac{s}{32n} - 1) \frac{nd_k}{s} - H$	177
	NEXUS	$2d^2H$	14155776
	Powerformer	$4 \frac{P}{n} (\frac{d}{n})^2$	1656
	THOR	$\frac{d_k}{2} + 2 (\frac{d}{n} - 1) \frac{n}{c}$	118

A. Therefore, if the target function $f(X)$ is applied row by row to A and can be approximated by a function $\rho(X)$, this function can be used to approximate f for each row of A . Specifically, let $\rho = f(A)$ be the row-wise application of f to A . The homomorphically evaluated result $\rho(\text{ct}:A_k)$ corresponds to an approximate result for the upper diagonal vectors of ρ . This approach can be extended to a packed implementation, making it applicable for approximations such as the softmax and LayerNorm in secure transformer inference.

4.5. Optimization

Lazy homomorphic operations. We apply a lazy rescaling technique during plaintext-ciphertext multiplications by adjusting a scaling factor only once within the resulting ciphertext after multiple multiplications and additions (e.g., lines 5 and 6 in Algorithm 1). Additionally, we rely on a lazy key-relinearization technique during ciphertext multiplications, which enables tensor products between ciphertexts and summation over the extended ciphertexts, along with only a single relinearization via key-switching operation. For example, we perform key-switching operations after the computations in lines 14 and 15 of Algorithm 2. Consequently, by using this optimized approach, we reduce the number of key-switching operations during homomorphic multiplications from $\mathcal{M}n$ to $2\mathcal{M}$ (Table 2).

Complexification. For simplicity of presentation, we focus on Proposition 1 given that m is even. Eq. (1) can be

reformulated as follows:

$$\begin{aligned} L_r(C) &= \sum_{\substack{r \\ 2|m=2}} r(U_{[r]_m}(A)) \cdot L_r(B) \\ &\quad + r(U_{[r+m=2]_m}(A)) \cdot L_{r+m=2}(B) \\ &= \sum_{\substack{r \\ 2|m=2}} \text{Re}(r(U_{[r]_m}(A)) \cdot L_r(B)) \\ &= \text{Re}(\sum_{\substack{r \\ 2|m=2}} r(U_{[r]_m}(A)) \cdot L_r(B)); \end{aligned}$$

where $\bar{U}_k(A) = U_k(A) + \rho^{-1}U_{k+m=2}(A)$ and $\bar{L}_r(B) = L_r(B) + \rho^{-1}L_{r+m=2}(B)$. Since $\text{Re}(2x) = x + \text{conj}(x)$, we have $L_r(AB) = X + \text{conj}(X)$ where

$$X = \sum_{\substack{r \\ 2|m=2}} r(\bar{U}_{[r]_m}(A)) \cdot \bar{L}_r(B);$$

Thus, the number of scalar multiplications is reduced by half. A similar approach can be applied to all the proposed matrix multiplication algorithms.

4.6. Comparison

In Table 3, we compare the computational complexity of the proposed matrix multiplication algorithms for multi-head attention in the BERT-base model with BOLT [6], NEXUS [14], and Powerformer [13]. Since the key-switching operation turns out to be the dominant process in homomorphic operations [18], we count the number of the required key-switching operations during ciphertext multiplications and rotations in each matrix computation. To

ensure a fair comparison, we assume that all methods use the same HE parameters, with the ring dimension of $N = 2^{16}$ and the number of plaintext slots $S = 2^{15}$. Specifically, THOR uses the diagonal packing parameter $c = 16$ as noted in Section 4.1.2. THOR consistently outperforms the state-of-the-art HE-based approaches, BOLT and Powerformer, across all matrix operations.

5. Secure Non-linear Computations

5.1. Softmax

For a vector $\mathbf{x} = (x_i) \in \mathbb{R}^n$, the softmax function is defined as

$$\text{Softmax}(\mathbf{x}) = \left(\frac{\exp(x_i)}{\sum_{j \in [n]} \exp(x_j)} \right)_{i \in [n]}.$$

5.1.1. Our Approach at a High Level. To approximate the exponential function e^x over a given interval $[a; b]$, we can directly approximate it by minimizing the least squares error. However, for larger intervals, direct polynomial approximation becomes less effective due to the rapid growth of the exponential function. To address this issue, we approximate the scaled exponential function e^{x^*} over the scaled interval, and the result is then raised to the power of α .

After computing the exponential function, we aggregate the results and compute their inverse using Goldschmidt's division algorithm [19]. Specifically, Goldschmidt's algorithm uses an iterative process of repeatedly refining the inverse approximation using the update formula $f(x) = x(2 - \alpha x)$. The Adaptive Successive Over-Relaxation (aSOR) method [20] enhances this process by dynamically adjusting relaxation factors k_i in each iteration. These factors are applied as $x_{i+1} = k_i x_i (2 - k_i x_i)$, where $k_i = 2/(1 + \epsilon_i)$ is chosen to maximize the rate of convergence by efficiently reducing the error range. Here, ϵ_i represents the minimum value in the current iteration. We consistently apply this approach to the inverse operations throughout the paper.

5.1.2. The Two-phase Method. A large input range can introduce another significant computational challenges beyond stability issues, as it requires a substantial number of fractional bits. This makes it computationally intensive to maintain a necessary high precision within ciphertext. To address this issue and achieve accurate softmax results efficiently, we propose a new *two-phase method* which consists of the following steps.

We begin by shifting the input vector to enhance numerical stability during exponentiation. Given that the elements of the input vector are within the range $[a; b]$, each element is adjusted by subtracting the constant b . As a result, the shifted elements are within the interval $[-M; 0]$, where $M = b - a$. This adjustment prevents the exponential of large positive numbers, making it more manageable to approximate numerically. We then introduce three scaling factors: α , β , and γ . The factor α is selected to define the range for the

Algorithm 3 ASOFTMAX(\mathbf{x})

Input: $\mathbf{x} = (x_i) \in [-M; 0]^n$, $\alpha \in \mathbb{Z}^+$
Output: An approximated softmax of \mathbf{x}
1: $z_i = (\text{AExp}(x_i - \alpha))^\alpha$; $i = 0::n - 1$
2: $S^{(0)} = \sum_{i \in [n]} z_i$
3: $y_i^{(0)} = z_i / \text{Alnv}(S^{(0)})$; $i = 0::n - 1$
4: **for** $j = 1$ to $\log_2(\gamma)$ **do**
5: $z_i = (y_i^{(j-1)})^2$; $i = 0::n - 1$
6: $S^{(j)} = \sum_{i \in [n]} z_i$
7: $y_i^{(j)} = z_i / \text{Alnv}(S^{(j)})$; $i = 0::n - 1$
8: **return** $\mathbf{y} = (y_i^{(\log_2(\gamma))})$

exponential approximation, α is used to manage the power of the initial scaled exponentiation, and γ determines the number of repeated squaring-and-normalization operations. In practice, γ is calculated as $\gamma = 2^{\lceil \log_2(M - \alpha) \rceil}$.

Phase 1: The input vector is scaled by a factor of $1 - \alpha$. Then each element of the scaled input vector lies within the interval $[-M(1 - \alpha); 0]$, enabling stable exponential computations in the subsequent steps. Specifically, the exponentials of the scaled inputs are approximated to a small-degree polynomial using the least squares method. Then, we compute the α -th power of these values. Finally, the resulting exponentials are normalized to form the initial probability distribution by dividing each value by the sum of all exponentials.

Phase 2: To enhance numerical stability and accuracy, the probabilities are refined through an iterative process of squaring and normalization for $\log_2 \gamma$ iterations.

Algorithm 3 provides an explicit description of our softmax approximation. Here, AExp and Alnv refer to the approximated exponential and inverse functions. The output $\mathbf{y}^{(0)} = (y_i^{(0)})$ from the first phase is an approximation of $\text{Softmax}(\mathbf{x} - \alpha)$ over $[-M; 0]^n$ (line 3). Since

$$\text{Softmax}(2\mathbf{x})_i = \frac{y_i^2}{\sum_{j \in [n]} y_j^2}$$

for $\mathbf{y} = (y_i) = \text{Softmax}(\mathbf{x})$, the j -th update $\mathbf{y}^{(j)} = (y_i^{(j)})$ in the second phase is approximate to $\text{Softmax}(2^j \mathbf{x} - \alpha)$. After $\log_2 \gamma$ iterations, this process yields an approximate result of $\text{Softmax}(\mathbf{x})$.

5.2. LayerNorm

In LayerNorm, the normalized output is computed as follows: for $i \in [n]$ and $j \in [d]$,

$$\text{LayerNorm}(X)_{i,j} = \frac{X_{i,j} - \mu_j}{\sqrt{\frac{\sigma_j^2}{d} + \epsilon}}$$

where $X_{i,j}$ is the input vector, μ_j is the mean, σ_j^2 is the variance, ϵ is a small constant for numerical stability, and

α_j and β_j are parameters that scale and shift the normalized output, respectively. We integrate the Goldschmidt’s algorithm with the aSOR technique to approximate the inverse square root operation. Goldschmidt’s inverse square root algorithm refines the approximation iteratively using the update formula $f(x) = x(3 - x)^2 = 4$. Similar to the softmax approximation, aSOR is applied to these factors as $X_{i+1} = k_i X_i (3 - k_i X_i)^2 = 4$, where k_i is selected to maximize the rate of convergence by reducing the error range efficiently.

5.3. GELU & Tanh

We employ a two-step approach to approximate the GELU and Tanh functions using a composition of low-degree polynomials. We first derive a small-degree polynomial $f_1(x)$ that approximates the function by minimizing the least squares error. While this polynomial effectively captures the overall behavior of the input function, it may lack precision in certain regions, especially near the boundaries. To enhance the accuracy, the output of $f_1(x)$ is used as the input for the polynomial $f_2(x)$. This second polynomial is optimized through least squares fitting to refine the approximation, correcting any residual errors. The resulting composite function $f_2(f_1(x))$ offers a close approximation of the target function over large ranges, providing an efficient and accurate approximation method. We apply the Paterson-Stockmeyer algorithm to efficiently evaluate the polynomials in terms of both running time and depth consumption. In our implementation, GELU is approximated using degree-31 and degree-27 polynomials, while Tanh is approximated using two degree-15 polynomials.

6. Secure Transformer Inference

6.1. System Setup & Threat Model

THOR is a secure inference protocol for transformer models. Our system involves two parties: the client and the service provider (or server). The service provider holds a fine-tuned transformer model, while the client seeks inference services from this model. We use BERT to demonstrate the efficiency of our system. Given that the original BERT model is trained on a public dataset, its tokenization process is robust and effectively maps sentences into a BERT-optimized space. We freeze the tokenization process during private fine-tuning.

The threat model of this paper is similar to that of previous private transformer inference works of NEXUS and Powerformer. We assume that the server is honest-but-curious. In other words, it follows the defined protocol precisely but may attempt to infer information from legitimately received messages during the interaction. Our security goal is to ensure that the server does not learn any information about the client’s private input.

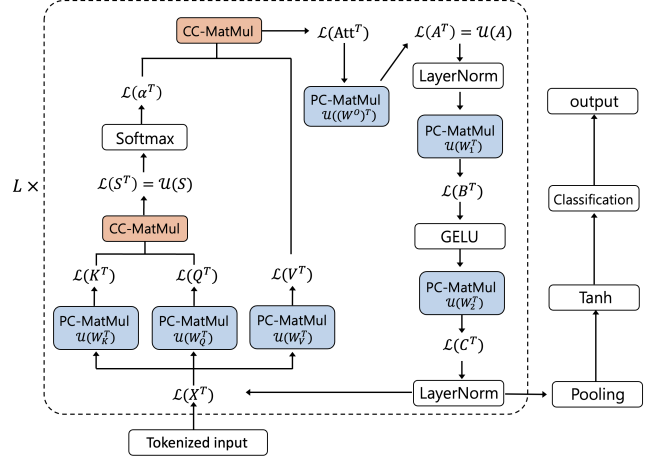


Figure 3. Overview of secure transformer inference in BERT.

6.2. THOR

Figure 3 provides an overview of secure transformer inference within the BERT architecture. In this setup, the client first computes the embeddings for a tokenized sentence using the embedding tables. The resulting output is then encrypted and securely fed into the transformer model for inference.

We consider the computation of $C = AB$, where the matrix B can be provided in the clear. According to the original matrix multiplication algorithm for $C = AB$, the encrypted lower diagonal vectors of C are computed using the rotated ciphertexts of the upper diagonal vectors of A and the plaintexts of the lower diagonal vectors of B . For square matrix multiplication as in Eq. (1), this process requires n^2 rotations to compute $r(U \cdot r(A))$ from the encryptions of $U_k(A)$. Alternatively, the encrypted lower diagonal vectors of the transposed matrix C^T can be obtained using the rotated plaintexts of the upper diagonal vectors of B^T and the ciphertexts of the lower diagonal vectors of A^T . This approach can be implemented without rotations by precomputing $r(U \cdot r(B^T))$ in cleartext. As a result, the transposed approach achieves better computation efficiency than the original method. Throughout the paper, the THOR protocol relies on this encoding format in attention mechanism during transformer inference.

Attention layer: query, key, and value. Given the weight matrices $W_{Q,h} \in \mathbb{R}^{d \times d_k}$ for $h \in [H]$, each query matrix is computed as $Q_h = XW_{Q,h} \in \mathbb{R}^{n \times d_k}$, and these matrices are horizontally stacked to form the query matrix $Q \in \mathbb{R}^{n \times d}$. In other words, we have $Q^T = W_Q^T X^T$, where $W_Q^T \in \mathbb{R}^{d \times d}$ is a matrix formed by vertically stacking $W_{Q,h}^T$ s for $h \in [H]$. This formulation enables secure computation using the diagonal block matrix multiplication method described in Section 4.2.2. This process generates the compactly packed ciphertexts $ct:Q_j$ for $j \in [d_k=c]$, which correspond to the stacked lower diagonals of Q^T . We repeat this procedure to generate ciphertexts $ct:K_j$ and $ct:V_j$ for $j \in [d_k=c]$, which represent the lower diagonals of the transposed key matrix

$K^l = fK_h^l = W_{K,h}^l X^l g_{h2[H]}$ and the transposed value matrix $V^l = fV_h^l = W_{V,h}^l X^l g_{h2[H]}$, respectively.

Attention score. As mentioned in Section 4.2.2, the resulting ciphertext $\text{ct}:K_j$ corresponds to the plaintext vector of $(\tilde{L}_{c_j}(K^l); \tilde{L}_{c_{j+1}}(K^l); \dots; \tilde{L}_{c_{(j+1)-1}}(K^l))$, where each lower diagonal entries of K_h^l are interlaced within $\tilde{L}_k(K^l)$. Using the property in Lemma 3 that $L_{\cdot}(A) = \tilde{L}_{n^+}(\cup_m \cdot(A)) = \tilde{L}_{m^+}(L_m \cdot(A))$ for $A \in \mathbb{R}^{n \times d_k}$, we first securely compute the lower diagonals of K_h^l 's, denoted the output ciphertexts as $\text{ct}:K_j^l$ for $j \in [d_k=c]$ (see Appendix B.1 for details). On the other hand, the ciphertexts of $\text{ct}:Q_j$ are transformed to contain c copies of the lower diagonal vectors of Q_h^l 's, denoted as the output ciphertexts $\text{ct}:Q^l$ for $l \in [d_k]$ (see Appendix B.2 for more detail). Finally, to compute the transposed dot-product $S_h^l = K_h^l Q_h^l \in \mathbb{R}^{n \times n}$ in parallel, we apply the ciphertext-ciphertext matrix multiplication on $\text{ct}:K_j^l$ and $\text{ct}:Q^l$, yielding the ciphertexts $\text{ct}:S_r$ for $r \in [n=c]$. These ciphertexts represent the stacked lower diagonals of S_h^l 's, which are equivalent to its upper diagonals of S_h^l 's.

Soft weights. Let $h = \text{Softmax}(S_h) \in \mathbb{R}^{n \times n}$ be the row-wise application of the softmax operation of the matrix S_h . The row-wise homomorphic evaluation described in Section 4.4 can be extended to compute on compactly packed upper diagonals across multiple input matrices. Specifically, for an approximation ρ of the softmax function, the homomorphically evaluated result of $\rho(\text{ct}:S_0; \text{ct}:S_1; \dots; \text{ct}:S_{n=c-1})$ represents an approximate result of the upper diagonal vectors of h 's.

Attention heads. From the softmax evaluation, we obtain compactly packed ciphertexts of $\text{ct}:h_r$ for $r \in [n=c]$, which represent the lower diagonal vectors of h . Let $\text{Att}_h = h V_h \in \mathbb{R}^{n \times d_k}$ be the h -th context matrix. First, we transform $\text{ct}:h_r$ to contain c copies of the lower diagonal vectors of h 's, resulting in the output ciphertexts $\text{ct}:h^l$ for $l \in [n]$. We then perform ciphertext-ciphertext matrix multiplication on the ciphertexts $\text{ct}:V_j$ and $\text{ct}:h^l$. The resulting ciphertexts, $\text{ct}:\text{Att}_j$ for $j \in [d_k=c]$, correspond to the compactly packed lower diagonals of Att_h .

Multi-head attention. The multi-head attention is computed as $A = \text{Concat}(\text{Att}_0; \dots; \text{Att}_{H-1}) W^O \in \mathbb{R}^{n \times d}$ with $W^O \in \mathbb{R}^{d \times d}$. Here, the projection parameter W^O can be considered as a horizontal concatenation of $W_{O,h} \in \mathbb{R}^{d_k \times d}$ for $h \in [H]$. Let $H_h = \text{Att}_h W_{O,h} \in \mathbb{R}^{n \times d}$. Then the multi-head attention A can be computed by summing H_h over $h \in [H]$.

Let Att be the matrix obtained by vertically stacking the matrices Att_h . We apply diagonal block matrix multiplication on $\text{ct}:\text{Att}_j$ to compute $(W^O)^l \text{Att} \in \mathbb{R}^{d \times n}$, resulting in a matrix formed by vertically concatenating H_h^l 's. Finally, we sum up the output ciphertexts $\text{ct}:H_r$ over $r \in [d=c]$, which correspond to the upper diagonals of the matrix A .

LayerNorm1. We start with computation of the row-wise mean of the matrix A , given the encryptions of its upper diagonals. Similar operations can be applied to compute the mean of squares, and the variance can be obtained by subtracting the square of the means from the mean of squares.

The subsequent computation, involving an approximation of the inverse square root, is performed entrywise over the encrypted data. In the end, we obtain the ciphertexts $\text{ct}:A_r$ for $r \in [d=c]$, which correspond to the upper diagonal vectors of the normalized matrix A .

FC1. Let $B = AW_1$ with $W_1 \in \mathbb{R}^{d \times fd}$. The weight matrix W_1 is partitioned into f submatrices $W_{1;t} \in \mathbb{R}^{d \times d}$. Let $B_t = AW_{1;t} \in \mathbb{R}^{n \times d}$. We apply diagonal block plaintext-ciphertext multiplication on $\text{ct}:A_r$ to compute the compactly packed lower diagonals of B_t^l , denoted as $\text{ct}:B_{t,r}$ for $r \in [d=c]$.

GELU. The approximation of the GELU is performed entrywise over the encrypted data, enabling it to be implemented using the SIMD homomorphic operations on the input ciphertexts, denoted as $\text{ct}:G_{t,r}$.

FC2. For $t \in [f]$, let $C_t = G_t W_{2,t}$, where $G_t = \text{GELU}(B_t)$, $W_{2,t} \in \mathbb{R}^{d \times d}$, and $C = \sum_{t \in [f]} C_t \in \mathbb{R}^{n \times d}$. We apply diagonal block matrix multiplication on $\text{ct}:G_{t,r}$ to compute the compactly packed lower diagonals of C_t^l . Finally, we sum up the ciphertexts $\text{ct}:C_{t,r}$ over $t \in [f]$ to generate the ciphertexts $\text{ct}:C_r$, which correspond to the lower diagonals of C .

LayerNorm2. Similar operations can be applied to obtain the normalized ciphertexts $\text{ct}:X_r$ for $r \in [d=c]$, which represent the lower diagonals of X^l . These are used as input for the next attention layer.

Pooler & classification. Finally, the dense layers are computed on $\text{ct}:X_r$ by applying the plaintext-ciphertext multiplication algorithm.

We remark that our architecture is carefully designed to eliminate the need for repacking or reshaping the output from one encoder layers to match the input format of the subsequent layer. This seamless alignment enables our matrix multiplication algorithms to perform efficiently during secure transformer inference without incurring additional computational costs.

7. Evaluation

7.1. Experimental Setup

BERT Model and Dataset. We follow the BERT-base model with the number of encoding blocks $L = 12$, the number of multi-heads $H = 12$, the embedding dimension $d = 768$, the number of input tokens $t = 128$. We fix the number of self-attention heads to $H = d=64$ (i.e., $d_k = 64$) and the size of the feed-forward features to $4d$ (i.e., $f = 4$). We follow the default fine-tuning hyper-parameters in [21], e.g., batch size 16, learning rate 5×10^{-5} and epoch 5.

We use the MRPC (Microsoft Research Paraphrase Corpus) dataset [22] from the GLUE tasks as our dataset. We performed fine-tuning on the training set and evaluated accuracy on the official validation set.

Libraries and Configurations. THOR is evaluated using the Liberate.FHE library [23], which implements the RNS variant of the CKKS scheme [24]. The ciphertext dimension

TABLE 4. SOFTMAX COMPARISON BETWEEN [26] AND OUR METHOD.

Method	# Iter.	Depth	# Mult	Precision (bits)
[26]	5	40	106	11
Ours	3	30	30	8

was set as $N = 2^{16}$ to support the bootstrapping operation, providing plaintext slots of length $s = 2^{15}$. The ciphertext modulus q was chosen to ensure a 128-bit security level [25]. Under these HE parameters, 13 multiplicative levels are available before bootstrapping is needed. All experiments were conducted on an Intel Platinum 8462 CPU at 2.8GHz and a single NVIDIA A100 GPU.

7.2. Micro-benchmark Evaluation of Softmax

In Table 4, we compare our softmax approximation algorithm with that of Cho et al. [26] in terms of the number of iterations (i.e., the number of approximate normalization operations), multiplicative depth, and the number of homomorphic multiplications required for softmax evaluation excluding initial exponentiation computations. We also measured the worst-case accuracy of the approximations. The input data for this evaluation was extracted from 12,800 samples from the third layer of the MRPC dataset, with $x_i \in [70;70]$ and $n = 128$.

It follows from [26, Theorem 3.4] that the softmax approximation requires $k = \lceil \log_{140} \log \ln 128e \rceil = 5$ iterations of repeated normalization-and-squaring operations. Since their experimental results were based on random inputs, we made slight adjustments to the inverse square root approximations used in the normalization procedure to fit the MRPC dataset. Specifically, the first and last iterations require depths of 7 and 11 for the inverse square roots, while the intermediate iterations require a depth of 4. In contrast, our algorithm achieves similar accuracy with only 3 iterations of squaring and normalization by the inverse root approximation. For this evaluation, the scaling factors are chosen as follows: $\alpha_0 = \frac{140}{8} = 17.5$, $\alpha_1 = 2$, and $\alpha_2 = 4$.

Since each iteration in both algorithms requires a single bootstrapping, the reduced number of iterations in our method lowers the required bootstrapping operations, thereby reducing the overall computational time and resource usage. Additionally, we optimize depth consumption and minimize multiplications for normalization by applying the aSOR-based Goldschmidt inverse algorithm. The experimental results show that our approach achieves substantial efficiency gains without compromising accuracy, making it suitable for applications that requires rapid and resource-efficient computation, particularly in Machine Learning as a Service (MLaaS) models.

7.3. End-to-enc Inference Evaluation

Performance Breakdown. Table 5 shows a detailed breakdown of the execution time for the BERT-base model. Due to the significant cost of bootstrapping, THOR’s architecture

TABLE 5. BREAKDOWN OF THE EXECUTION TIME OF THOR.

Operation	Input	Time (sec)
Attention layer	3 ($R^{128 \times 768}$, $R^{768 \times 64}$)	49.77
Attention score	12 ($R^{128 \times 64}$, $R^{64 \times 128}$)	32.53
Softmax	12 ($R^{128 \times 128}$)	15.53
Attention heads	12 ($R^{128 \times 128}$, $R^{128 \times 64}$)	20.63
Multi-head attention	$R^{128 \times 768}$, $R^{768 \times 768}$	27.43
LayerNorm1	$R^{128 \times 768}$	7.13
FC1	$R^{128 \times 768}$, $R^{768 \times 3072}$	49.80
GELU	$R^{128 \times 3072}$	29.42
FC2	$R^{128 \times 3072}$, $R^{3072 \times 768}$	49.19
LN2	$R^{128 \times 768}$	4.10
Pooler & Classification	$R^{128 \times 768}$	2.70
Bootstrappings	-	337.86
Total	-	104.35

TABLE 6. PERFORMANCE FOR THE BASELINE AND OUR METHOD.

Metrics	Baseline	Unencrypted			Encrypted
		G	G-LN	G-LN-S	
Accuracy	85.29	85.54	85.54	85.78	84.80
F1-score	89.90	90.05	90.05	90.24	89.49

is designed to perform bootstrapping operations when the number of intermediate or output ciphertexts is minimal. Nevertheless, it turns out that more than 53.96% of the total runtime is attributed to bootstrapping operations. Since the complexity of basic homomorphic operations, such as additions, multiplication, and rotations, scales almost linearly with the ciphertext level, we perform homomorphic matrix multiplications at the lowest possible level. The linear layers in the attention mechanism and the fully connected layers account for 36.70% of the total runtime. Our optimized matrix multiplication algorithms provide performance gains, resulting in a relatively lower runtime. While the total inference time of THOR is 10.43 minutes in a single GPU setting, NEXUS takes approximately 37.3 sec 128 inputs 78.93 minutes with four A100 GPUs to process a single prediction request on the BERT-base model.

Accuracy. In Table 6, we present the performance of the BERT-base model across five different settings to evaluate the effectiveness of our approximation components: (a) Baseline: fine-tuning the pre-trained model with no approximation, (b) G: fine-tuning the model with GELU replaced by our GELU approximation, (c) G-LN: additionally fine-tuning the model with layer normalization operation replaced by our layer normalization approximation, (d) G-LN-S: implementing full approximation, including the softmax replaced by our approximation model, and (e) Encrypted: applying the full approximation to encrypted data for secure inference. As in previous works, we report the accuracy and the F1 score. Overall, the accuracy achieved by THOR matches that of the plaintext model, with an accuracy drop of 0.49% compared with the baseline. We attribute this not to approximations of the non-linear functions but rather to computational error from homomorphic operations.

8. Related Work

8.1. Privacy-preserving Transformer Inference

Several studies have investigated how to preserve data privacy during transformer inference. In the MPC-based approach, MPCFormer [8] introduces a private inference system for the BERT-base model using secret-sharing techniques. More recently, PrivFormer [9] and PUMA [10] have proposed secure inference for transformers in a 3PC setting. In the hybrid approach using HE, Iron [5] and BOLT [6] present HE-based matrix multiplication algorithms that utilize ciphertext packing techniques. Despite these advancements in the protocol, they still incur a high communication cost; for example, BOLT requires 25.74 GB of communication for an end-to-end inference on the BERT-base model under one of the WAN settings.

In the HE-based approach, THE-X [11] proposed crypto-friendly approximation methods for non-linear functions in transformer inference. Specifically, THE-X replaces the softmax function with a combination of ReLU and a three-layer linear neural network, substitutes GELUs with ReLUs, and replaces LayerNorm with an affine transformation with learnable parameters. However, their scenario is different from ours: intermediate results of non-linear layers (input of ReLU) are sent to the client, who then performs computation on the decrypted results and re-encrypts them before sending them back to the server. This approach introduces two significant issues: (i) the frequent back-and-forth data exchange between the client and server increases communication costs significantly and (ii) the inputs to RELU are exposed to the client during this process, leading to potential privacy leakages. Recent works of [12, 13] have focused on the BERT-tiny model. Thus, cryptographic techniques including homomorphic matrix multiplication algorithms are not scalable to larger transformer models with high embedding dimensions and a large number of layers.

The intersection of HE and BERT was investigated in [27], where BERT embeddings are used for training a HE-based logistic regression model. However, since the primary focus is on implementing the logistic regression model, the BERT inference is performed by the client rather than the server, and only the text classification is conducted on encrypted data.

8.2. Homomorphic Matrix Multiplication

Various approaches have been investigated to optimize plaintext-ciphertext matrix multiplication (PC-MM). Gazelle [28] and Iron [5] rely on inner dot products of input matrices, resulting in a non-compact output packing on plaintext slots or coefficients after computation. BOLT [6], the current state-of-the-art for PCMM, still suffers from computational complexity in large transformer models. Furthermore, it aims to minimize the multiplicative depth for ciphertext-ciphertext matrix multiplication due to limitations in the HE parameters when integrated with MPC. As a

result, the computation is suboptimal, and its computational complexity is significantly higher than ours.

The algorithms proposed in [12, 29] perform inner dot products over packed ciphertexts, requiring numerous slot permutations to accumulate values located in different plaintext slots. This approach relies on repetitive rotate-and-sum operations, leading to wasted slots in the resulting ciphertext after homomorphic matrix multiplication. Consequently, this increases computational overhead due to the need for repacking the non-compact output into multiple ciphertexts.

NEXUS [14] employs various packing methods similar to BOLT [6], including component-wise, row-wise, column-wise, and diagonal-wise packings. However, the details of changes in the packing structure during secure inference are not fully explained. On the other hand, PowerFormer [13] leverages the state-of-the-art ciphertext-ciphertext matrix multiplication algorithms of Jiang et al. [30]. Despite these algorithmic improvements, the matrix multiplication algorithms remain suboptimal for rectangular matrix multiplication, as they result in redundant entries after computation and these sparsely packed ciphertexts are then used as input for the subsequent matrix computations.

The recent work of Bae et al. [31] proposed the MaMBo (Matrix Multiplication Bootstrapping) framework, which performs PC-MM by leveraging multi-secret variants of RLWE to compactly represent multiple ciphertexts. However, this approach does not address the challenge of performing multiplication between ciphertexts, which is a necessary operation for tasks such as transformer inference. In contrast, THOR introduces a unified concept that not only retains the efficiency of multiplication between plaintext and ciphertext, but also supports multiplication between ciphertexts in an efficient manner, along with other complex operations required for implementing large language models in homomorphic encryption.

8.3. Softmax Computation

In NEXUS [14], the softmax is approximated using the Taylor series expansion of the exponential function and Goldschmidt’s division algorithm. To improve numerical stability of the exponential function, the maximum value of each row in the input matrix is subtracted from each element. In their implementation, this maximum value is taken as a constant for efficiency, which is impractical in real-world applications. As explained in [32], this value should be computed through a substantial number of comparison operations on encrypted data, leading to significant computational overhead.

An alternative method for computing softmax on encrypted data is to replace it with a more HE-friendly function, avoiding unstable approximation methods. For instance, PowerFormer [13] replaced the softmax function with a ReLU-based function, while [33] substituted it with Gaussian kernels to enable secure computation without division. However, these methods achieve only comparable performance to the original softmax in small BERT models (e.g., two-layer BERT model) on encrypted data.

9. Conclusion

In this study, we present THOR, a fast and secure transformer inference system on encrypted data. THOR leverages efficient homomorphic matrix multiplication algorithms, achieving a significant speedup in attention mechanisms compared to previous approaches. With recent advancements in algorithmic efficiency of HE and implementation optimization through hardware accelerators, we believe that performance can be significantly accelerated and parallelized, further enhancing THOR's performance. The cryptographic techniques in THOR are applicable to the other transformer encoders that utilize different non-linear approximations from ours, such as [13, 33], as well as transformer decoders such as GPT [34]. Additionally, our fundamental developments can be extended to other privacy-preserving machine learning applications, including neural network models and transfer learning.

References

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [2] J. Devlin, M.-W. Chang, K. Lee, and K. 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: Human Language Technologies*, 2019, pp. 4171–4186.
- [3] Openai, "chatgpt," <https://openai.com/blog/chatgpt>, 2022.
- [4] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "LLaMA: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [5] M. Hao, H. Li, H. Chen, P. Xing, G. Xu, and T. Zhang, "Iron: Private inference on transformers," *Advances in neural information processing systems*, vol. 35, pp. 15 718–15 731, 2022.
- [6] Q. Pang, J. Zhu, H. Möllering, W. Zheng, and T. Schneider, "BOLT: Privacy-preserving, accurate and efficient inference for transformers," in *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2024, pp. 4753–4771.
- [7] W.-j. Lu, Z. Huang, Z. Gu, J. Li, J. Liu, C. Hong, K. Ren, T. Wei, and W. Chen, "BumbleBee: Secure two-party inference framework for large transformers," *Manuscript, to appear in Network and Distributed System Security Symposium 2025*, p. available at <https://eprint.iacr.org/2023/1678>.
- [8] D. Li, H. Wang, R. Shao, H. Guo, E. Xing, and H. Zhang, "MPCFORMER: Fast, performant and private transformer inference with MPC," in *The Eleventh International Conference on Learning Representations*, 2023.
- [9] Y. Akimoto, K. Fukuchi, Y. Akimoto, and J. Sakuma, "Privformer: Privacy-preserving transformer with MPC," in *2023 IEEE 8th European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2023, pp. 392–410.
- [10] Y. Dong, W.-j. Lu, Y. Zheng, H. Wu, D. Zhao, J. Tan, Z. Huang, C. Hong, T. Wei, and W. Chen, "PUMA: Secure inference of LLaMA-7B in five minutes," *arXiv preprint arXiv:2307.12533*, 2023.
- [11] T. Chen, H. Bao, S. Huang, L. Dong, B. Jiao, D. Jiang, H. Zhou, J. Li, and F. Wei, "THE-X: Privacy-preserving transformer inference with homomorphic encryption," in *Findings of the Association for Computational Linguistics: ACL 2022*, 2022, pp. 3510–3520.
- [12] L. Roviada and A. Leporati, "Transformer-based language models and homomorphic encryption: An intersection with bert-tiny," in *Proceedings of the 10th ACM International Workshop on Security and Privacy Analytics*, 2024, pp. 3–13.
- [13] D. Park, E. Lee, and J.-W. Lee, "Powerformer: Efficient privacy-preserving transformer with batch rectifier-power max function and optimized homomorphic attention," *Cryptology ePrint Archive, Paper 2024/1429*, 2024.
- [14] J. Zhang, J. Liu, X. Yang, Y. Wang, K. Chen, X. Hou, K. Ren, and X. Yang, "Secure transformer inference made non-interactive," *Manuscript, to appear in Network and Distributed System Security Symposium 2025*, p. available at <https://eprint.iacr.org/2024/136>, 2025.
- [15] C. Gentry, "Fully homomorphic encryption using ideal lattices," in *Proceedings of the forty-first annual ACM symposium on Theory of computing*, 2009, pp. 169–178.
- [16] J. H. Cheon, A. Kim, M. Kim, and Y. Song, "Homomorphic encryption for arithmetic of approximate numbers," in *International Conference on the Theory and Application of Cryptology and Information Security*. Springer, 2017, pp. 409–437.
- [17] J. H. Cheon, K. Han, A. Kim, M. Kim, and Y. Song, "Bootstrapping for approximate homomorphic encryption," in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 2018, pp. 360–384.
- [18] M. Kim, D. Lee, J. Seo, and Y. Song, "Accelerating HE operations from key decomposition technique," in *Annual International Cryptology Conference*. Springer, 2023, pp. 70–92.
- [19] R. E. Goldschmidt, "Applications of division by convergence," Ph.D. dissertation, Massachusetts Institute of Technology, 1964.
- [20] J. Moon, Z. Omarov, D. Yoo, Y. An, and H. Chung, "Adaptive successive over-relaxation method for a faster iterative approximation of homomorphic operations," *Cryptology ePrint Archive, Paper 2024/1366*, 2024.
- [21] Bert, <https://github.com/google-research/bert>, 2019, google.

- [22] B. Dolan and C. Brockett, "Automatically constructing a corpus of sentential paraphrases," in *Third international workshop on paraphrasing (IWP2005)*, 2005.
- [23] DESILO, "Liberate.FHE: A new FHE library for bridging the gap between theory and practice with a focus on performance and accuracy," <https://github.com/Desilo/liberate-fhe>, 2023.
- [24] J. H. Cheon, K. Han, A. Kim, M. Kim, and Y. Song, "A full RNS variant of approximate homomorphic encryption," in *International Conference on Selected Areas in Cryptography*. Springer, 2018, pp. 347–368.
- [25] J.-P. Bossuat, R. Cammarota, J. H. Cheon, I. Chillotti, B. R. Curtis, W. Dai, H. Gong, E. Hales, D. Kim, B. Kumara *et al.*, "Security guidelines for implementing homomorphic encryption," *Cryptology ePrint Archive*, 2024.
- [26] W. Cho, G. Hanrot, T. Kim, M. Park, and D. Stehlé, "Fast and accurate homomorphic softmax evaluation," in *Proceedings of the 2024 ACM SIGSAC conference on computer and communications security*, 2024.
- [27] G. Lee, M. Kim, J. H. Park, S.-w. Hwang, and J. H. 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*, 2022, pp. 3169–3175.
- [28] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan, "GAZELLE: A low latency framework for secure neural network inference," in *27th USENIX security symposium (USENIX security 18)*, 2018, pp. 1651–1669.
- [29] E. Crockett, "A low-depth homomorphic circuit for logistic regression model training," *Cryptology ePrint Archive*, 2020.
- [30] X. Jiang, M. Kim, K. Lauter, and Y. Song, "Secure outsourced matrix computation and application to neural networks," in *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*, 2018, pp. 1209–1222.
- [31] Y. Bae, J. H. Cheon, G. Hanrot, J. H. Park, and D. Stehlé, "Plaintext-ciphertext matrix multiplication and FHE bootstrapping: Fast and fused," in *Annual International Cryptology Conference*. Springer, 2024, pp. 387–421.
- [32] S. Lee, G. Lee, J. W. Kim, J. Shin, and M.-K. Lee, "HETAL: Efficient privacy-preserving transfer learning with homomorphic encryption," in *International Conference on Machine Learning*. PMLR, 2023, pp. 19 010–19 035.
- [33] D. Rho, T. Kim, M. Park, J. W. Kim, H. Chae, J. H. Cheon, and E. K. Ryu, "Encryption-friendly LLM architecture," *arXiv preprint arXiv:2410.02486*, 2024.
- [34] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.

Appendix A.

Proofs of Main Results

A.1. Proof of Proposition 1

Proof. Let $C = AB \in \mathbb{R}^{n \times n}$. For $\ell \in [m]$ and $r, t \in [n]$, the elements $A[r + t; \ell + t]$ and $B[\ell + t; t]$ correspond to the t -th entries of the r -th shifted vector $r(U_{\cdot-r}(A))$ and the lower diagonal vector $L_{\cdot}(B)$, respectively. Therefore, the t -th entry of the vector $\sum_{\ell \in [m]} r(U_{\cdot-r}(A)) \cdot L_{\cdot}(B)$ is as follows:

$$\begin{aligned} \left(\sum_{\ell \in [m]} r(U_{\cdot-r}(A)) \cdot L_{\cdot}(B) \right)[t] &= \sum_{\ell \in [m]} A_{r+t; \ell+t} \cdot B_{\ell+t; t} \\ &= C_{r+t; t} \\ &= L_r(AB)[t]. \end{aligned}$$

□

A.2. Proof of Proposition 2

Proof. Let $C = AB \in \mathbb{R}^{m \times n}$. We first consider the t -th entry of the r -th shifted vector $r(U_{\cdot-r}^{\ell}(A))$ for $r \in [m]$ and $\ell \in [n]$. If $\ell - r < m$, this entry corresponds to the t -th entry of $r(U_{\cdot-r}(A))$, which is equal to the element $A[r + t \pmod{m}; \ell + t \pmod{n}]$. Otherwise, suppose that $\ell - r = mQ + R$ for some integers $Q \in \mathbb{Z}$ and $R \in [m]$. In this case, the entry corresponds to the t -th entry of

$$\begin{aligned} r(U_{\cdot-r}^{\ell}(A)) &= r(m b(\ell - r) = mc)(U_{(\ell - r) \pmod{m}}(A)) \\ &= r(mQ)(U_R(A)) \\ &= r+mQ(U_R(A)). \end{aligned}$$

We note that

$$\begin{aligned} \rho^{r+mQ}(U_R(A)) &= \rho^{r+mQ}(A_{0;R}, A_{0;R+1}, \dots) \\ &= (A_{r;r+mQ+R}, A_{r+1;r+mQ+R+1}, \dots) \\ &= (A_{r;}, A_{r+1;+1}, \dots). \end{aligned}$$

Thus, the t -th entry of $r+mQ(U_R(A))$ is $A[r + t \pmod{m}; \ell + t \pmod{n}]$. Similarly, the element $B[\ell + t \pmod{n}; t \pmod{n}]$ corresponds to the t -th entry of the lower diagonal vector $L_{\cdot}(B)$. Therefore, for each $r \in [m]$, the t -th entry of $\sum_{\ell \in [n]} r(U_{\cdot-r}^{\ell}(A)) \cdot L_{\cdot}(B)$ is as follows:

$$\begin{aligned} \left(\sum_{\ell \in [n]} r(U_{\cdot-r}^{\ell}(A)) \cdot L_{\cdot}(B) \right)[t] &= \sum_{\ell \in [n]} A_{r+t; \ell+t} \cdot B_{\ell+t; t} \\ &= C_{r+t; t} \\ &= L_r(AB)[t]. \end{aligned}$$

□

A.3. Proof of Corollary 1

Proof. By definition 3, we have

$$L_n^{\ell} \cdot k = {}^{m(n-m-1)}(L_m \cdot k) = {}^{n-m}(L_m \cdot k)$$

for $k \in [m]$. It follows from Lemma 3 that

$$\begin{aligned} U_k(A) &= U^{m+k}(L_{m-k}(A)) \\ &= U^{m+k}(U^{n+m}(L_n^0(L_{m-k}(A)))) \\ &= U^k(L_n^0(L_{m-k}(A))) \end{aligned}$$

for $k \in [m]$. Therefore, we have

$$r(U_{\cdot,r}) = r(U_{\cdot,r}(L_n^0(L_{\cdot+r}(A)))) = U_{\cdot+r}^0(L_{\cdot+r}(A));$$

as desired. \square

A.4. Proof for Corollary 2

Proof. The upper diagonal vector $U_k^0(A) \in \mathbb{R}^n$ is equal to $U_k(A)$ when $\cdot + r < m$. By Lemma 3, we have $r(U_{\cdot+r}^0(A)) = r(U_{\cdot+r}(A)) = r(U_{\cdot+r}(L_{m-\cdot+r}(A))) = U_{[m-\cdot+r]_m}^0(A)$. Now, suppose that $\cdot + r \geq m$. Then we have

$$\begin{aligned} r(U_{\cdot+r}^0(A)) &= r(U_{[m-\cdot+r]_m}^0(A)) \\ &= r(U_{[m-\cdot+r]_m}(L_{m-[m-\cdot+r]_m}(A))) \\ &= r(U_{\cdot+r}(L_{m-[m-\cdot+r]_m}(A))) \\ &= U_{[m-\cdot+r]_m}^0(A); \end{aligned}$$

\square

Appendix B. Homomorphic Matrix Operations

B.1. Matrix Transposition

Algorithm 4 performs matrix transposition which transforms the upper-diagonally stacked vectors of a matrix B into the upper-diagonally stacked vectors of B^\top . Equivalently, the output ciphertexts can be also considered as the lower-diagonally stacked vectors of B . We define two masking vectors $\kappa_{j;k,0}, \kappa_{j;k,1} \in \mathbb{R}^S$ which have entries as follows:

$$\begin{aligned} \kappa_{j;k,0}[i] &= \begin{cases} 1 & \text{if } a \leq i < a+n-b; \\ 0 & \text{otherwise;} \end{cases} \\ \kappa_{j;k,1}[i] &= \begin{cases} 1 & \text{if } [a-b]_S \leq i < [a-b]_S+b; \\ 0 & \text{otherwise;} \end{cases} \end{aligned}$$

where $a = [kn]_S$ and $b = [n-k-c]_n$. This transposition procedure requires $(2n-1)$ scalar multiplications and $(n-1+n-c)$ rotations in total. By leveraging both real and complex numbers for rotation, the total computational cost changes to $4n$ scalar multiplications and $\frac{(c+2)n}{2c}$ rotations.

Algorithm 4 TRANSPOSE

Input: Ciphertexts $\text{ct}_{\cdot,j}$ of the upper-diagonally stacked vectors of $B \in \mathbb{R}^{n \times n}$ for $j \in [n=c]$

Output: Ciphertexts $\text{ct}_{\cdot,j}^\top$ of the upper-diagonally stacked vectors of B^\top for $\cdot \in [n=c]$

```

1: for  $j = 0$  to  $n=c-1$  do
2:    $\text{ct}_{\cdot,j}^\top \leftarrow [n-c]_n(\text{ct}_{\cdot,j})$ 
3:    $l \leftarrow j \pmod{n=c}$ 
4:    $\text{ct}_{\cdot,0} \leftarrow \text{SMult}(\text{ct}_{\cdot,0}^\top; j;k,0)$ 
5:    $\text{ct}_{\cdot,1} \leftarrow \text{SMult}(\text{ct}_{\cdot,1}^\top; j;k,1)$ 
6: for  $j = 0$  to  $n=c-1$  do
7:    $\cdot \leftarrow n=c-1-j$ 
8:   for  $k = 1$  to  $c$  do
9:      $\text{ct}_{\cdot,0}^\top \leftarrow [n-c]_{k+2kn-cn}(\text{ct}_{\cdot,j})$ 
10:     $\text{ct}_{\cdot,0} \leftarrow \text{Add}(\text{ct}_{\cdot,0}; \text{SMult}(\text{ct}_{\cdot,0}^\top; j;k,0))$ 
11:     $\text{ct}_{\cdot,1} \leftarrow \text{Add}(\text{ct}_{\cdot,1}; \text{SMult}(\text{ct}_{\cdot,1}^\top; j;k,1))$ 
12: for  $\cdot = 0$  to  $n=c-1$  do
13:    $\text{ct}_{\cdot,j}^\top \leftarrow \text{Add}(\text{ct}_{\cdot,0}; [n](\text{ct}_{\cdot,1}))$ 
14: return  $\text{ct}_{\cdot,j}^\top$ 

```

Algorithm 5 MERGE-COPY

Input: Ciphertexts $\text{ct}_{\cdot,j}$ of the lower-diagonally stacked vectors of $B \in \mathbb{R}^{n \times n}$ for $j \in [n=c]$

Output: Ciphertexts $\text{ct}_{\cdot,j}^\top$ of the replicated lower diagonal vectors of B for $\cdot \in [n]$

```

1: for  $j = 0$  to  $n=c-1$  do
2:   for  $k = 0$  to  $c-2-1$  do
3:      $\text{ct}_{j;k} \leftarrow \text{SMult}(\text{ct}_{\cdot,j}; k)$ 
4:     for  $r = 0$  to  $\log_2(c-2)-1$  do
5:        $\text{ct}_{j;k} \leftarrow \text{ct}_{j;k} + 2^{n-2^r}(\text{ct}_{j;k})$ 
6:      $\text{ct}_{\cdot,j+2k}^\top \leftarrow \text{SMult}(\text{ct}_{j;k}; l)$ 
7:      $\text{ct}_{\cdot,j+2k+1}^\top \leftarrow \text{ct}_{j;k} \oplus \text{ct}_{\cdot,j+2k}$ 
8:   for  $k = 0$  to  $c-1$  do
9:      $\text{ct}_{\cdot,j+k}^\top \leftarrow \text{ct}_{\cdot,j+k}^\top + [n](\text{ct}_{\cdot,j+k})$ 
10: return  $\text{ct}_{\cdot,j}^\top$ 

```

B.2. Merge-and-copy

Algorithm 5 provides the detailed procedure for generating replicated lower diagonal vectors from the stacked lower diagonal vectors of a matrix B . We define two masking vectors $\kappa_{\cdot}, \iota_{\cdot} \in \mathbb{R}^S$ which have entries as follows:

$$\begin{aligned} \kappa[l] &= \begin{cases} 1 & \text{if } 2kn \leq i < 2(k+1)n; \\ 0 & \text{otherwise;} \end{cases} \\ \iota[l] &= \begin{cases} 1 & \text{if } 0 \leq i < n; \\ 0 & \text{otherwise;} \end{cases} \end{aligned}$$

This merge-and-copy procedure requires n scalar multiplication and $\frac{n}{2}(\log c + 1)$ rotations in total.