# Efficient Homomorphic Matrix Computation for Secure Transformer Inference

**Miran Kim**

Hanyang University

Joint work with Jungho Moon, Dongwoo Yoo, Xiaoqian Jiang

# Motivation

- Privacy issue with AI
  - **Personalized services** require personal information
  - Cloud computing services: Google Gemini, Meta LLaMA, OpenAI ChatGPT
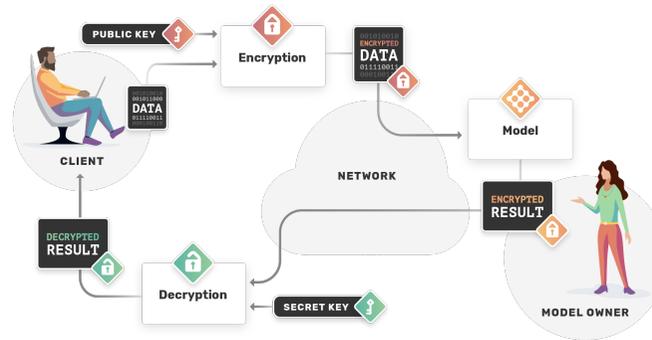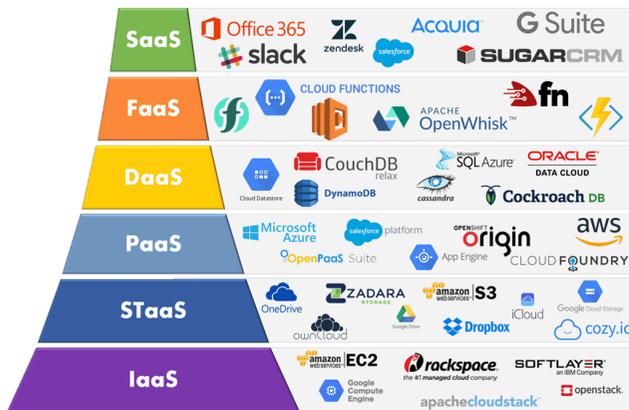  - **Data abuse**: if data is available, it will be used





Italy's privacy watchdog fines OpenAI for ChatGPT's violations in collecting users personal data

# Motivation

**Question**. Can ***Prediction as a Service*** be made trustworthy and efficient?

(Privacy-preserving Personalized Service)

- Traditional encryption protects only storage and transmission—not the computation.

- HE enables computations to be performed without decrypting the data.

# Transformer-based Model

- What is *Transformer*?
  - o  **Self-attention** based architecture [V+18] (how strongly each token is related to every other token)
  - o  Foundation of modern NLP models like BERT, GPT, T5, BART
  - o  *Parallelizable* and efficient processing of sequences



[V+18] Attention is all you need, NeurIPS 2018

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d}}\right)V$$

4

# Table of Contents

# What is Homomorphic Encryption (HE)?

- HE can evaluate arithmetic functions on encrypted data.
    - For $x, y \in \mathcal{M} = \mathbb{Z}_q[X]/(X^N + 1)$,

$$Dec\big(Enc(x) + Enc(y)\big) = Dec\big(Enc(x)\big) + Dec\big(Enc(y)\big) \qquad = x + y$$
$$Dec\big(Enc(x) * Enc(y)\big) = Dec\big(Enc(x)\big) * Dec\big(Enc(y)\big) \qquad = x * y$$

- HE offers lots of flexibility in parameter selection
    - Plaintext modulus, vector-size, polynomial degree …

- ***Ciphertext packing technique*** to improve performance
    - Each ciphertext encrypts multiple plaintext elements and perform ciphertext computation in a *SIMD* manner.
    - Even with a fixed parameter set, determining how to put each data element into the ciphertexts requires careful design to achieve optimal efficiency.

# Ciphertext Packing Technique

- Slot-encoding
  - Use *DFT-like algorithm* to transform a plaintext vector into an element of a cyclotomic ring (BGV,BFV,CKKS)
  - Support *vectorized* operations in a *SIMD* manner
    - Entry-wise addition/multiplication
    - Rotation: moving values between slots
  - Computing the same function on $l$ inputs at the price of one computation (better amortized size and timing)

- Coefficient-encoding
  - Encode plaintext values *directly* as an element of a cyclotomic ring (FHEW, TFHE)

# Homomorphic Matrix Computation

*Question.* How to efficiently perform matrix computation
over encrypted data?

### Inside FHE

- **Slot**-encoding & HE-friendly expression of MM
  - Matrix $\Rightarrow$ $1d$-vector
  - Matrix operations $\Rightarrow$ composite of vectorized operations
- Related Work
  - Row-major: JKLS'18
  - Column-major: BOLT (PC-MM)
  - Diagonal-major: HS'14 (Mat×vec)

### Outside FHE

- **Coefficient**-encoding & reduction from encrypted MM to plaintext MMs
  - Multiply matrix directly on the decryption equation
  - Enable regular matrix computation optimizations
  - Related work: [B+24, P25]

# Inside FHE

- <span style="color:red">Row-major</span> matrix encoding: identify $(d \times d)$ matrix to $d^2$-dimensional vector $(d^2 \leq s)$

- **(Vectorized) matrix computation** [JKLS'18]
  - $\mathrm{vec}(A + B) = \mathrm{vec}(A) \oplus \mathrm{vec}(B)$
  - $\mathrm{vec}(AB) = \sum \mathrm{vec}(A_i) \odot \mathrm{vec}(B_i)$ where $A_i, B_i$ are permuted matrices of $A$ and $B$
    - $\mathrm{vec}(A), \mathrm{vec}(B) \xRightarrow{\text{Linear Transf}} \mathrm{vec}(A_0), \mathrm{vec}(B_0)$
    - $A_i = \texttt{Column-shift}(A_0, i)$; $B_i = \texttt{Row-shift}(B_0, i)$

$A$

| $a_{00}$ | $a_{01}$ | $a_{02}$ |
|---|---|---|
| $a_{10}$ | $a_{11}$ | $a_{12}$ |
| $a_{20}$ | $a_{21}$ | $a_{22}$ |

$\mathrm{vec}(\boldsymbol{A})$    $\| \| \|$

| $a_{00}$ | $a_{01}$ | $a_{02}$ | $a_{10}$ | $a_{11}$ | $a_{12}$ | $a_{20}$ | $a_{21}$ | $a_{22}$ |
|---|---|---|---|---|---|---|---|---|

- **Limitation**
  - Small matrix size assumption: $d^2 \leq s$
    - But need large-scale MM (e.g., (768x768)x(768x128))
    - Block-wise computation [C+20]
  - Originally proposed for CC-MM; Not optimal for plaintext-ciphertext multiplication (PC-MM)
    - $A_i$ or $B_i$ can be provided in cleartext (e.g., feed-forward, linear projection)
  - Matrix transposition: $\mathcal{O}(d)$ complexity

[JKLS18] X. Jiang, **M. Kim**, K. Lauter, Y. Song. "Secure outsourced matrix computation and application to neural networks", CCS 2018
[C+20] H. Chen, **M. Kim,** I. Razenshteyn, D. Rotaru, Y. Song, S. Wagh. "Maliciously secure MM with applications to Private DL". Asiacrypt 2020.

# Outside FHE

- **Reduction from encrypted MM to plaintext MMs** [B+24, P25]
  - A RLWE ciphertext of $m(X) = \sum m_j X^j \Leftrightarrow$ plaintext vector of $(m_j)$
  - RLWE ciphertexts of $m_i(X) = \sum m_{ij} X^j \Leftrightarrow$ plaintext matrix of $(m_{ij})$
  - $AS + B = M \ (mod\ q) \Rightarrow (WA)S + (WB) \approx WM (mod\ q)$

- **Limitation**
  - Optimal when $\#(\text{matrices}) = N$
  - Need a conversion to slot-encoded ciphertext when performing coefficient-wise operations (e.g., nonlinear function evaluation)

$$
\begin{bmatrix} \text{---} & a_1^t & \text{---} \\ \text{---} & a_2^t & \text{---} \\ & A \vdots & \\ \text{---} & a_N^t & \text{---} \end{bmatrix}
\begin{bmatrix} s_0 & s_1 & \cdots & s_{N-1} \\ -s_{N-1} & s_0 & \cdots & s_{N-2} \\ \vdots & \vdots & \ddots & \vdots \\ -s_1 & -s_2 & \cdots & s_0 \end{bmatrix}
+
\begin{bmatrix} \text{---} & b_1^t & \text{---} \\ \text{---} & b_2^t & \text{---} \\ & B \vdots & \\ \text{---} & b_N^t & \text{---} \end{bmatrix}
=
\begin{bmatrix} \text{---} & m_1^t & \text{---} \\ \text{---} & m_2^t & \text{---} \\ & M \vdots & \\ \text{---} & m_N^t & \text{---} \end{bmatrix}
$$

[B+24] Y. Bae, J.H. Cheon, G. Hanrot, J.H. Park, D. Stehle. "Plaintext-ciphertext MM and FHE BTS: fast and fused", CRYPTO 2024
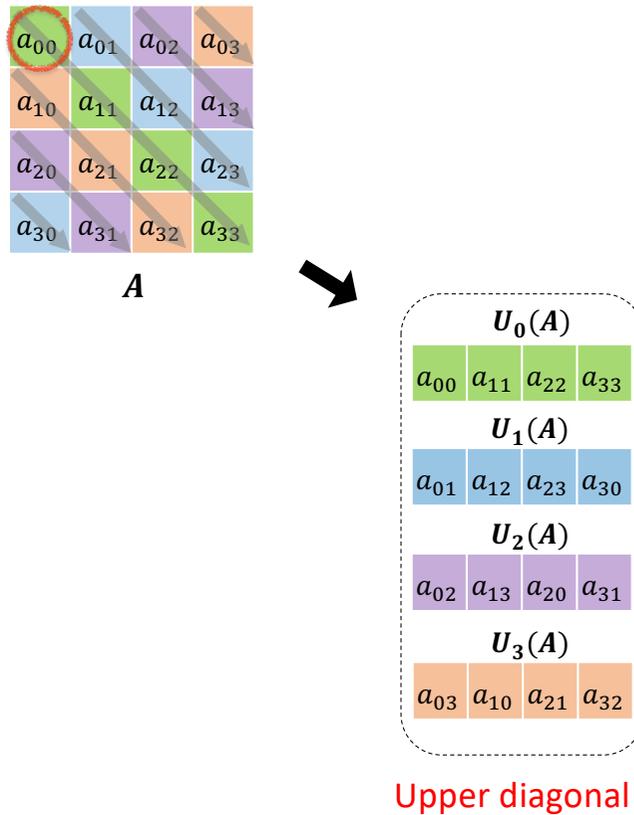[P25] J.H. Park. "Ciphertext-ciphertext MM: fast for large matrices", Eurocrypt 2025

# Table of Contents

# Diagonal-major Matrix Encoding

$$
A = \begin{bmatrix} a_{00} & a_{01} & a_{02} & a_{03} \\ a_{10} & a_{11} & a_{12} & a_{13} \\ a_{20} & a_{21} & a_{22} & a_{23} \\ a_{30} & a_{31} & a_{32} & a_{33} \end{bmatrix}
$$

$U_0(A)$

| $a_{00}$ | $a_{11}$ | $a_{22}$ | $a_{33}$ |
|---|---|---|---|

$U_1(A)$

| $a_{01}$ | $a_{12}$ | $a_{23}$ | $a_{30}$ |
|---|---|---|---|

$U_2(A)$

| $a_{02}$ | $a_{13}$ | $a_{20}$ | $a_{31}$ |
|---|---|---|---|

$U_3(A)$

| $a_{03}$ | $a_{10}$ | $a_{21}$ | $a_{32}$ |
|---|---|---|---|

Upper diagonal

[HS14] S. Halevi, V. Shoup. "Algorithms in HElib", CRYPTO 2014
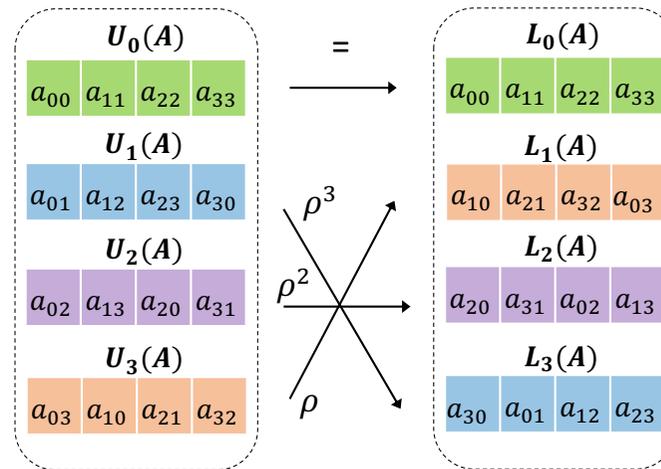
# Matrix-Vector Multiplication

- [HS14] "If an input matrix is provided in plaintext or we can process it before the multiplication, the best solution is to put it in diagonal order, which would lets us use the parallel *systolic* multiplication algorithm."

  ○ $w = U_0(A) \odot v + U_1(A) \odot \rho(v) + \cdots + U_{d-1}(A) \odot \rho^{d-1}(v)$

  ○ Perform in a systolic array

    ▪ Multiply between diagonal and rotated vector

    ▪ Accumulate

    ▪ Pass data to neighbor

$$
\begin{array}{c}
A \\
\begin{array}{|c|c|c|c|}
\hline
a_{00} & a_{01} & a_{02} & a_{03} \\
\hline
a_{10} & a_{11} & a_{12} & a_{13} \\
\hline
a_{20} & a_{21} & a_{22} & a_{23} \\
\hline
a_{30} & a_{31} & a_{32} & a_{33} \\
\hline
\end{array}
\end{array}
\times
\begin{array}{c}
v \\
\begin{array}{|c|}
\hline
v_0 \\
\hline
v_1 \\
\hline
v_2 \\
\hline
v_3 \\
\hline
\end{array}
\end{array}
=
\begin{array}{c}
w \\
\begin{array}{|c|}
\hline
w_0 \\
\hline
w_1 \\
\hline
w_2 \\
\hline
w_3 \\
\hline
\end{array}
\end{array}
$$

[HS14] S. Halevi, V. Shoup. "Algorithms in HElib", CRYPTO 2014

# Diagonal-major Matrix Encoding
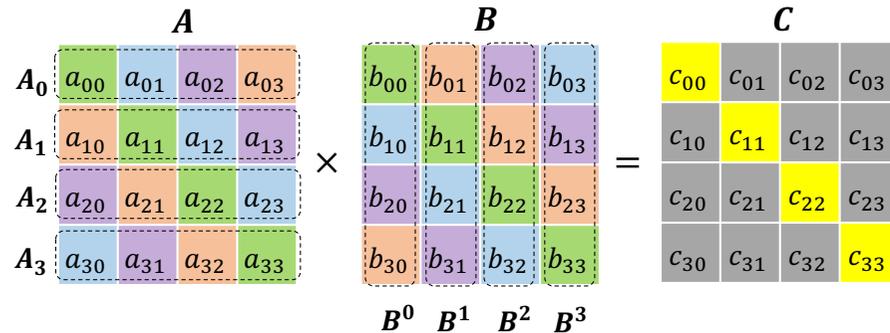


$$L_r(A) = U_r(A^T)$$
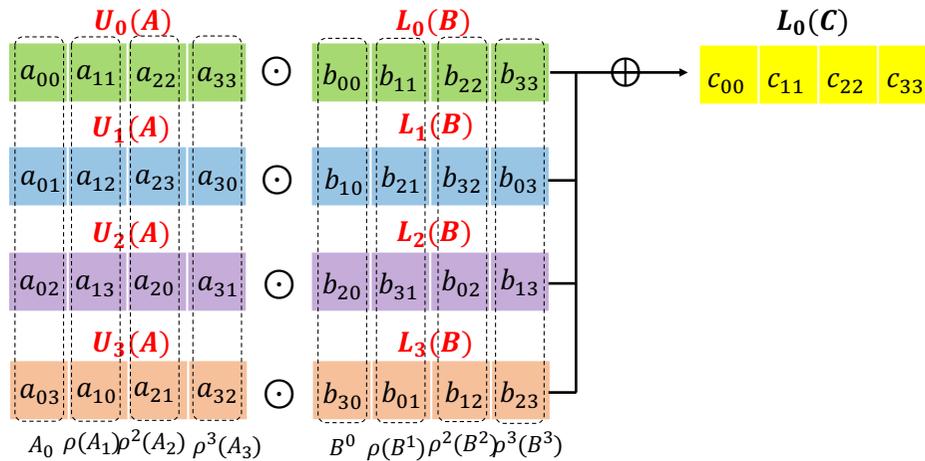
*$\rho^k$: left-rotation by $k$ positions

# Main Idea



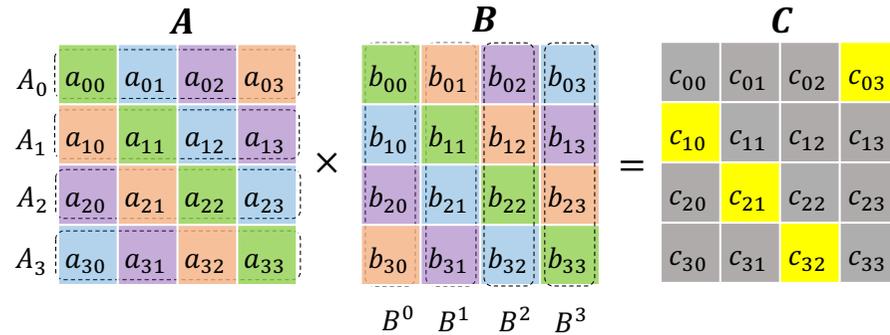- $L_0(C)$: the 0th lower diagonal
  - $c_{00} = \langle A_0, B^0 \rangle$
  - $c_{11} = \langle A_1, B^1 \rangle = \langle \rho(A_1), \rho(B^1) \rangle$
  - $c_{22} = \langle A_2, B^2 \rangle = \langle \rho^2(A_2), \rho^2(B^2) \rangle$
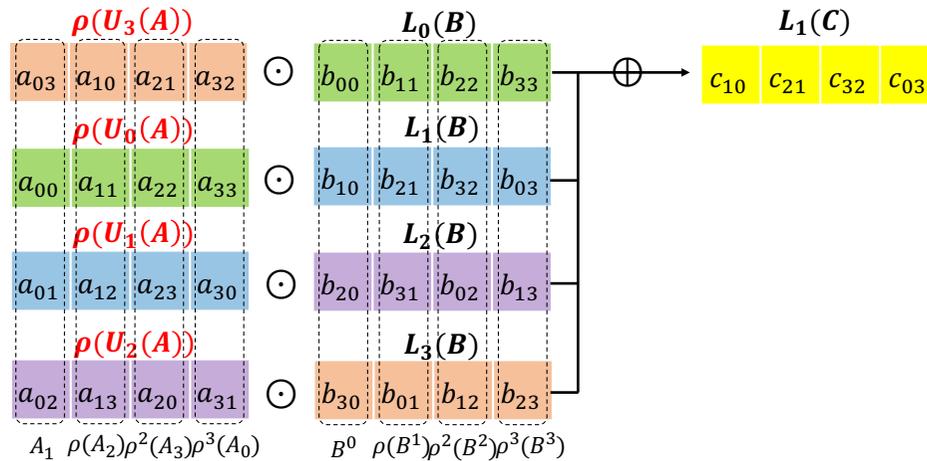  - $c_{33} = \langle A_3, B^3 \rangle = \langle \rho^3(A_3), \rho^3(B^3) \rangle$

# Main Idea



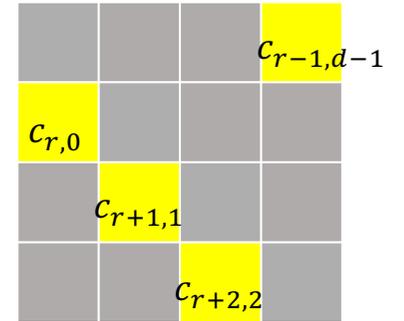- $L_1(C)$: the 1st lower diagonal
  - $c_{10} = \langle A_1, B^0 \rangle$
  - $c_{21} = \langle A_2, B^1 \rangle = \langle \rho(A_2), \rho(B^1) \rangle$
  - $c_{32} = \langle A_3, B^2 \rangle = \langle \rho^2(A_3), \rho^2(B^2) \rangle$
  - $c_{03} = \langle A_0, B^3 \rangle = \langle \rho^3(A_0), \rho^3(B^3) \rangle$

# Basic Square Matrix Multiplication

- Given $A, B \in \mathbb{R}^{d \times d}$, the $r$-th lower diagonal $L_r(C)$ can be obtained as follows:

  - $c_{r,0} = \langle A_r, B^0 \rangle$

  - $c_{r+1,1} = \langle A_{r+1}, B^1 \rangle = \langle \rho(A_{r+1}), \rho(B^1) \rangle$

    $\vdots$

  - $c_{r-1,d-1} = \langle A_{r-1}, B^{d-1} \rangle = \langle \rho^{d-1}(A_{r-1}), \rho^{d-1}(B^{d-1}) \rangle$

$$L_r(C) = (c_{r,0}, c_{r+1,1}, \ldots, c_{r-1,d-1}) = \sum_{0 \le l < d} \rho^r\big(U_{l-r}(A)\big) \odot L_l(B)$$

- In the context of HE computation,

  - PC-MM: When $L_l(B)$ is encrypted, $[\![L_r(C)]\!] = \sum_{l=0}^{d-1} \text{PMult}(\rho^r\big(U_{l-r}(A)\big), [\![L_l(B)]\!])$

  - CC-MM: When both are encrypted, $[\![L_r(C)]\!] = \sum_{l=0}^{d-1} \text{Mult}(\rho^r([\![U_{l-r}(A)]\!]), [\![L_l(B)]\!])$

# HE-friendly Matrix Multiplication

$$L_r(C) = (C_{r,0}, C_{r+1,1}, \ldots, C_{r-1,d-1})$$
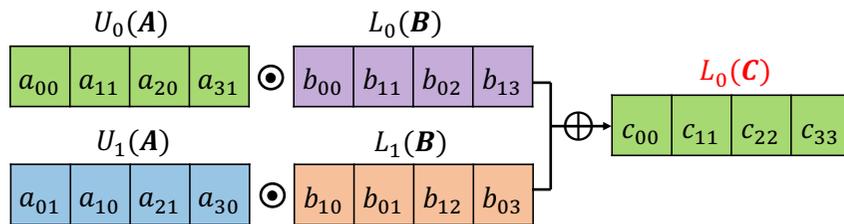$$= \sum_{l=0}^{d-1} \rho^r(U_{l-r}(A)) \odot L_l(B)$$

(PC-MM) $[\![L_r(C)]\!] = \sum_{l=0}^{d-1} \mathrm{Mult}(\rho^r(U_{l-r}(A)), [\![L_l(B)]\!])$

(CC-MM) $[\![L_r(C)]\!] = \sum_{l=0}^{d-1} \mathrm{Mult}(\rho^r([\![U_{l-r}(A)]\!]), [\![L_l(B)]\!])$

✓ Easy to implement

✓ Optimized for both PC-MM and CC-MM

✓ Matrix transposition for free but a different format ($L_r(A) = U_r(A^T)$)

✓ *Unified* encrypted matrix representation (reusable for subsequent computations)

✓ Easily extended to parallel matrix multiplications
   (interlace multiple matrices & batch the computation)

✓ Easily extended to matrix multiplication of various sizes

# Non-Square Matrix Multiplication

# Multi-diagonal Batched Encryption

$A$

| $a_{00}$ | $a_{01}$ | $a_{02}$ | $a_{03}$ |
|---|---|---|---|
| $a_{10}$ | $a_{11}$ | $a_{12}$ | $a_{13}$ |
| $a_{20}$ | $a_{21}$ | $a_{22}$ | $a_{23}$ |
| $a_{30}$ | $a_{31}$ | $a_{32}$ | $a_{33}$ |

$B$

| $b_{00}$ | $b_{01}$ | $b_{02}$ | $b_{03}$ |
|---|---|---|---|
| $b_{10}$ | $b_{11}$ | $b_{12}$ | $b_{13}$ |
| $b_{20}$ | $b_{21}$ | $b_{22}$ | $b_{23}$ |
| $b_{30}$ | $b_{31}$ | $b_{32}$ | $b_{33}$ |

- Diagonal packing capacity $c := s/(nH)$

  ($n$: dimension, $H$: number of parallelized matrices)

- **Batch ciphertext computation**

- **Case 1**: $s = 4$

| $a_{00}$ | $a_{11}$ | $a_{22}$ | $a_{33}$ |
|---|---|---|---|

| $a_{10}$ | $a_{21}$ | $a_{32}$ | $a_{03}$ |
|---|---|---|---|

| $a_{20}$ | $a_{31}$ | $a_{02}$ | $a_{13}$ |
|---|---|---|---|

| $a_{30}$ | $a_{01}$ | $a_{12}$ | $a_{23}$ |
|---|---|---|---|

- **Case 2**: $s = 8$

$\bar{L}_0(A, B)$

| $a_{00}$ | $b_{00}$ | $a_{11}$ | $b_{11}$ | $a_{22}$ | $b_{22}$ | $a_{33}$ | $b_{33}$ |
|---|---|---|---|---|---|---|---|

$\bar{L}_1(A, B)$

| $a_{10}$ | $b_{10}$ | $a_{21}$ | $b_{21}$ | $a_{32}$ | $b_{32}$ | $a_{03}$ | $b_{03}$ |
|---|---|---|---|---|---|---|---|

$\bar{L}_2(A, B)$

| $a_{20}$ | $b_{20}$ | $a_{31}$ | $b_{31}$ | $a_{02}$ | $b_{02}$ | $a_{13}$ | $b_{13}$ |
|---|---|---|---|---|---|---|---|

$\bar{L}_3(A, B)$

| $a_{30}$ | $b_{30}$ | $a_{01}$ | $b_{01}$ | $a_{12}$ | $b_{12}$ | $a_{23}$ | $b_{23}$ |
|---|---|---|---|---|---|---|---|

- **Case 3**: $s = 16 \Rightarrow c = 2$

$\bar{L}_0(A, B)$       $\bar{L}_1(A, B)$

| $a_{00}$ | $b_{00}$ | $a_{11}$ | $b_{11}$ | $a_{22}$ | $b_{22}$ | $a_{33}$ | $b_{33}$ | $a_{10}$ | $b_{10}$ | $a_{21}$ | $b_{21}$ | $a_{32}$ | $b_{32}$ | $a_{03}$ | $b_{03}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$\hat{L}_0(A, B)$

$\bar{L}_2(A, B)$       $\bar{L}_3(A, B)$

| $a_{20}$ | $b_{20}$ | $a_{31}$ | $b_{31}$ | $a_{02}$ | $b_{02}$ | $a_{13}$ | $b_{13}$ | $a_{30}$ | $b_{30}$ | $a_{01}$ | $b_{01}$ | $a_{12}$ | $b_{12}$ | $a_{23}$ | $b_{23}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$\hat{L}_1(A, B)$

e.g., $s = 2^{15}, H = 16, n = 2^7$

$\Rightarrow c = \dfrac{2^{15}}{2^7 \cdot 16} = 2^4$ different diagonals in one ciphertext

20

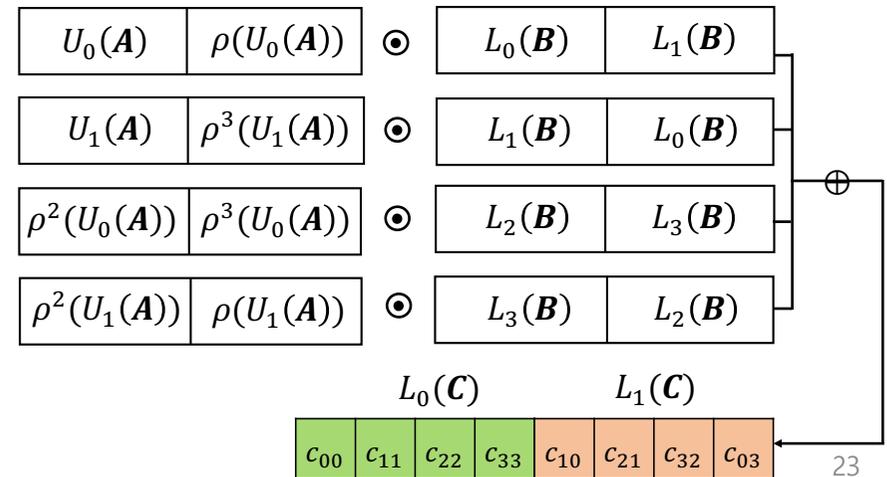# Non-Square Matrix Multiplication (Packed Version)

# Table of Contents

# M1. CC-MM: U-L Approach

- Given encrypted $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times n}$ with $m \leq n$, compute $AB$
  - Each ciphertext of $A$ contains the *same diagonal* of $A$ but *rotated differently* (need $mn$ rotations)



$$(m+1)n \text{ rotations}$$
for $A^{(i)} \in \mathbb{R}^{m \times n}$ and $B^{(i)} \in \mathbb{R}^{n \times n}$

$$
\begin{array}{c}
U_0(A) \odot L_0(B) \\
U_1(A) \odot L_1(B) \\
\rho^2(U_0(A)) \odot L_2(B) \\
\rho^2(U_1(A)) \odot L_3(B)
\end{array}
\oplus \rightarrow L_0(C)
\qquad
\begin{array}{c}
\rho^3(U_1(A)) \odot L_0(B) \\
\rho(U_0(A)) \odot L_1(B) \\
\rho(U_1(A)) \odot L_2(B) \\
\rho^3(U_0(A)) \odot L_3(B)
\end{array}
\oplus \rightarrow L_1(C)
$$

| $U_0(A)$ | $\rho(U_0(A))$ | $\odot$ | $L_0(B)$ | $L_1(B)$ |
| --- | --- | --- | --- | --- |
| $U_1(A)$ | $\rho^3(U_1(A))$ | $\odot$ | $L_1(B)$ | $L_0(B)$ |
| $\rho^2(U_0(A))$ | $\rho^3(U_0(A))$ | $\odot$ | $L_2(B)$ | $L_3(B)$ |
| $\rho^2(U_1(A))$ | $\rho(U_1(A))$ | $\odot$ | $L_3(B)$ | $L_2(B)$ |

$L_0(C)$   $L_1(C)$

| $c_{00}$ | $c_{11}$ | $c_{22}$ | $c_{33}$ | $c_{10}$ | $c_{21}$ | $c_{32}$ | $c_{03}$ |
| --- | --- | --- | --- | --- | --- | --- | --- |

# M2. CC-MM: L-L Approach

- Given encrypted $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times n}$ with $m \leq n$, compute $AB$

  o How about packing *replicated* diagonals of $B$?

  o Use $U_0(A) = L_0(A)$, $U_1(A) = \rho(L_1(A))$

  o Each ciphertext of $A$ contains the *different lower diagonal* but rotated by *the same amount*.

  $\Rightarrow \mathbf{2mn/c}$ rotations for $\left(A^{(i)}B^{(i)}\right)_{1 \leq i \leq H}$ (where $c = s/nH$).



| $L_0(A)$ | $L_1(A)$ | | $U_0(A)$ | $\rho^3(U_1(A))$ | $\odot$ | $L_0(B)$ | $L_0(B)$ |
| --- | --- | --- | --- | --- | --- | --- | --- |
| $\rho(L_1(A))$ | $\rho(L_0(A))$ | $=$ | $U_1(A)$ | $\rho(U_0(A))$ | $\odot$ | $L_1(B)$ | $L_1(B)$ |
| $\rho^2(L_0(A))$ | $\rho^2(L_1(A))$ | $=$ | $\rho^2(U_0(A))$ | $\rho(U_1(A))$ | $\odot$ | $L_2(B)$ | $L_2(B)$ |
| $\rho^3(L_1(A))$ | $\rho^3(L_0(A))$ | $=$ | $\rho^2(U_1(A))$ | $\rho^3(U_0(A))$ | $\odot$ | $L_3(B)$ | $L_3(B)$ |

$L_0(C)$ $\qquad$ $L_1(C)$

| $c_{00}$ | $c_{11}$ | $c_{22}$ | $c_{33}$ | $c_{10}$ | $c_{21}$ | $c_{32}$ | $c_{03}$ |
| --- | --- | --- | --- | --- | --- | --- | --- |

# M2. CC-MM: BSGS-integrated L-L Approach

- Rotate the ciphertexts of $A$ internally $\Rightarrow$ rotate them at first, and apply *only a single rotation* to the summed intermediate result.
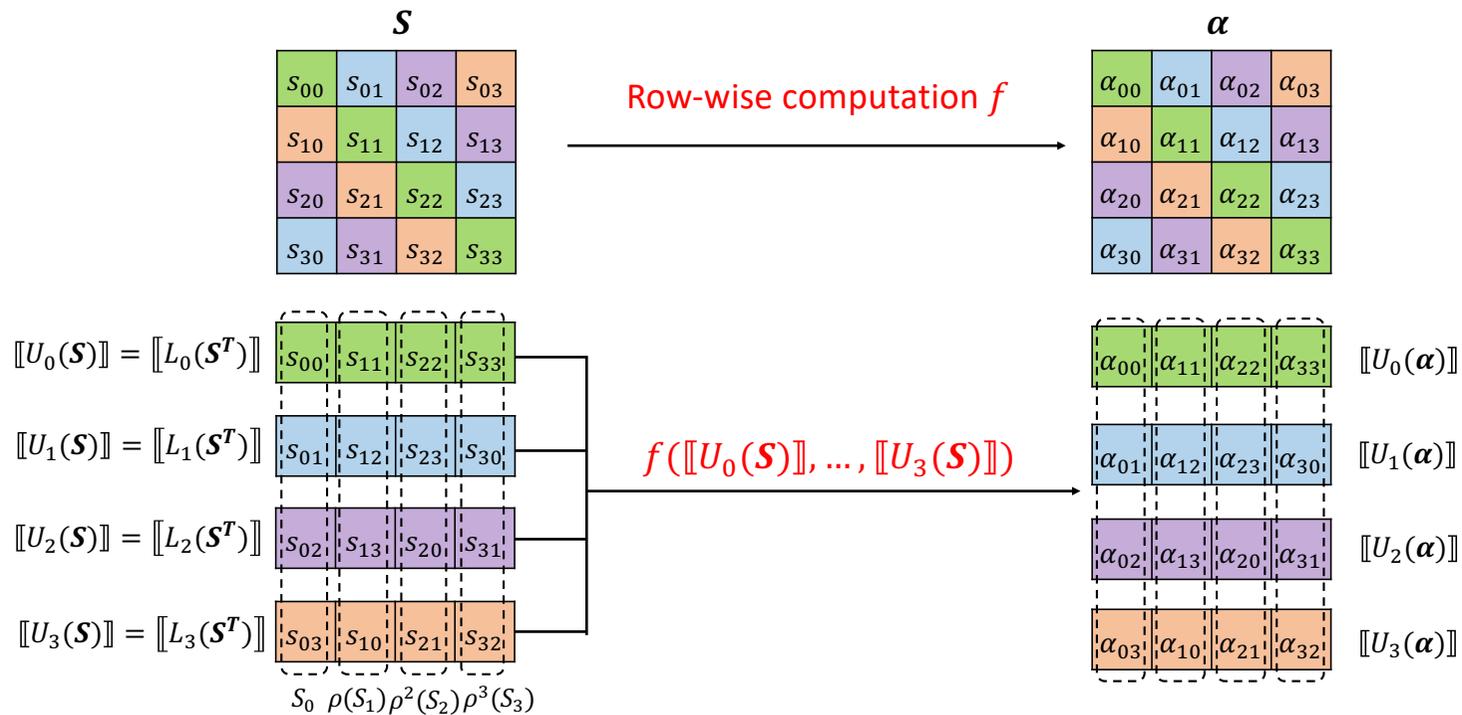


$$mn \Rightarrow 2mn/c \Rightarrow mn/c \text{ Rotations}$$
$$\text{for } A^{(i)} \in \mathbb{R}^{m \times n}, B^{(i)} \in \mathbb{R}^{n \times n}$$

# Comparison

| | JKLS'18 | Our work |
|---|---|---|
| Encoding format | Row-major | Diagonal-major |
| *Parallel* large-scale matrix multiplication | ✗ | ✔ |
| PC-MM | Inefficient | Efficient |
| Transposition | Non-trivial | For free;<br>but a different format |

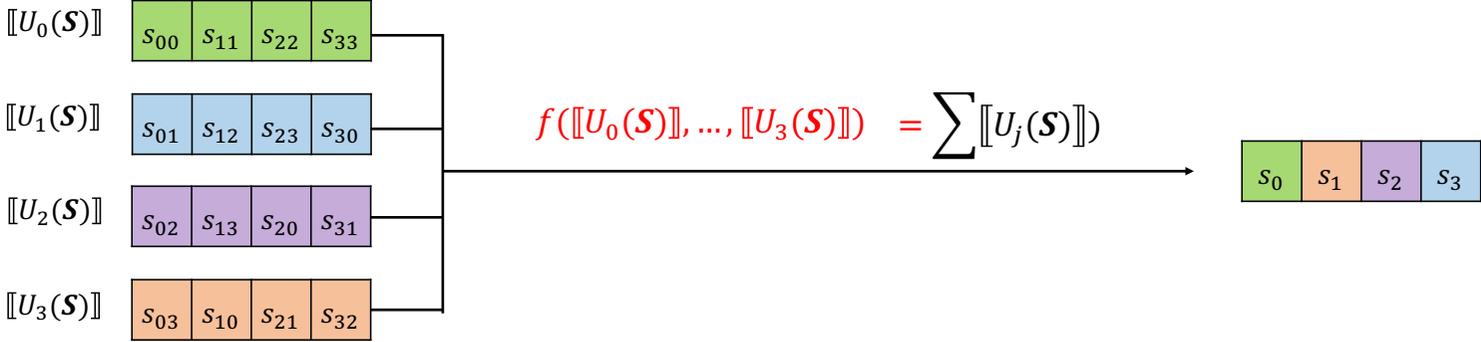| | Equation | $H$ | | | | |
|---|---|---|---|---|---|---|
| | | 1 | 2 | 4 | 8 | 16 |
| JKLS'18 | $H(5\sqrt{n}+4n)$ | 572 | 1144 | 2272 | 4576 | 9152 |
| Our work | $H(n/2+3)+n/4(\log(n/H)+3)$ | 387 | 422 | 520 | 760 | 1264 |

CC-MM computation of $A^{(z)}B^{(z)}$ for $1 \le z \le H$, where $A^{(z)}, B^{(z)} \in \mathbb{R}^{n \times n}$ and $s = n^2$
$(n = 2^7, s = 2^{14})$

# Row-wise Computation over Diagonals

- If the upper diagonals of $S$ are provided (= lower diagonals of $S^T$)
  - Each row of $S$ is distributed across the same position in different slots.
  - Apply $f$ row-wise on $S \Rightarrow$ apply $f$ directly to ciphertexts (no computation is needed between slots)

# Example



$$f(s_{i0}, s_{i1}, s_{i2}, s_{i3}) = \sum_j s_{ij} = s_i$$

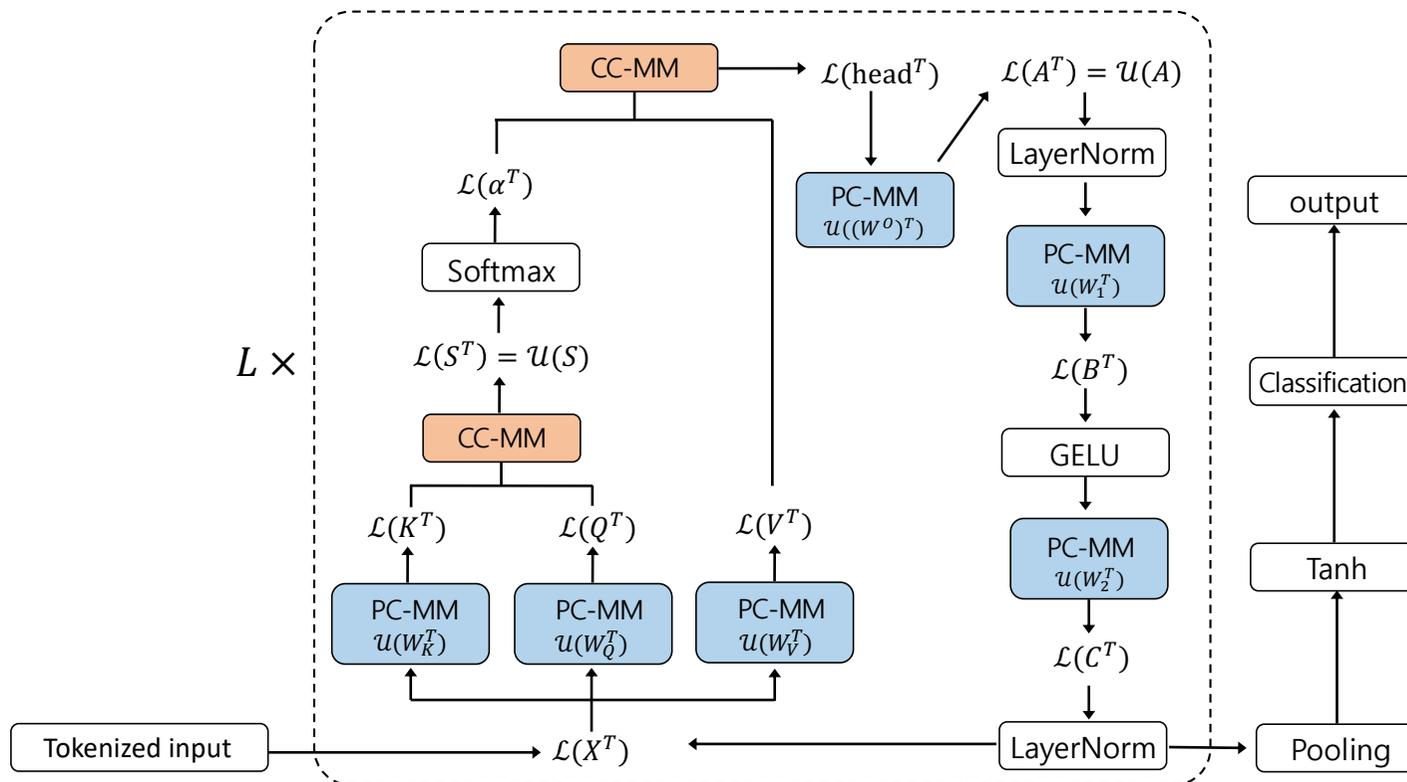$$f(\llbracket U_0(S) \rrbracket, \ldots, \llbracket U_3(S) \rrbracket) = \sum \llbracket U_j(S) \rrbracket)$$

# Application to Transformer-based Inference

- THOR: secure Transformer inference framework with HomomOrphic encRyption

# Experimental Results

- BERT-base model ($L$ =12 layers, 110M parameters, $n$ =128 tokens)

- **10 minutes** with only a 0.8% accuracy drop (2.7hours in NEXUS): **16x** improvement

  o   Softmax: square-and-normalization (see paper)

  o   Use the Goldschmidt's algorithm + Adaptive Successive over-relaxation method (aSOR)

*Intel Xeon Platinum 8462 at 2.8GHz,*
*A100 GPU (DESILO FHE library)*

| Operation | Input | Time (sec) |
|---|---|---|
| Attention layer | $3 \times (\mathbb{R}^{128\times768} \times \mathbb{R}^{768\times64})$ | 49.77 |
| Attention score | $12 \times (\mathbb{R}^{128\times64} \times \mathbb{R}^{64\times128})$ | 16.25 |
| Softmax | $12 \times (\mathbb{R}^{128\times128})$ | 15.53 |
| Attention head | $12 \times (\mathbb{R}^{128\times128} \times \mathbb{R}^{128\times64})$ | 13.08 |
| Multi-head attention | $\mathbb{R}^{128\times768} \times \mathbb{R}^{768\times768}$ | 27.43 |
| LayerNorm1 | $\mathbb{R}^{128\times768}$ | 7.13 |
| FC1 | $\mathbb{R}^{128\times768} \times \mathbb{R}^{768\times3072}$ | 49.80 |
| GELU | $\mathbb{R}^{128\times3072}$ | 29.42 |
| FC2 | $\mathbb{R}^{128\times3072} \times \mathbb{R}^{3072\times768}$ | 49.19 |
| LayerNorm2 | $\mathbb{R}^{128\times768}$ | 4.10 |
| Pooler & Classification | $\mathbb{R}^{128\times768}$ | 2.70 |
| Bootstrappings | - | 337.86 |
| Total | - | **602.26** |

| Dataset | #Test | Metric | Unencrypted | | | | Encrypted |
|---|---|---|---|---|---|---|---|
| | | | Baseline | G | G-LN | G-LN-S | |
| MPRC | 408 | Accuracy | 85.29 | 85.54 | 85.54 | 85.78 | **84.80** |
| | | F1-score | 89.90 | 90.05 | 90.05 | 90.24 | **89.49** |
| RTE | 277 | Accuracy | 72.20 | 71.48 | 71.84 | 72.20 | **71.12** |
| SST-2 | 872 | Accuracy | 91.51 | 91.40 | 91.40 | 91.63 | **90.71** |

MRPC: Sentence pair 2-class paraphrase

RTE: Premise/hypothesis 2-class entailment

SST-2: Single-sentence 2-class paraphrase

[Zhang+25] "Secure transformer inference made noninteractive", NDSS'25

# Table of Contents

# Summary

- **New Efficient parallel matrix computation**
  - Block-wise PC-MM
  - Optimized CC-MM using the BSGS strategy
  - Row-wise computation over diagonals

- **Question.** Can the proposed matrix computation be further optimized?
  - Level consumption (vs. 3 levels in JKLS18)
    - PC-MM: 2 levels (one for masking and one for plaintext mult)
    - CC-MM: 2 levels (one for making and one for ciphertext mult)
  - Need *replication* during CC-MM
    - Generate replicated (batched) lower diagonals
    - $\frac{n}{2}(\log c + 1)$ rotations for $\boldsymbol{B}^{(i)} \in \mathbb{R}^{n \times n}$ for $1 \leq i \leq H$ with one level ($c = s/nH$)

# Future Works

- Extension
    - Low Latency vs *high-throughput* batch processing
    - Extend secure inference to larger parameter LLMs (e.g., LLaMA)
    - Privacy-preserving training for LLMs

- Potential applications beyond LLM
    - Genomic analysis (e.g., genotype imputation)
    - Machine learning (e.g., neural networks, transfer learning)

Implementation available at:
`https://github.com/crypto-starlab/THOR`    ia.cr/2024/1881