Ciphertext-Ciphertext Matrix Multiplication: Fast for Large Matrices

> Jai Hyun Park jaihyunp@gmail.com



Madrid, Spain – May 6, 2025

Summary

- Fast ciphertext-ciphertext matrix multiplication (CCMM)
 - 85.2 s for CCMM of 4096 × 4096 matrices in a single thread CPU
 - How? Reduce CCMM to plaintext matrix multiplications
- Fast ciphertext matrix transpose (CMT)
 - 0.76 s for CMT of a 2048×2048 matrices in a single thread CPU
- Lightweight CCMM and CMT algorithms with smaller key sizes



Matrix Multiplication on Encrypted Data



- PPMM: plaintext-plaintext matrix multiplication
- PCMM: plaintext-ciphertext matrix multiplication
- CCMM: ciphertext-ciphertext matrix multiplication
- PCMMs and CCMMs with diverse dimensions
 e.g., PCMM of dimension 128 ~ 16384 for GPT-3.5



PPMM vs. PCMM vs. CCMM

- PPMM BLAS libraries
 - highly optimized open libraries
 - Can be 30x faster than a naïve implementation



PPMM vs. PCMM vs. CCMM

- PPMM BLAS libraries
 - highly optimized open libraries
 - Can be 30x faster than a naïve implementation
- PCMM BCH<u>P</u>S'24
 - Reduction from PCMM to PPMM
 - Optimizations with shared-a, truncation, and others



PPMM vs. PCMM vs. CCMM

- PPMM BLAS libraries
 - highly optimized open libraries
 - Can be 30x faster than a naïve implementation
- PCMM BCH<u>P</u>S'24
 - Reduction from PCMM to PPMM
 - Optimizations with shared-a, truncation, and others
- CCMM JKLS'18
 - cubic bit complexity
 - 0.6 seconds for matrix dimension 64

For matrix dimension 2¹²:

PPMM (OpenBLAS) 1.47 seconds



Ciphertext-Ciphertext Matrix Multiplication (CCMM)

Ciphertext-Ciphertext Matrix Multiplication (CCMM)

- CCMM with RLWE-based (fully) homomorphic encryption schemes
 - Compatibility with the other machine learning tasks
 - High efficiency

Encrypted Matrix Multiplication with CKKS

• CKKS

- Plaintext: <u>vector</u> of real numbers
- Native operations: // add, // mult, and rotate

Encrypted Matrix Multiplication with CKKS

• CKKS

- Plaintext: vector of real numbers
- Native operations: // add, // mult, and rotate
- With the native operations, PCMM/CCMM requires lots of rotates
 - For example, [JKLS'18] has a cubic bit complexity, but is orders of magnitude slower than PPMM

Encrypted Matrix Multiplication with CKKS

• CKKS

- Plaintext: vector of real numbers
- Native operations: // add, // mult, and rotate
- With the native operations, PCMM/CCMM requires lots of rotates
 - For example, [JKLS'18] has a cubic bit complexity, but is orders of magnitude slower than PPMM

Q. How to utilize PPMM BLAS libraries?

A. Reduction from PCMM/CCMM to PPMM

 $[a_0$

RLWE-based Encryption of Matrices

• In the ring $\mathbb{Z}_Q[X]/(X^N+1)$, an RLWE ciphertext $(a, b = -a \cdot s + m)$ is:

$$\begin{bmatrix} s_{0} & s_{1} & \cdots & s_{N-1} \\ -s_{N-1} & s_{0} & \cdots & s_{N-2} \\ \vdots & \vdots & \ddots & \vdots \\ -s_{1} & -s_{2} & \cdots & s_{0} \end{bmatrix} + \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix} = \begin{bmatrix} m_{0} & m_{1} & \cdots & m_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix} = \begin{bmatrix} m_{0} & m_{1} & \cdots & m_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix} = \begin{bmatrix} m_{0} & m_{1} & \cdots & m_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

$$+ \begin{bmatrix} b_{0} & b_{1} & \cdots & b_{N-1} \end{bmatrix}$$

RLWE-based Encryption of Matrices

• In the ring $\mathbb{Z}_Q[X]/(X^N + 1)$, an RLWE ciphertext $(a, b = -a \cdot s + m)$ is:

• *N* RLWE ciphertexts are:

Jai Hyun Park

Jai Hyun Park

♦ N × N × N PCMM ≤ two N × N × N PPMMs modulo Q
♦ We use fast PPMM BLAS libraries for N × N × N PCMM
♦ For PCMMs with other dimensions, see BCH<u>P</u>S'24

Ciphertext-Ciphertext Matrix Multiplication: Fast for Large Matrices

Ciphertext-Ciphertext Matrix Multiplication: Fast for Large Matrices

Ciphertext-Ciphertext Matrix Multiplication: Fast for Large Matrices

Encrypting each row: $a_i s + b_i = m_i$

Encrypting each row: $a_i s + b_i = m_i$

Encrypting each column: $a_j s + b_j = m'_j$

Encrypting each row: $a_i s + b_i = m_i$

Encrypting each column: $a_j s + b_j = m'_j$

Trace

$$\forall i, j, \qquad M_{i,j} = N^{-1}$$

$$\cdot \quad \sum_{\sigma \in \operatorname{Aut}} \sigma(m_i(X) \cdot X^{-j})$$

N ciphertexts

Jai Hyun Park

∀i,j,

∀*j*,

CMT with *N* keyswitchings

N ciphertexts

Trace

Sharing

automorphisms

Trace
$$\forall i, j, \quad M_{i,j} = N^{-1} \cdot \sum_{\sigma \in Aut} \sigma(m_i(X) \cdot X^{-j})$$

Sharing
automorphisms $\forall j, \quad m'_j(X) = N^{-1} \cdot \sum_{i \in [N]} \sum_{\sigma \in Aut} \sigma(X^{-j} \cdot m_i(X)) \cdot X^i$
 $= N^{-1} \cdot \sum_{\sigma \in Aut} \sum_{i \in [N]} \sigma(X^{-j}) \cdot \sigma(m_i(X)) \cdot X^i$
 $= N^{-1} \cdot \sum_{\sigma \in Aut} \sigma(X^{-j}) \cdot \sigma\left(\sum_{i \in [N]} m_i(X) \cdot \sigma^{-1}(X^i)\right)$
 $m'_j(X) = \sum_{i \in [N]} M_{i,j} X^i$
 N ciphertexts

Trace
$$\forall i, j, \quad M_{i,j} = N^{-1} \cdot \sum_{\sigma \in Aut} \sigma(m_i(X) \cdot X^{-j})$$

Sharing
automorphisms $\forall j, \quad m'_j(X) = N^{-1} \cdot \sum_{i \in [N]} \sum_{\sigma \in Aut} \sigma(X^{-j} \cdot m_i(X)) \cdot X^i$
 $= N^{-1} \cdot \sum_{\sigma \in Aut} \sum_{i \in [N]} \sigma(X^{-j}) \cdot \sigma(m_i(X)) \cdot X^i$
 $= N^{-1} \cdot \sum_{\sigma \in Aut} \sigma(X^{-j}) \left[\sigma\left(\sum_{i \in [N]} m_i(X) \cdot \sigma^{-1}(X^i)\right) \right]$
 $m'_j(X) = \sum_{i \in [N]} M_{i,j} X^i$
 $M_{i,j} X^i$
 $M_{i,j} X^i$
 N ciphertexts

Trace
$$\forall i, j, \quad M_{i,j} = N^{-1} \cdot \sum_{\sigma \in Aut} \sigma(m_i(X) \cdot X^{-j})$$

Sharing
automorphisms $\forall j, \quad m'_j(X) = N^{-1} \cdot \sum_{i \in [N]} \sum_{\sigma \in Aut} \sigma(X^{-j} \cdot m_i(X)) \cdot X^i$
 $= N^{-1} \cdot \sum_{\sigma \in Aut} \sum_{i \in [N]} \sigma(X^{-j}) \cdot \sigma(m_i(X)) \cdot X^i$
 $= N^{-1} \cdot \sum_{\sigma \in Aut} \sum_{i \in [N]} \sigma(X^{-j}) \cdot \sigma(m_i(X)) \cdot X^i$
 $= N^{-1} \cdot \sum_{\sigma \in Aut} \sigma(X^{-j}) \int \sigma\left(\sum_{i \in [N]} m_i(X) \cdot \sigma^{-1}(X^i)\right)$
 $[m_i]_i \longrightarrow N^2 \text{ add} [\widetilde{m_\sigma}]_\sigma \longrightarrow [\sigma(\widetilde{m_\sigma})]_\sigma \longrightarrow N^2 \text{ add} [m'_j]_j$
 $M \text{ ciphertexts}$
 $M \text{ ciphertexts}$

$$\checkmark \quad \widetilde{m_{\sigma}}(X) = \sum_{i} \sigma^{-1} (X^{i}) \cdot m_{i}$$

$$\checkmark \quad m'_{j}(X) = \sum_{\sigma} \sigma (X^{-j}) \cdot \sigma (\widetilde{m_{\sigma}})$$

$$\checkmark \quad \widetilde{m_{\sigma}}(X) = \sum_{i} \sigma^{-1} (X^{i}) \cdot m_{i}$$

$$\checkmark \quad m'_{j}(X) = \sum_{\sigma} \sigma (X^{-j}) \cdot \sigma (\widetilde{m_{\sigma}})$$

$$\operatorname{Tweak}(\{m_i\}_{i\in[n]}) \mapsto \left\{\sum_{i\in[n]} X^{\frac{2N}{n}ij} \cdot m_i\right\}_{j\in[n]}$$

$$\checkmark \quad \widetilde{m_{\sigma}}(X) = \sum_{i} \sigma^{-1} (X^{i}) \cdot m_{i}$$

$$\checkmark \quad m'_{j}(X) = \sum_{\sigma} \sigma (X^{-j}) \cdot \sigma (\widetilde{m_{\sigma}})$$

$$\operatorname{Tweak}(\{m_i\}_{i \in [n]}) \mapsto \left\{ \sum_{i \in [n]} X^{\frac{2N}{n}ij} \cdot m_i \right\}_{j \in [n]}$$

- Tweak $({m_i}_{i \in [n]})$ can be done with
 - Tweak $(\{m_{2i}\}_{i \in [n/2]})$
 - Tweak $({m_{2i+1}}_{i \in [n/2]})$
 - *n* ring additions

$$\checkmark \quad \widetilde{m_{\sigma}}(X) = \sum_{i} \sigma^{-1} (X^{i}) \cdot m_{i}$$

$$\checkmark \quad m'_{j}(X) = \sum_{\sigma} \sigma (X^{-j}) \cdot \sigma (\widetilde{m_{\sigma}})$$

$$\operatorname{Tweak}(\{m_i\}_{i\in[n]}) \mapsto \left\{\sum_{i\in[n]} X^{\frac{2N}{n}ij} \cdot m_i\right\}_{j\in[n]}$$

- Tweak $(\{m_i\}_{i \in [n]})$ can be done with
 - Tweak $(\{m_{2i}\}_{i \in [n/2]})$
 - Tweak $\overline{(\{m_{2i+1}\}_{i\in[n/2]})}$
 - *n* ring additions

$$\checkmark \quad \widetilde{m_{\sigma}}(X) = \sum_{i} \sigma^{-1} (X^{i}) \cdot m_{i}$$

$$\checkmark \quad m'_{j}(X) = \sum_{\sigma} \sigma (X^{-j}) \cdot \sigma (\widetilde{m_{\sigma}})$$

$$\operatorname{Tweak}(\{m_i\}_{i\in[n]}) \mapsto \left\{\sum_{i\in[n]} X^{\frac{2N}{n}ij} \cdot m_i\right\}_{j\in[n]}$$

- Tweak $(\{m_i\}_{i \in [n]})$ can be done with
 - Tweak $(\{m_{2i}\}_{i \in [n/2]})$
 - Tweak $({m_{2i+1}}_{i \in [n/2]})$
 - *n* ring additions

$$\checkmark \quad \widetilde{m_{\sigma}}(X) = \sum_{i} \sigma^{-1} (X^{i}) \cdot m_{i}$$

$$\checkmark \quad m'_{j}(X) = \sum_{\sigma} \sigma (X^{-j}) \cdot \sigma (\widetilde{m_{\sigma}})$$

$$\operatorname{Tweak}(\{m_i\}_{i\in[n]}) \mapsto \left\{\sum_{i\in[n]} X^{\frac{2N}{n}ij} \cdot m_i\right\}_{j\in[n]}$$

N is the ring degree of RLWE

- Tweak $({m_i}_{i \in [n]})$ can be done with
 - Tweak $(\{m_{2i}\}_{i \in [n/2]})$
 - Tweak $({m_{2i+1}}_{i \in [n/2]})$
 - *n* ring additions

♦ The cost of Tweak $({m_i}_{i \in [n]})$ is $Nn \log n$

CMT Algorithm

$$m_i(X) = m_{i,0} + m_{i,1}X + m_{i,2}X^2 + m_{i,3}X^3$$

Jai Hyun Park

Reduction from CCMM to PPMM $(1/2)^{\dagger}$

Jai Hyun Park

Experimental Results on CCMM

Algorithm	Matrix Dimension	$(\log N, \log Q)$	CMTs	PPMMs	Relin. & Resc.	Total (s)	Prec. (bit)	Key size (MB)
Basic	4096	(12, 36 + 28)	25.5	57.1	2.58	85.2	18.7	436
Basic	8192	(13, 38 + 28)	104	481	11.8	596	18.5	1960
Lightweight	8192	(13, 38 + 28)	186	474	11.8	672	18.5	1.57

All experiments are measured on Intel[®] Xeon[®] Gold 6242 CPU at 2.80GHz with a single-thread

All parameters are 128-bit secure

HEaaN library for HE, FLINT library (based on OpenBLAS) for PPMM

Experimental Results on CMT

Algorithm	Matrix Dimension	$(\log N, \log Q)$	Latency (s)	Prec. (bit)	Key size (MB)
Basic	2048	(11,26)	0.764	10.7	27.3
Basic	4096	(12,28)	3.04	16.3	134
Lightweight	4096	(12,28)	4.92	14.2	0.246

All experiments are measured on Intel[®] Xeon[®] Gold 6242 CPU at 2.80GHz with a single-thread

All parameters are 128-bit secure

HEaaN library for HE

Wrapping up!

- Fast CCMM
 - Leverage efficiency of BLAS libraries
- Fast CMT
 - Useful beyond being as a tool for CCMM
- Lightweight algorithms
 - CCMM with keys less than 2 MB

Wrapping up!

- Fast CCMM
 - Leverage efficiency of BLAS libraries
- Fast CMT
 - Useful beyond being as a tool for CCMM
- Lightweight algorithms
 - CCMM with keys less than 2 MB

eprint: 2025/448 Thank you!

References

[BCH<u>P</u>S'24] Y. Bae, J. H. Cheon, G. Hanrot, <u>J. H. Park</u>, D. Stehlé. "Plaintext-Ciphertext Matrix Multiplication and FHE Bootstrapping: Fast and Fused." Crypto 2024

[JKLS'18] X. Jiang, M. Kim, K. Lauter, Y. Song. "Secure Outsourced Matrix Computation and Application to Neural Networks." CCS 2018

[LZ'22] J. Liu, L. F. Zhang. "*Privacy-preserving and publicly verifiable matrix multiplication."* IEEE Trans. on Services Computing, 2022

Jai Hyun Park