

Plaintext-Ciphertext Matrix Multiplication and FHE Bootstrapping: Fast and Fused

Youngjin Bae, Jung Hee Cheon, Guillaume Hanrot,
Jai Hyun Park, Damien Stehlé

Summary

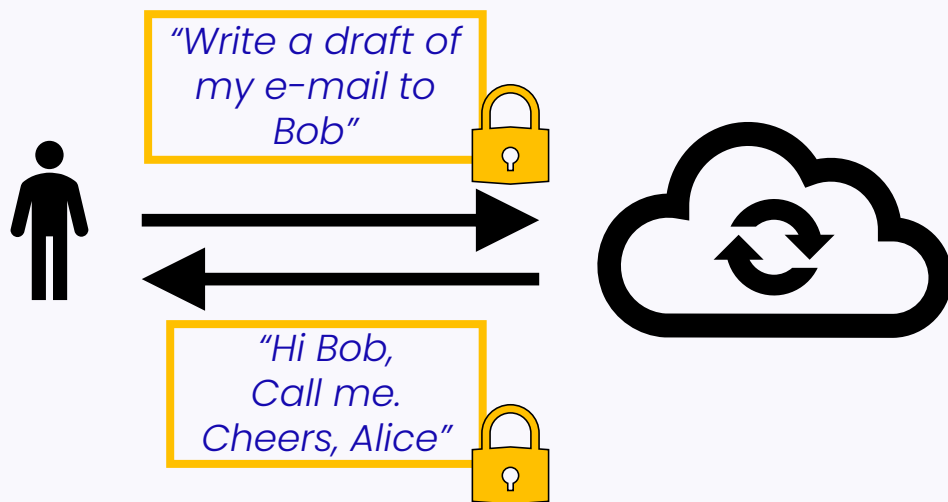
- Fast plaintext-ciphertext matrix multiplication (PCMM)
 - 0.31s for PCMM with 256×256 matrices in a single thread CPU.
 - How? Reduce PCMM to plaintext matrix multiplications.
- Batch bootstrapping (BTS) with high throughput
 - 2x for 32 batches
- Fused PCMM with batch BTS
 - 41% higher throughput than current BTS without PCMM

Matrix Multiplication

- Matrix multiplication is central in high-performance computing
 - highly optimized libraries for basic linear algebra subprograms (BLAS)
 - Can be 10x faster than a naïve implementation for large matrices

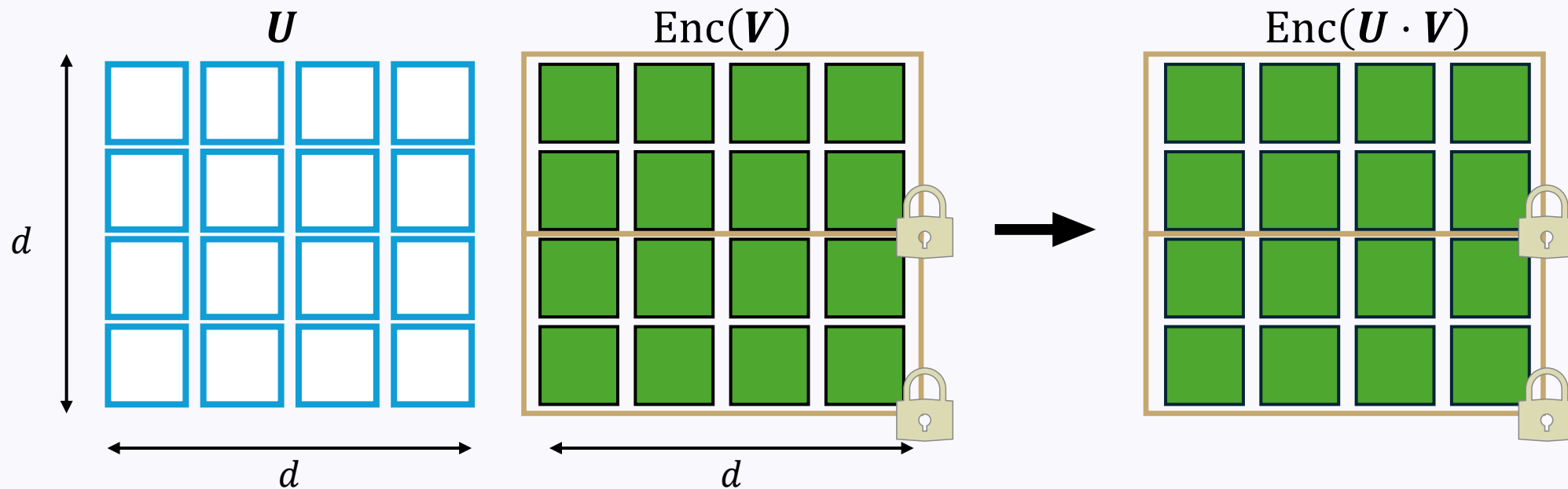
What about matrix multiplication on encrypted data?

Privacy-preserving machine learning as a service



- 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

Plaintext-Ciphertext Matrix Multiplication (PCMM)



- Multiplication between a plaintext matrix and a ciphertext matrix.
 - $d \times d \times d$ PCMM
- PCMM with RLWE-based (fully) homomorphic encryption schemes (CKKS)
 - Compatibility with the other machine learning tasks
 - High efficiency

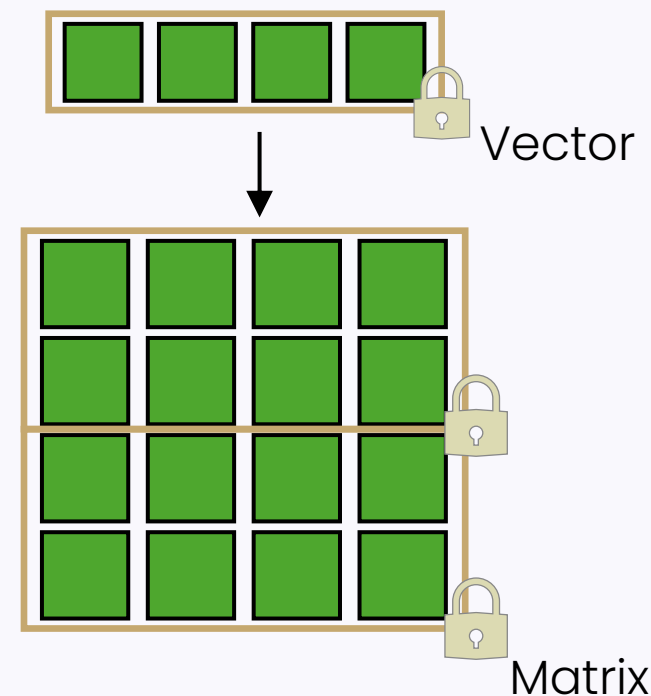
PCMM with CKKS

- CKKS
 - Plaintext: vector of real numbers
 - Native operations: // add, // mult, and rotate.
- With the native operations, PCMM requires lots of rotates.
 - For example, [JKLS18] has a cubic bit complexity, but is orders of magnitude slower than PPMM.

Questions

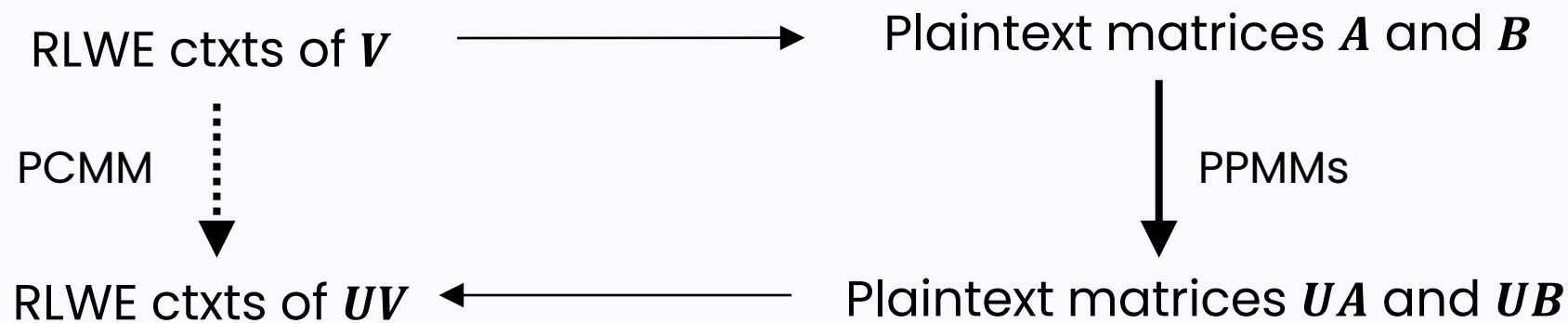
1. How to utilize PPMM BLAS libraries?
2. How to handle multiple ciphertexts?

- ✓ Reduction from PCMM to PPMM
- ✓ Batch ciphertext computation



Reduction from PCMM to PPMs

- [LZ22] considers verifiable PCMM
 - Performs one PCMM using two PPMs
 - Restriction: $d \geq N$



- This is a great idea for fast PCMM: we can use BLAS libraries

RLWE-based Encryption of Matrices

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

$$\begin{bmatrix} a_0 & a_1 & \cdots & a_{N-1} \end{bmatrix} + \begin{bmatrix} s_0 & s_1 & \cdots & s_{N-1} \\ -s_{N-1} & s_0 & \cdots & s_{N-2} \\ \vdots & \vdots & \ddots & \vdots \\ -s_1 & -s_2 & \cdots & s_0 \end{bmatrix} = \begin{bmatrix} m_0 & m_1 & \cdots & m_{N-1} \end{bmatrix}$$

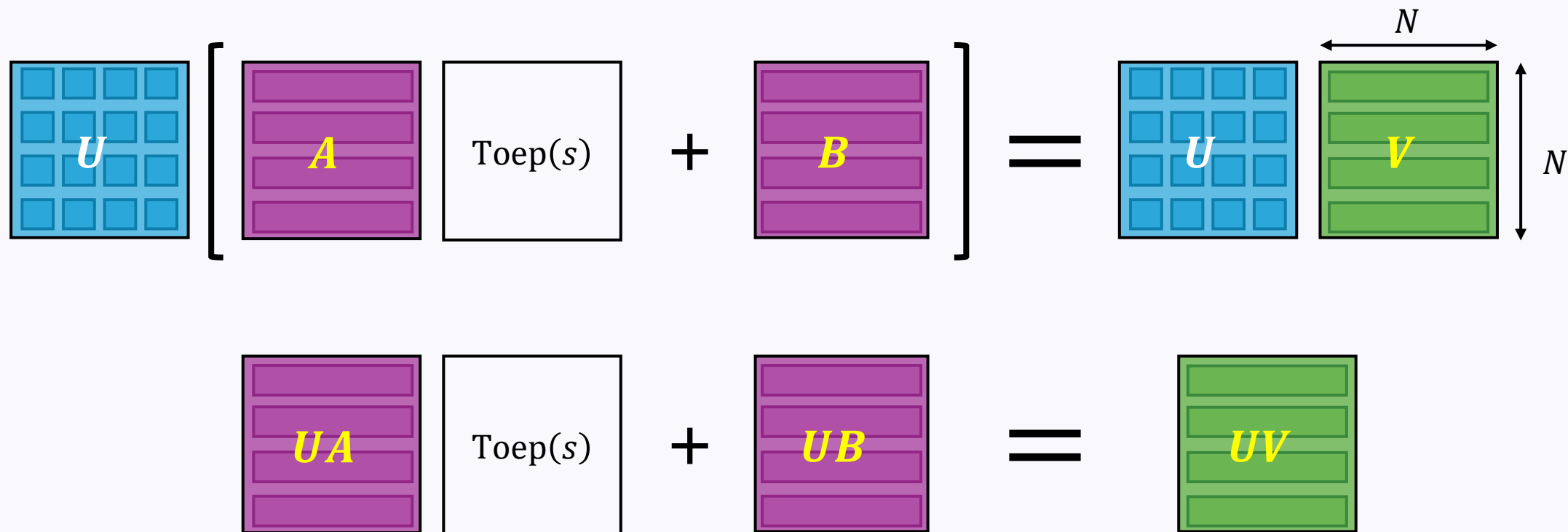
Toep(s)

$\checkmark a_i, b_i, s_i, m_i$ are coeffs of a, b, s, m

- N RLWE ciphertexts (with a shared secret) are:

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

RLWE PCMM \leq PPMMs

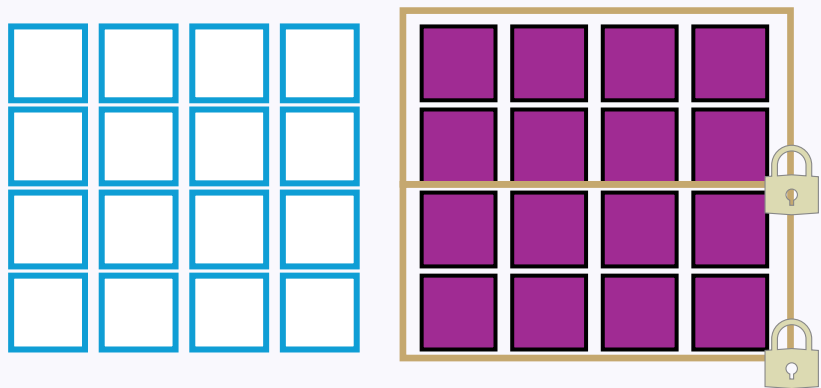


- ❖ $N \times N \times N$ PCMM \leq two $N \times N \times N$ PPMMs modulo Q
- ❖ We use **PPMM BLAS libraries** for PCMM

PCMM with Small/Large Matrices

Small matrices

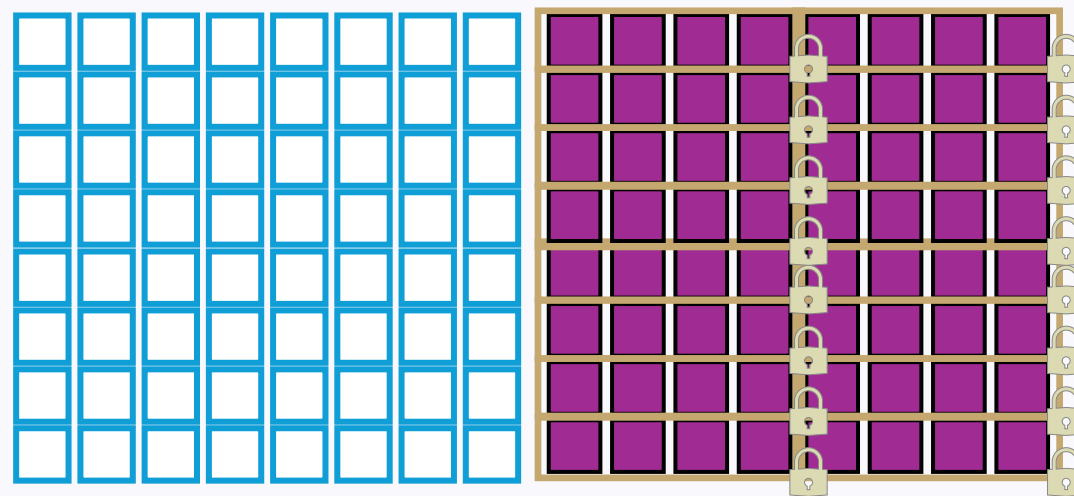
- Data moves within ciphertexts



Can we reduce PCMM to PPMMs?

Large matrices

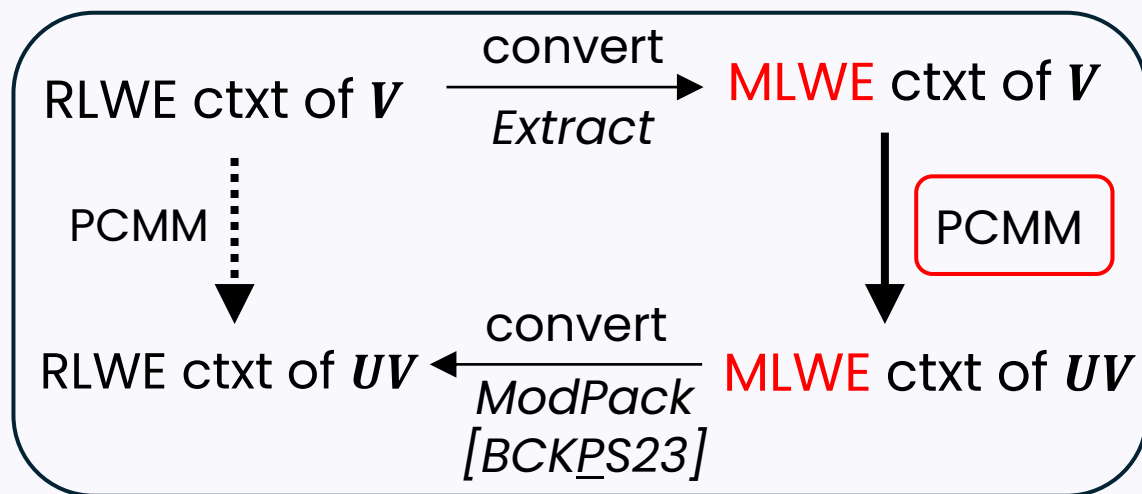
- Latency of 2 PPMMs is not satisfactory



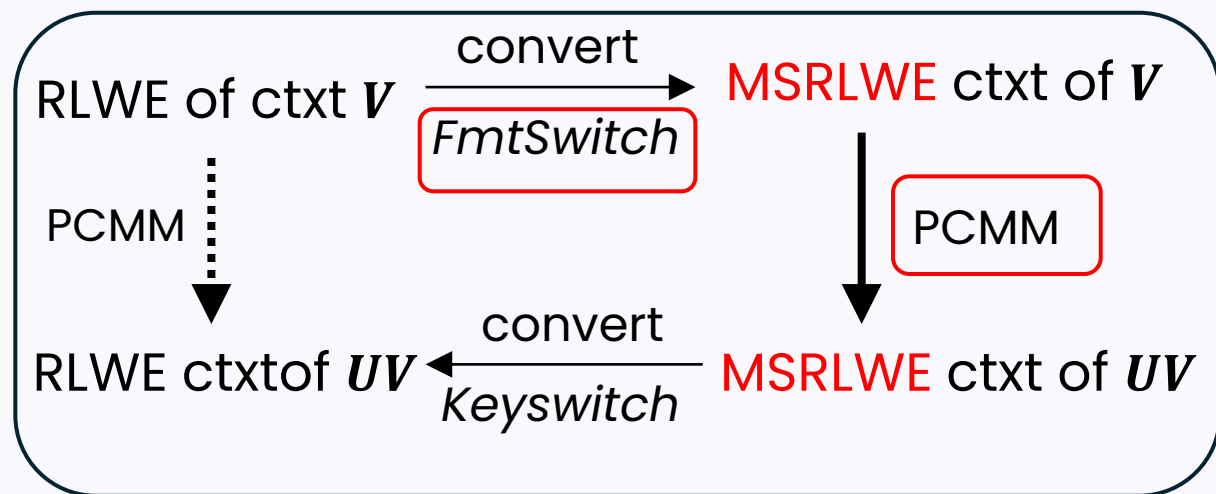
Can we do better than 2 PPMMs?

Different Formats for Various Dimensions

- **Module LWE** and **Multi-secret RLWE** to reduce various dimensional PCMMs to PPMs.
 (MLWE) (MSRLWE)



Small matrices



Large matrices

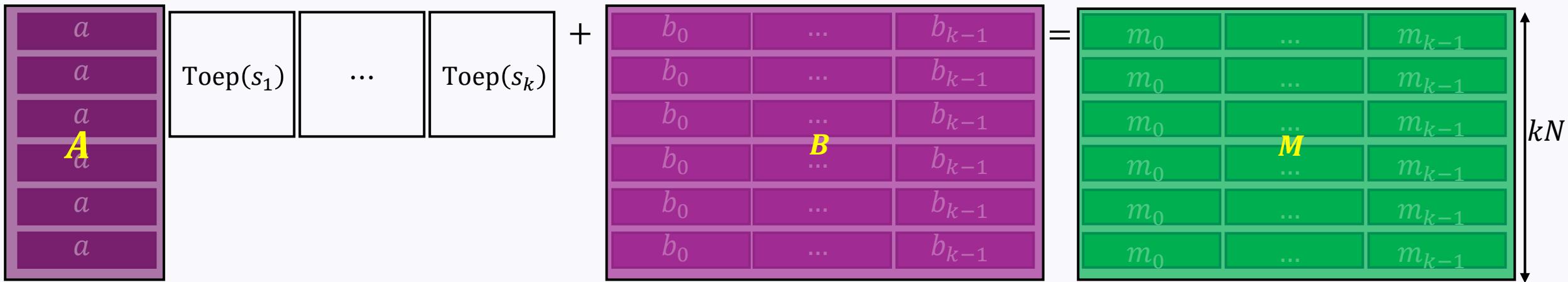
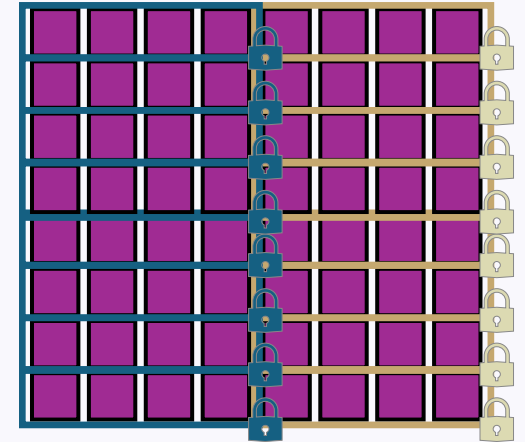
- Format conversions are negligible unless matrices are small.

MSRLWE for Large Matrices

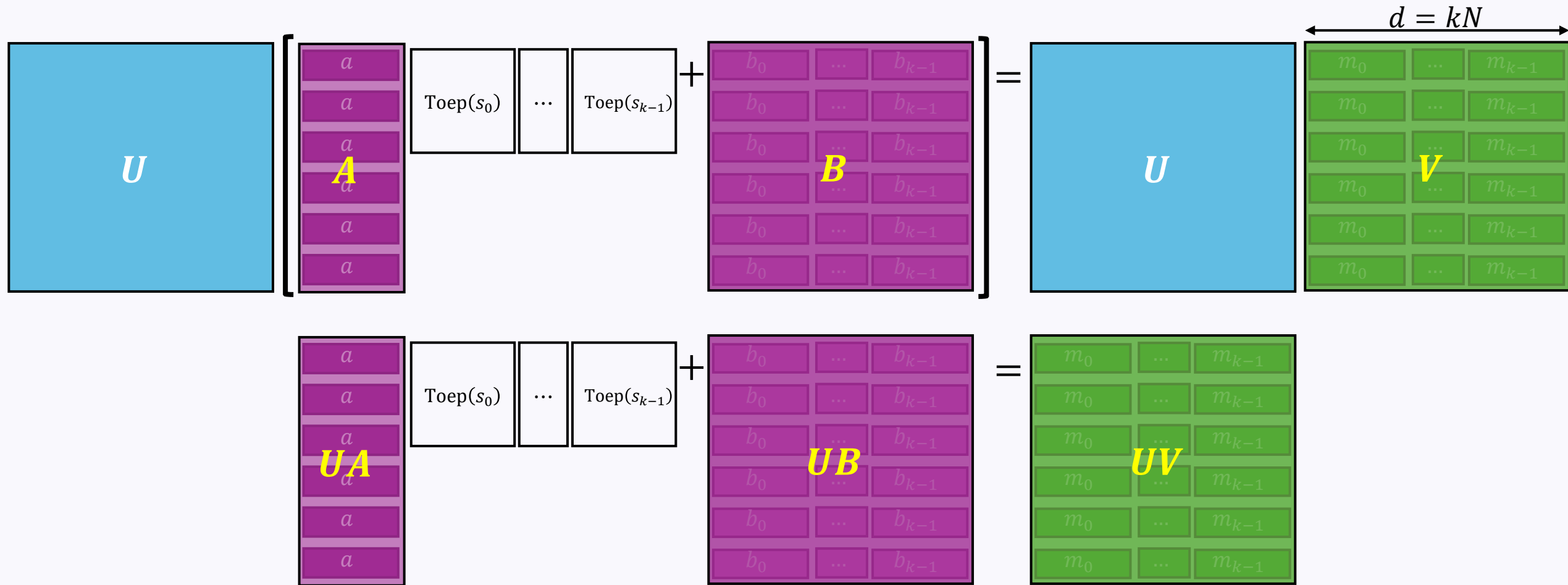
- MSRLWE ciphertexts share a -part rather than secret.

$$as_j + b_j = m_j \quad \forall j \in [k]$$

where s_j is a different secret for each j .

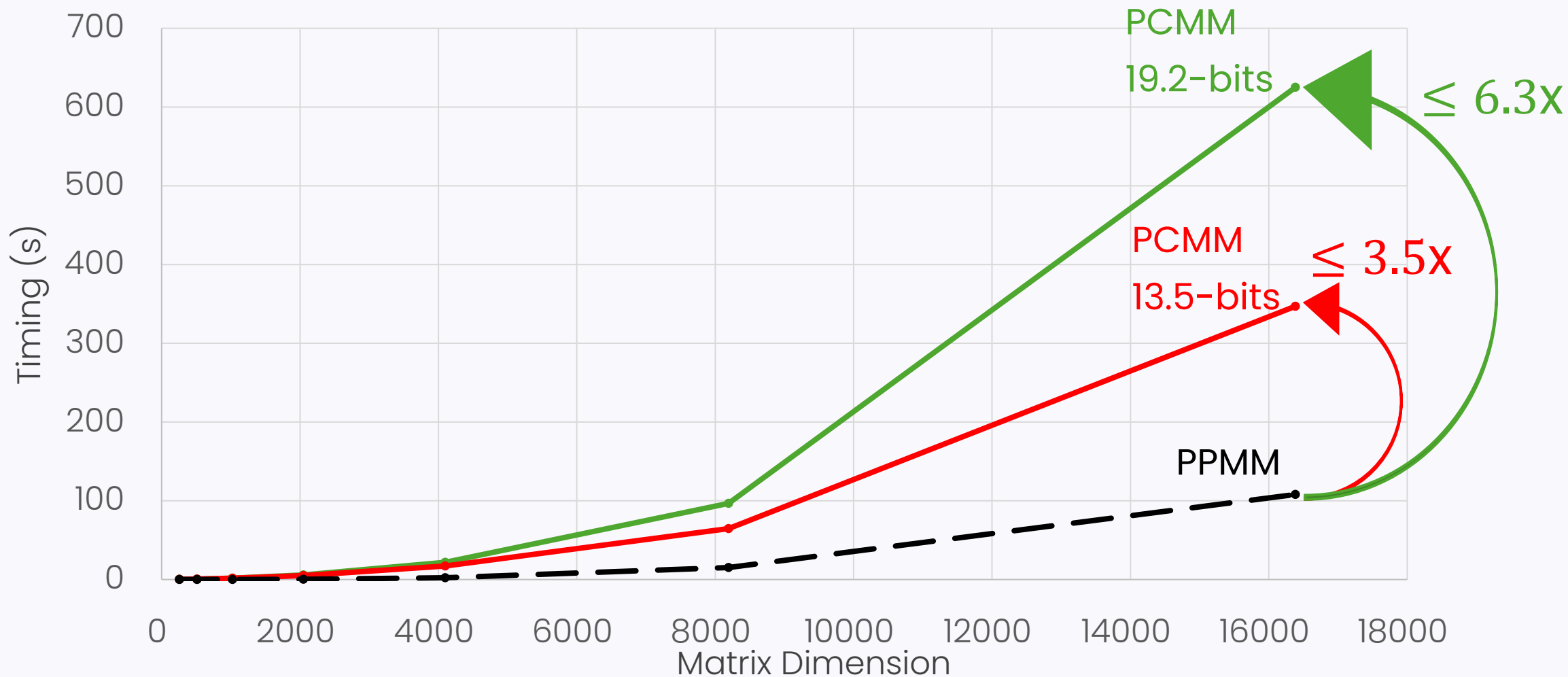


MSRLWE PCMM \leq PPMs



- ❖ $(d > N)$ $d \times d \times d$ MSRLWE PCMM \leq two PPMs modulo Q
- ❖ PPM UA is easier than UV

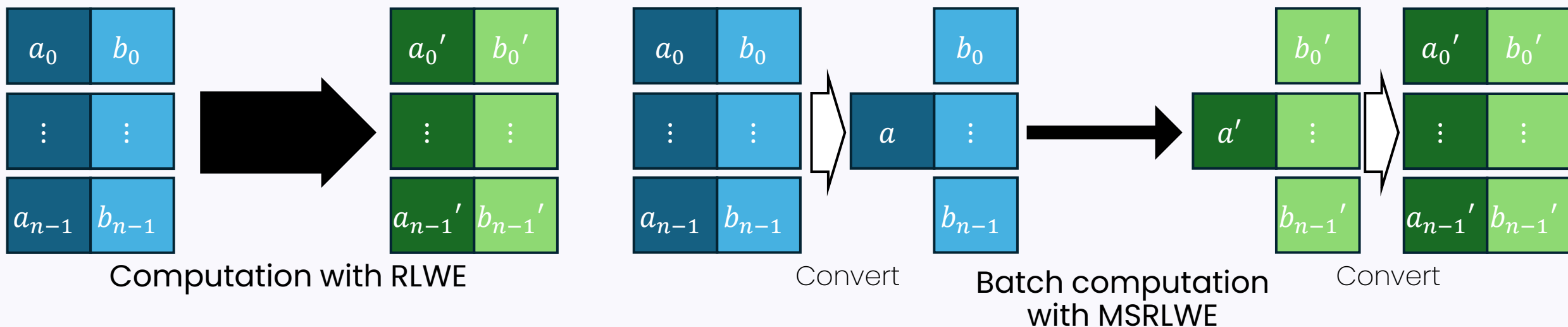
Experimental Results



Intel® Xeon® Gold 6242 CPU at 2.80GHz, single thread
HEaaN library for HE, OpenBLAS for PPM.

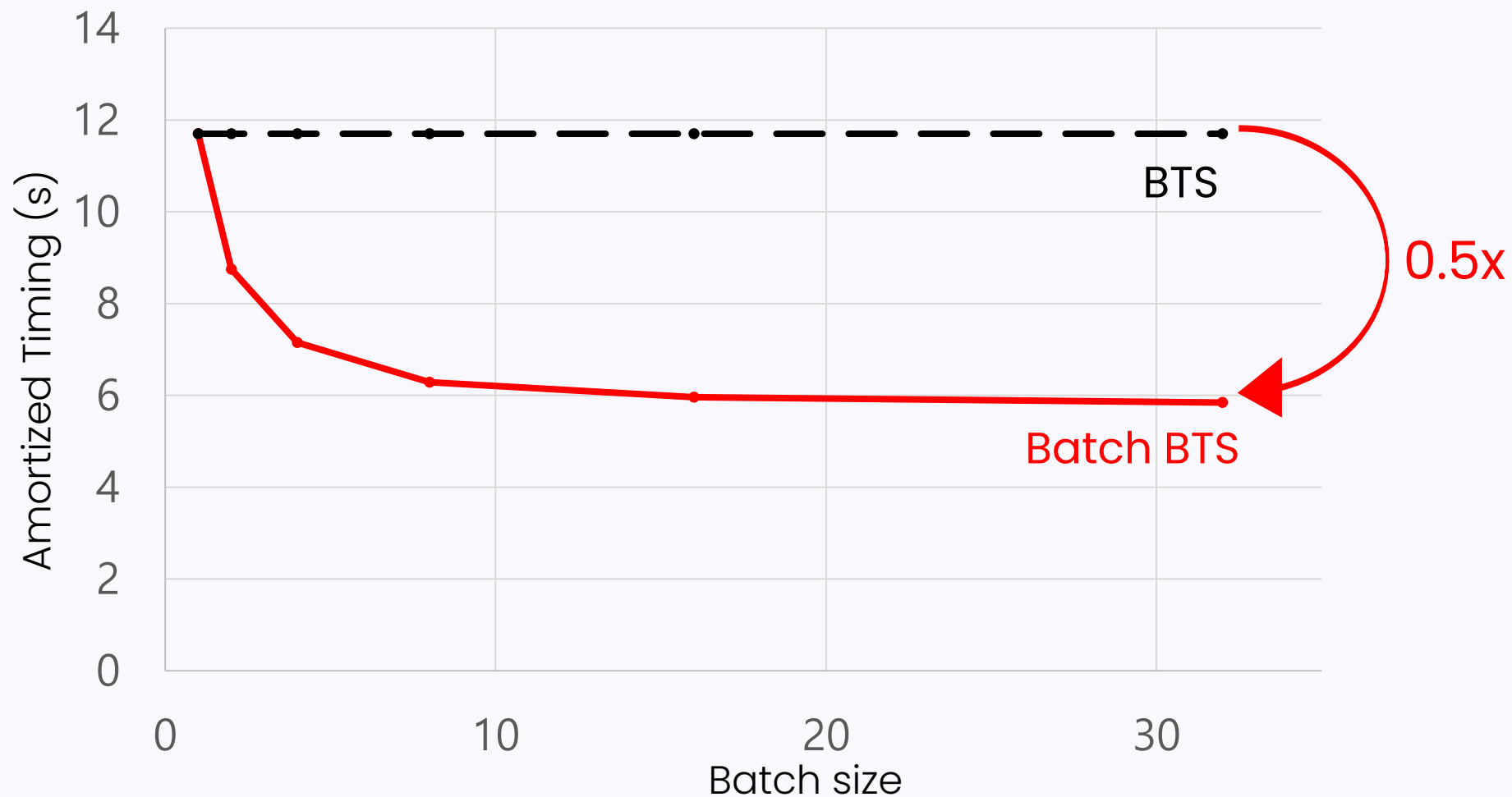
MSRLWE for Batch Computation

- Another question: **how to manage the multiple ciphertexts involved in PCMM?**



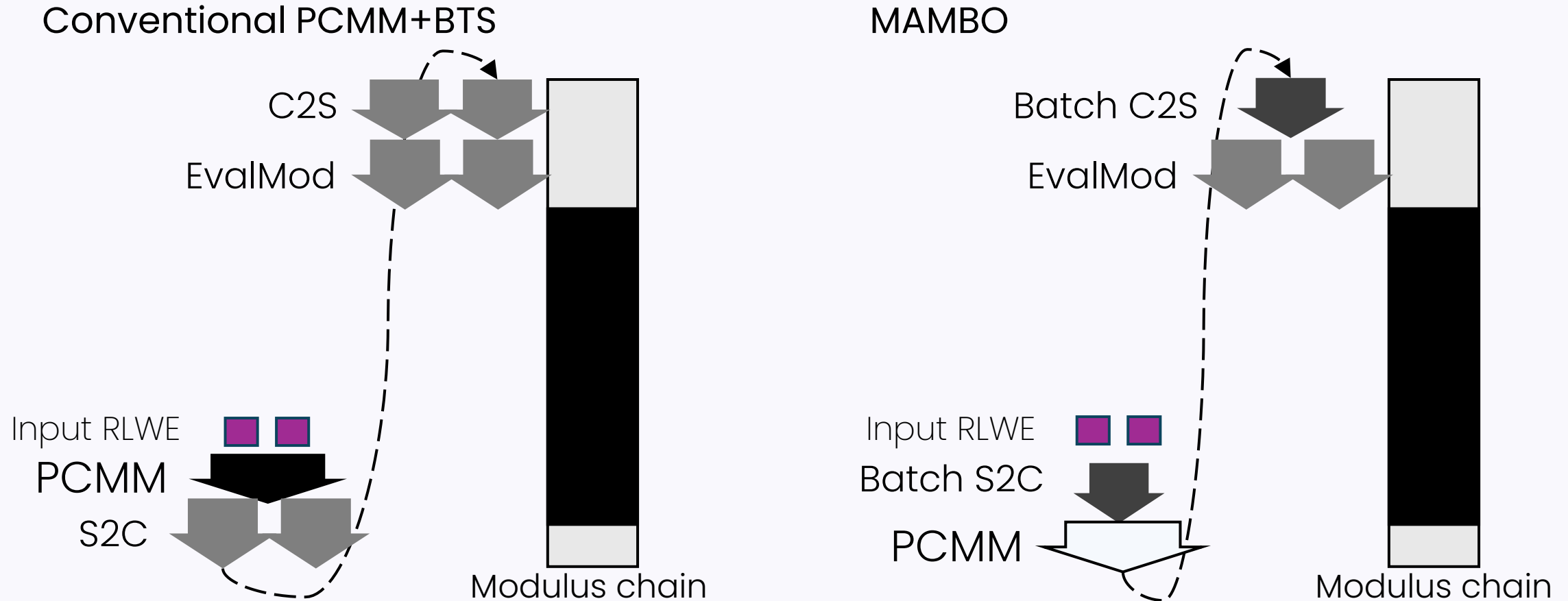
- Native batch operations using MSRLWE
 - // add, // ptxt-ctxt mult, rotate
 - No // ctxt-ctxt mult
- ✓ Batch FHE bootstrapping: batch C2S and S2C during bootstrapping
 - S2C and C2S do not require ctxt-ctxt mult.

Experimental Results



Intel® Xeon® Gold 6242 CPU at 2.80GHz, single thread

MAMBO: Fused PCMM and Batch Bootstrapping



- ✓ For large matrices, MAMBO is 41% faster than current bootstrapping without PCMM.

Wrapping Up!

- **Fast PCMM**
 - Uses various ciphertext formats for various dimensions
 - Exploits efficiency of BLAS libraries
 - Our PCMM is only $\leq 3.5x$ slower than PPM
- **Even faster PCMM with a fixed plaintext matrix (see paper)**
- **Conversion from RLWE formats to MSRLWE formats (see paper)**
- **Batch FHE bootstrapping**
 - MAMBO: Fused PCMM and Batch Bootstrapping
 - **41%** faster than current bootstrapping without PCMM

eprint: 2024/1284

Thank you!