# High-Parallel In-Memory NTT Engine with Hierarchical Structure and Even-Odd Data Mapping

## Bing Li<sup>1</sup>, Huaijun Liu<sup>2</sup>, Yibo Du<sup>3,4</sup>, Ying Wang<sup>3,4</sup>

Institute of Microelectronics, Chinese Academy of Sciences<sup>1</sup> Capital Normal University<sup>2</sup> Institute of Computing Technology, Chinese Academy of Sciences<sup>3</sup>, University of Chinese Academy Sciences<sup>4</sup>







### Outline

- Background and Motivation
- Proposed Method
  - Overview
  - Architecture & Data Mapping
- Evaluation and Results
- Conclusion

# **Fully Homomorphic Encryption**



- Data Security
- Powerful Functionality
- High Computational Overhead

> FHE Review

[Viand A, et al., S&P 2021]

#### **Classic NTT Challenges & Advantages**



**Serially stage computation Complex data transfer μω with n/2 size** 

#### **MVM-Based NTT Challenges & Advantages**



 $\mathbf{I}$   $\boldsymbol{\omega}$  with  $\mathbf{n}^2$  size

High Parallelism and Efficiency

# **Digital in-SRAM Computing**

- High-Efficiency Parallel Computation
- Computing-in-Memory Capacity
- Digital Matrix-Vector Multiplication



#### > All-Digital SRAM-Based ML Edge Applications

[Yu-Der Chih, et al., ISSCC 2021]

### Outline

- Background & Motivation
- Proposed Method
  - Overview
  - Architecture & Data Mapping
- Evaluation and Results
- Conclusion

#### Overview



### Outline

- Background & Motivation
- Proposed Method
  - Overview
  - Architecture & Data Mapping
- Evaluation and Results
- Conclusion

#### **MVM Module-Data Mapping**



#### **MVM Module-Data Mapping**



#### **MVM Module-Computation**



#### **MVM Module-Reduction**



#### **Large-scale NTT Operations**



> Karatsuba multiplication algorithm

 $a \cdot \omega = a_{h} \cdot \omega_{h} \ll 2^{\log_{2}q} + (a_{h} \cdot \omega_{l} + a_{l} \cdot \omega_{h}) \ll 2^{(\log_{2}q)/2} + a_{l} \cdot \omega_{l}$ 

## **Mod Algorithm Optimization**

> Adapt the original Barrett algorithm to the efficient implementation on CIM



### **Mod Algorithm Optimization**



### **Mod Module-Data Mapping**



#### **Mod Module-Computation**



### Outline

- Background & Motivation
- Proposed Method
  - Overview
  - Architecture & Data Mapping
- Evaluation and Results
- Conclusion

# **Evaluation Setup**

Design	Platform	Algorithm	NTT Parameters (n, log <sub>2</sub> q)
HP-CIM(Ours)	6T SRAM	MVM	(32K,32)
<b>BP-NTT</b>	6T SRAM	<b>CT Butterfly</b>	(1024,16)
MeNTT	6T SRAM	<b>CT Butterfly</b>	(32K,32)
<b>RM-NTT</b>	ReRAM	MVM	(1024,16)
CryptoPIM (Baseline)	ReRAM	Butterfly	(32K,32)
HP-CIM Settings			
<b>MVM Module</b>	16 PEs, 32 KB/PE 256 SubArrays/PE, 64×16 SubArray		
<b>MOD Module</b>	2 PEs, 8 KB/PE 128 SubArrays/PE, 8×64 SubArray		



HP-CIM achieves a latency reduction of up to 3.08× compared to the fastest existing CIM-based NTT accelerator, RM-NTT



HP-CIM provides significant energy savings of up to 4.96× over the most energy-efficient prior solution, MeNTT



n=32K,  $log_2q=32$ 

Under large-scale NTT parameter settings, HP-CIM outperforms other designs in terms of latency and energy



➢ HP-CIM reduces execution time by over 2.4× compared to CPU

## Conclusion

- 1. High Parallelism with Hierarchical SRAM Architecture
- Introduced a digital SRAM-based CIM NTT engine, utilizing a hierarchical structure to achieve high parallelism and scalability for large-scale NTT operations.

#### 2. Novel Even-Odd Data Mapping Strategy

• Proposed an even-odd data mapping approach to optimize memory utilization, enabling efficient reuse of intermediate computation results for better scalability.

#### **3. Integrated Mod Computation within CIM Arrays**

• Developed efficient mod operations directly within CIM arrays using SRAM read-write capabilities, eliminating the need for extra peripheral circuits and enhancing area and energy efficiency.

#### 4. Significant Performance and Energy Improvements

• Achieved up to 3.08× faster execution and 4.96× energy savings compared to prior CIMbased designs, validated through extensive comparisons with state-of-the-art methods.

## THANK YOU

# High-Parallel In-Memory NTT Engine with Hierarchical Structure and Even-Odd Data Mapping

### Bing Li<sup>1</sup>, Huaijun Liu<sup>2</sup>, Yibo Du<sup>3,4</sup>, Ying Wang<sup>3,4</sup>

Institute of Microelectronics, Chinese Academy of Sciences<sup>1</sup> Capital Normal University<sup>2</sup> Institute of Computing Technology, Chinese Academy of Sciences<sup>3</sup>, University of Chinese Academy Sciences<sup>4</sup>







### Reference

[1] Gentry C. Fully homomorphic encryption using ideal lattices[C]. Proceedings of the forty-first annual ACM symposium on Theory of computing, Bethesda, Maryland, 2009: 169-178.

[2] Fan J, Vercauteren F. Somewhat Practical Fully Homomorphic Encryption[J]. IACR Cryptology ePrint Archive,2012,2012(2012):144-162.

[3] Kim S, Kim J, Kim M J, et al. Bts: An accelerator for bootstrappable fully homomorphic encryption[C]. Proceedings of the 49th Annual International Symposium on Computer Architecture, New York, 2022: 711-725.

[4] Samardzic N, Feldmann A, Krastev A, et al. F1: A fast and programmable accelerator for fully homomorphic encryption[C]. MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture, Greece, 2021: 238-252.

[5] He Y, Qu S, Lin G, et al. Processing-in-SRAM acceleration for ultra-low power visual 3D perception[C]. Proceedings of the 59th ACM/IEEE Design Automation Conference, San Francisco California, 2022: 295-300.
[6] Li D, Pakala A, Yang K. MeNTT: A compact and efficient processing-in-memory number theoretic transform (NTT) accelerator[J]. IEEE Transactions on Very Large Scale Integration (VLSI) Systems, 2022, 30(5): 579-588.

[7] Albrecht M, Chase M, Chen H, et al. Homomorphic encryption standard[J]. Protecting privacy through homomorphic encryption, 2021: 31-62.