

Exploring the Advantages and Challenges of Fermat NTT in FHE Acceleration

CRYPTO 2024

Andrey Kim^{*†}, Ahmet Can Mert, Anisha Mukherjee, Aikata Aikata, Maxim Deryabin^{*}, Sunmin Kwon^{*}, Hyung Chul Kang^{*}, Sujoy Sinha Roy

IAIK – Graz University of Technology, Austria

^{*} Samsung Advanced Institute of Technology, Suwon, Republic of Korea; [†] Altbridge, Inc.



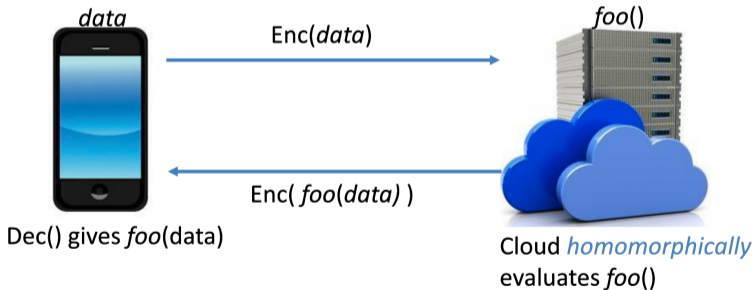
Fully Homomorphic Encryption

Fully Homomorphic Encryption

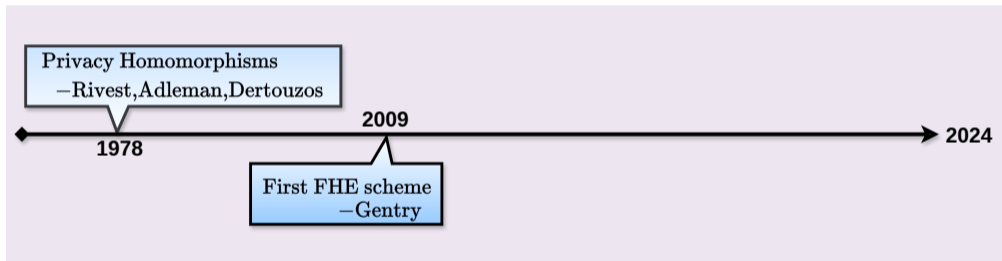
- Fully Homomorphic Encryption (FHE) allows computation on encrypted data.

Fully Homomorphic Encryption

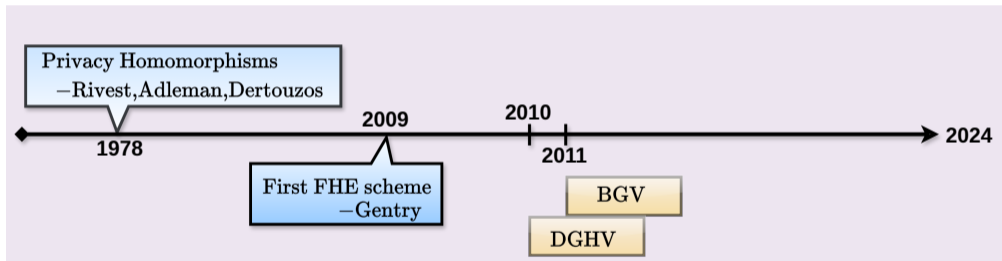
- Fully Homomorphic Encryption (FHE) allows computation on encrypted data.
- Enables secure processing of sensitive data in untrusted environments.



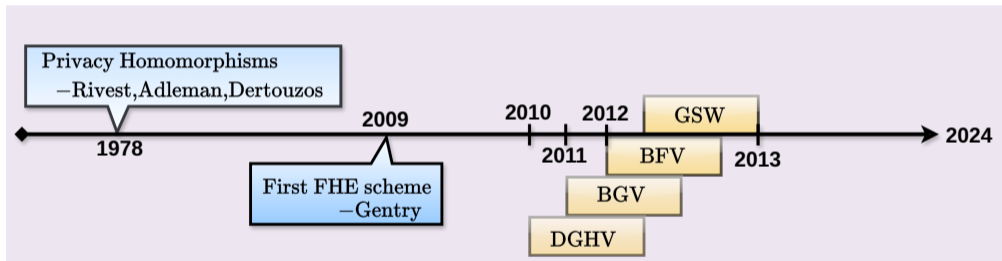
Fully Homomorphic Encryption



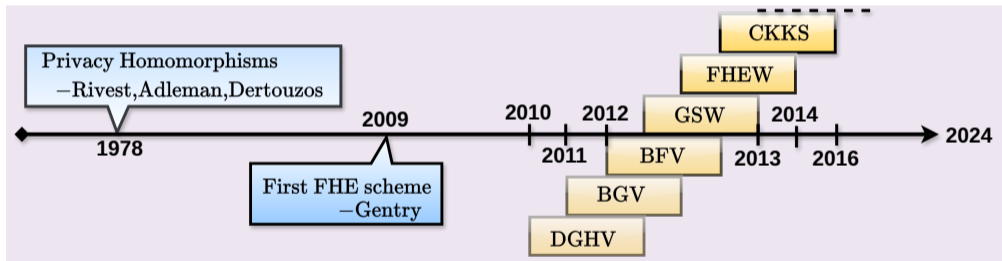
Fully Homomorphic Encryption



Fully Homomorphic Encryption



Fully Homomorphic Encryption



Fully Homomorphic Encryption

- FHE schemes are computationally intensive.

Fully Homomorphic Encryption

- FHE schemes are computationally intensive. 10,000-100,000× plain comp.

Fully Homomorphic Encryption

- FHE schemes are computationally intensive. 10,000-100,000× plain comp.

Solution: Hardware Acceleration.

Fully Homomorphic Encryption

- FHE schemes are computationally intensive. 10,000-100,000× plain comp.

Solution: Hardware Acceleration.

SHARP: A Short-Word Hierarchical Accelerator for Robust and Practical Fully Homomorphic Encryption

Jongmin Kim, Sangpyo Kim, Jaewan Choi
Seoul National University, Seoul National University, Seoul National University

Seoul, Republic of Korea
jongmin.kim@snu.ac.kr

Jaiyoung Park
Seoul National University
Seoul, Republic of Korea
jeff1275@snua.ac.kr

HEAX: An Architecture for Computing on Encrypted Data

M. Sadeq Riazi, Kim Laine, Blake Pelton, Wei Dai
UC San Diego, Microsoft Research, Microsoft, Microsoft Research

Seoul, Republic of Korea
mriazi@ucsd.edu

Microsoft Research
kim.laine@microsoft.com

Microsoft
blakep@microsoft.com

Microsoft Research
wei.dai@microsoft.com

REED: Chiplet-based Accelerator for Fully Homomorphic Encryption

Aikata Aikata¹, Ahmet Can Mert¹, Sunmin Kwon², Maxim Deryabin², Sujoy Sinha Roy¹
Graz University of Technology, Samsung Advanced Institute of Technology, Suwon, Republic of Korea

CraterLake: A Hardware Accelerator for Efficient Unbounded Computation on Encrypted Data

Nikola Samardzic, Axel Feldmann, Aleksandar Krastev
nsamar@csail.mit.edu, axelf@csail.mit.edu, alexalex@csail.mit.edu
Massachusetts Institute of Technology, Massachusetts Institute of Technology, Massachusetts Institute of Technology
Cambridge, MA, USA, Cambridge, MA, USA, Cambridge, MA, USA

BASALISC: Programmable Hardware Accelerator for BGV Fully Homomorphic Encryption

Distribution Statement A: Approved for Public Release, Distribution Unlimited

Robin Geelen^{1*}, Michiel Van Beirendonck^{1*}, Hilder V. L. Pereira¹, Brian Huffman², Tynan McAuley³, Ben Selfridge², Daniel Wagner², Georgios Dimou³, Ingrid Verbauwhede¹, Frederik Vercauteren¹ and David W. Archer²

Medha: Microcoded Hardware Accelerator for computing on Encrypted Data

Ahmet Can Mert¹, Aikata¹, Sunmin Kwon², Youngsam Shin², Donghoon Yoo², Yongwoo Lee² and Sujoy Sinha Roy¹

¹ IAIK, Graz University of Technology, Graz, Austria
{ahmet.mert, aikata, sujoy.sinharoy}@iaik.tugraz.at

² Samsung Advanced Institute of Technology, Suwon, Republic of Korea

FPT: a Fixed-Point Accelerator for Torus Fully Homomorphic Encryption

Michiel Van Beirendonck, Jan-Pieter D'Anvers
michiel.vanbeirendonck@esat.kuleuven.be, janpieter.danvers@esat.kuleuven.be
COSIC, KU Leuven, COSIC, KU Leuven
Leuven, Belgium, Leuven, Belgium

Furkan Turan, Ingrid Verbauwhede
furkan.turan@esat.kuleuven.be, ingrid.verbauwhede@esat.kuleuven.be
COSIC, KU Leuven, COSIC, KU Leuven
Leuven, Belgium, Leuven, Belgium

Fully Homomorphic Encryption

- FHE schemes are computationally intensive. $10,000-100,000\times$ plain comp.

Solution: Hardware Acceleration.

<p>SHARP: A Short-Word Hierarchical Accelerator for Robust and Practical Fully Homomorphic Encryption</p> <p>Jongmin Kim, Sangpyo Kim, Jaewon Choi Seoul National University, Seoul National University, Seoul National University Seoul, Republic of Korea, Seoul, Republic of Korea, Seoul, Republic of Korea jongmin.kim@snu.ac.kr, sangpyo.kim@snu.ac.kr, jaewon.choi@snu.ac.kr</p> <p>HEAX: An Architecture for Computing on Encrypted Data</p> <p>Jaiyoung Park, Donghyun Kim, Jung Ho Ahn, M. Sadegh Riazzi, Kin Lataief, Blake Fellous, Wei Dai Seoul National University, Seoul National University, Seoul National University, UC San Diego, Microsoft Research, Microsoft Research, Microsoft Research jaiyoung.park@snu.ac.kr, donghyun.kim@snu.ac.kr, jahn@seoul.ac.kr, m.riazzi@ucsd.edu, kin.lataief@microsoft.com, blake.fellous@microsoft.com, wei.dai@microsoft.com</p>	<p>REED: Chiplet-based Accelerator for Fully Homomorphic Encryption</p> <p>Aikata Aikata¹, Ahmet Can Mert¹, Sunmin Kwon², Maxim Deryabin², Sujoy Sinha Roy¹ IAIK, Graz University of Technology, Samsung Advanced Institute of Technology, Suwon, Republic of Korea</p>
<p>BASALISC: Programmable Hardware Accelerator for BGV Fully Homomorphic Encryption</p> <p>Distribution Statement A: Approved for Public Release, Distribution Unlimited</p> <p>Robin Geelen^{1*}, Michiel Van Beirendonck^{1*}, Hilder V. L. Pereira¹, Brian Huffman², Tynan McAuley³, Ben Selfridge², Daniel Wagner², Georgios Dimou³, Ingrid Verbauwhede¹, Frederik Vercauteren¹ and David W. Archer²</p>	<p>CraterLake: A Hardware Accelerator for Efficient Fully Homomorphic Encryption</p> <p>$\approx 3,000 - 6,000 \times$ speedup</p> <p>Nikola Sumitrozić, Axel Feldmann, Aleksandar Krstić nsamar@csail.mit.edu, axelf@csail.mit.edu, alexalex@csail.mit.edu Massachusetts Institute of Technology, Massachusetts Institute of Technology, Massachusetts Institute of Technology Cambridge, MA, USA, Cambridge, MA, USA, Cambridge, MA, USA</p>
<p>Medha: Microcoded Hardware Accelerator for computing on Encrypted Data</p> <p>Ahmet Can Mert¹, Aikata¹, Sunmin Kwon², Youngsam Shin², Donghoon Yoo², Yongwoo Lee² and Sujoy Sinha Roy¹</p> <p>¹ IAIK, Graz University of Technology, Graz, Austria {ahmet.mert, aikata, sujoy.sinharoy}@iaik.tugraz.at</p> <p>² Samsung Advanced Institute of Technology, Suwon, Republic of Korea</p>	<p>FPT: a Fixed-Point Accelerator for Torus Fully Homomorphic Encryption</p> <p>Nathan Manoj, Michiel Van Beirendonck, Chris Peikert, Jan-Pieter D'Anvers nmanohar@ibm.com, michiel.vanbeirendonck@esat.kuleuven.be, peikert@umich.edu, janpieter.danvers@esat.kuleuven.be IBM TJ Watson Research Center, COSIC, KU Leuven, University of Michigan, Massachusetts Institute of Technology Yorktown Heights, NY, USA, Menlo Park, CA, USA, Ann Arbor, MI, USA, Cambridge, MA, USA</p>

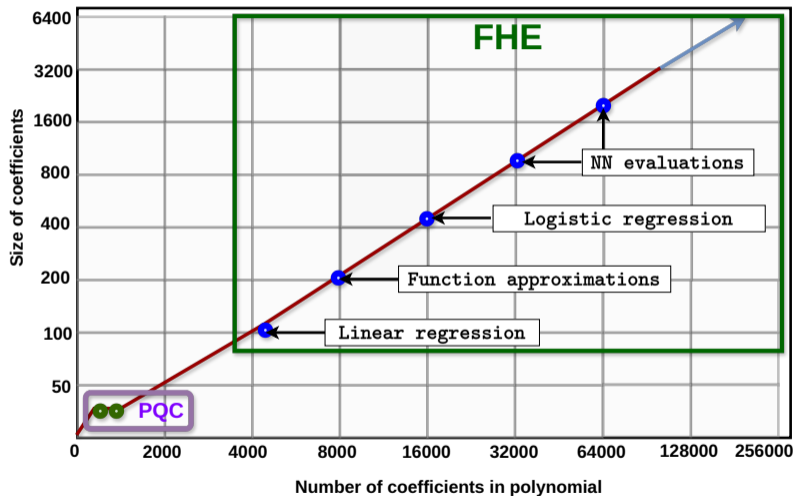
Acceleration Challenges

Challenges in accelerating Homomorphic Encryption

Many polynomial arithmetic operations

- Large degree polynomial arithmetic
- Long integer arithmetic

Challenges with practical realization of Homomorphic Encryption

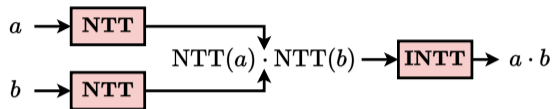


Challenges in accelerating Homomorphic Encryption

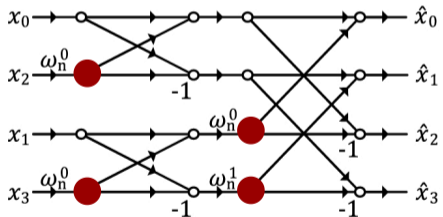
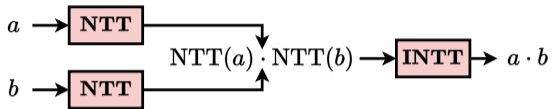
Many polynomial arithmetic operations-**Handled using NTT/INTT unit**

- Large degree polynomial arithmetic
- Long integer arithmetic

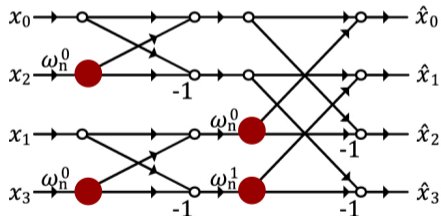
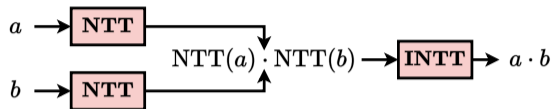
NTT Transformation



NTT Transformation

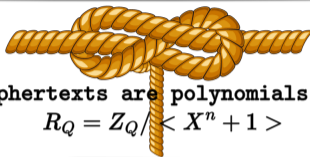


NTT Transformation



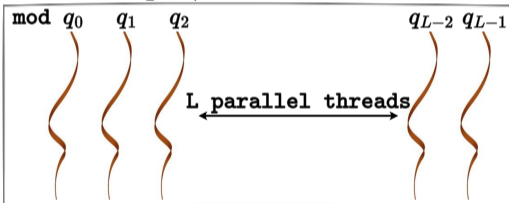
NTT transformation for 1600-bit Q is very expensive.

Residue Number System

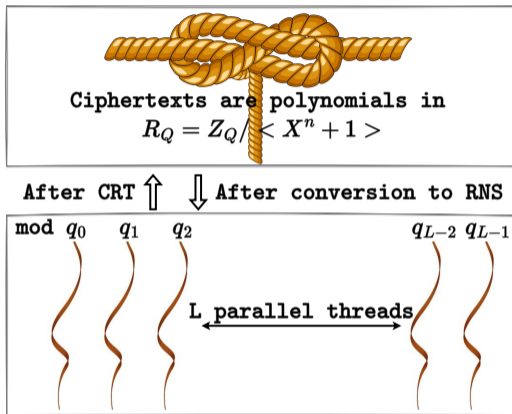


Ciphertexts are polynomials in
 $R_Q = \mathbb{Z}_Q / \langle X^n + 1 \rangle$

After CRT \Uparrow \Downarrow After conversion to RNS



Residue Number System



One NTT transformation is cheap, but it has to support multiple moduli.

REED – Arxiv'23

Components	28nm (mm ²)		7nm (mm ²)	
	1024×64	512×128	1024×64	512×128
REED	74.9	115	24	43.9
REED-PU	58.0	81.0	7.01	9.9
NTT/INTT	38.2	56.8	5.61	7.9
2×MAS	3.1	6.6	0.42	0.76
PRNG	0.15	0.28	0.02	0.04
2×AUT	0.14	0.32	0.02	0.04
Memory	16.1	16.1	1.2	1.2

ARK – ISCA'22

Component	Area (mm ²)	Peak power (W)
4 BConvUs	9.3	18.9
4 NTTUs	57.2	95.2
4 AutoUs	20.6	4.6
8 MADUs	8.9	24.7

F1 – MICRO'21

Component	Area [mm ²]	TDP [W]
NTT FU	2.27	4.80
Automorphism FU	0.58	0.99
Multiply FU	0.25	0.60
Add FU	0.03	0.05
Vector RegFile (512 KB)	0.56	1.67
Compute cluster (NTT, Aut, 2× Mul, 2× Add, RF)	3.97	8.75
Total compute (16 clusters)	63.52	140.0

CraterLake – ISCA'22

Component	Area [mm ²]
CRB FU	158.8
NTT FU	28.1
Automorphism FU	9.0
KSHGen FU	3.3
Multiply FU	2.2
Add FU	0.8
Total FUs (CRB, 2×NTT, Aut, KSHGen, 5×Mul, 5×Add)	240.5

*Can we make modular multiplications in the NTT/INTT units extremely **inexpensive** and ensure NTT **reusability**?*

The **Fermat**-number based Technique

The FNTT method:

Use the Fermat number ($P = 2^K + 1$) as an *auxiliary* modulus before NTT.

The Fermat-number based Technique

The FNTT method:

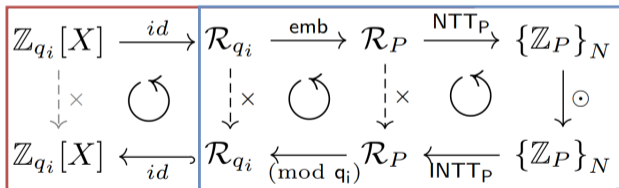
Use the Fermat number ($P = 2^K + 1$) as an *auxiliary* modulus before NTT.

$$\begin{array}{ccccc} \mathbb{Z}_{q_i}[X] & \xrightarrow{id} & \mathcal{R}_{q_i} & \xrightarrow{NTT_{q_i}} & \{\mathbb{Z}_{q_i}\}_N \\ \downarrow \times & \circlearrowleft & \downarrow \times & \circlearrowleft & \downarrow \odot \\ \mathbb{Z}_{q_i}[X] & \xleftarrow{id} & \mathcal{R}_{q_i} & \xleftarrow{INTT_{q_i}} & \{\mathbb{Z}_{q_i}\}_N \end{array}$$

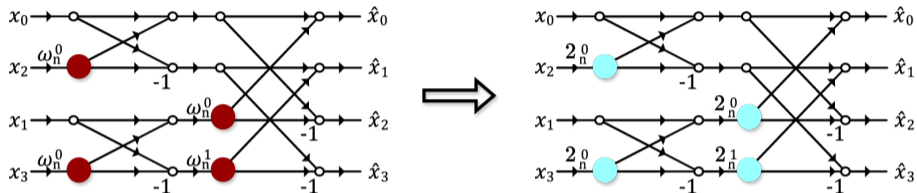
The Fermat-number based Technique

The FNTT method:

Use the Fermat number ($P = 2^K + 1$) as an *auxiliary* modulus before NTT.



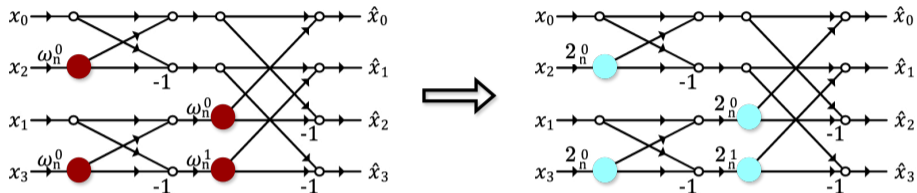
The Fermat-number based Technique



Advantages:

- All modular multiplications during NTT/INTT are transformed into simple shift operations. → Multiplier-less NTT

The Fermat-number based Technique



Advantages:

- All modular multiplications during NTT/INTT are transformed into simple shift operations. → Multiplier-less NTT
- The roots of unity modulo P are powers-of-two. → No storage required

Challenge:

- $P = 2^K + 1$ has K power-of-two twiddle factors.
- With K -th root of unity, we can perform $\frac{K}{2}$ -point negacyclic NTT.

Challenge:

- $P = 2^K + 1$ has K power-of-two twiddle factors.
- With K -th root of unity, we can perform $\frac{K}{2}$ -point negacyclic NTT.

Example: For $N = 2^{16}$, required $K = 2^{16} \rightarrow P = 2^{2^{16}} + 1$. (65,536-bit large modulus)

Challenge:

- $P = 2^K + 1$ K has K power-of-two twiddle factors.
- With K -th root of unity, we can perform $\frac{K}{2}$ -point negacyclic NTT.

Example: For $N = 2^{16}$, required $K = 2^{16} \rightarrow P = 2^{2^{16}} + 1$. (65,536-bit large modulus)
The auxiliary modulus is 1,214× larger than actual moduli q_i (e.g., 54-bit).

Solution: We propose utilizing univariate polynomial to multivariate switch.

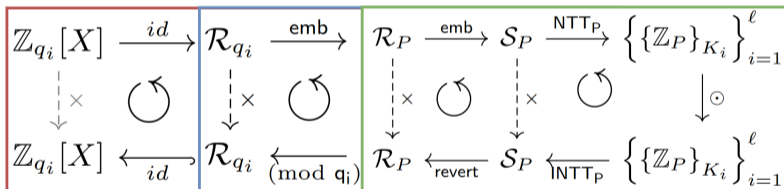
The Fermat-number based Technique

Solution: We propose utilizing univariate polynomial to multivariate switch.

$$\begin{array}{ccccccc}
 \mathbb{Z}_{q_i}[X] & \xrightarrow{id} & \mathcal{R}_{q_i} & \xrightarrow{emb} & \mathcal{R}_P & \xrightarrow{NTT_P} & \{\mathbb{Z}_P\}_N \\
 \downarrow \times & \circlearrowleft & \downarrow \times & \circlearrowleft & \downarrow \times & \circlearrowleft & \downarrow \odot \\
 \mathbb{Z}_{q_i}[X] & \xleftarrow{id} & \mathcal{R}_{q_i} & \xleftarrow{(\text{mod } q_i)} & \mathcal{R}_P & \xleftarrow{INTT_P} & \{\mathbb{Z}_P\}_N
 \end{array}$$

The Fermat-number based Technique

Solution: We propose utilizing univariate polynomial to multivariate switch.



The Fermat-number based Technique

Our technique: For $N = 2^{16}$, we take required $K = 2^7 \rightarrow P = 2^{2^7} + 1$.

The Fermat-number based Technique

Our technique: For $N = 2^{16}$, we take required $K = 2^7 \rightarrow P = 2^{2^7} + 1$.

■ With K -th root of unity, we can perform $\frac{K}{2}$ -point negacyclic NTT.

The Fermat-number based Technique

Our technique: For $N = 2^{16}$, we take required $K = 2^7 \rightarrow P = 2^{2^7} + 1$.

■ With K -th root of unity, we can perform $\frac{K}{2}$ -point negacyclic NTT.

We break $\mathcal{R}_q = \mathbb{Z}_q/(X^{2^{16}} + 1)$ as follows:

$$\Rightarrow (X^{2^{16}} + 1)$$

$$\Rightarrow (X^{2^6} + 1) \times (X^{2^{10}} + 1)$$

$$\Rightarrow (X^{2^6} + 1) \times (X^{2^6} + 1) \times (X^{2^4} + 1)$$

The Fermat-number based Technique

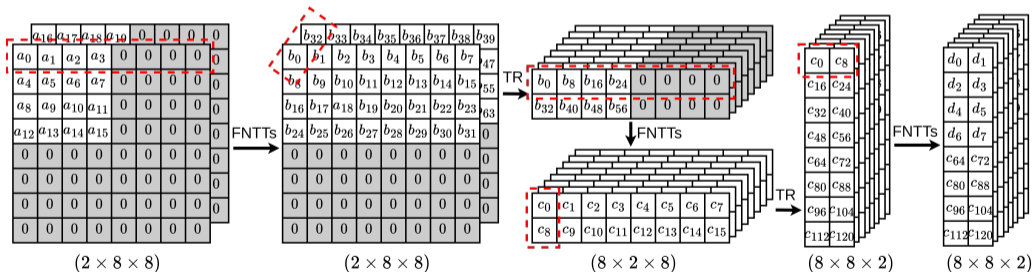
Toy Example: For $N = 2^5$, and $K = 2^4 \rightarrow P = 2^{2^4} + 1$.

$$\mathcal{R}_q = \mathbb{Z}_q / (X^{2^5} + 1) \quad \Rightarrow \quad (X^{2^2} + 1) \times (X^{2^2} + 1) \times (X^{2^1} + 1)$$

The Fermat-number based Technique

Toy Example: For $N = 2^5$, and $K = 2^4 \rightarrow P = 2^{2^4} + 1$.

$$\mathcal{R}_q = \mathbb{Z}_q/(X^{2^5} + 1) \quad \Rightarrow \quad (X^{2^2} + 1) \times (X^{2^2} + 1) \times (X^{2^1} + 1)$$



Implications:

- The auxiliary modulus is only $2\times$ larger than actual moduli q_i (e.g., 54-bit).
- The number of coefficients increase by $4\times$.

Implications:

- The auxiliary modulus is only $2\times$ larger than actual moduli q_i (e.g., 54-bit).
- The number of coefficients increase by $4\times$.

The multivariate has $8\times$ more data overhead than the prior $1,214\times$ overhead.

The **Implementation** Methodology

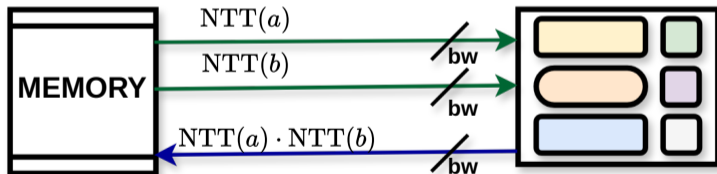
Building Blocks of an FHE accelerator:

- NTT/INTT Unit
- Multiply-and-Accumulate Unit
- Automorphism/Conjugation Unit

The Communication vs Computation trade-off: Data in FNTT form is $8\times$ larger than data in prior NTT forms.

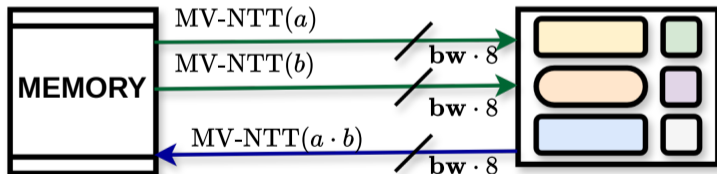
The Implementation Methodology

The Communication vs Computation trade-off: Data in FNTT form is $8\times$ larger than data in prior NTT forms.



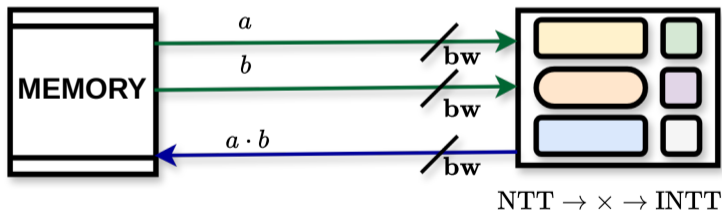
The Implementation Methodology

The Communication vs Computation trade-off: Data in FNTT form is $8\times$ larger than data in prior NTT forms.

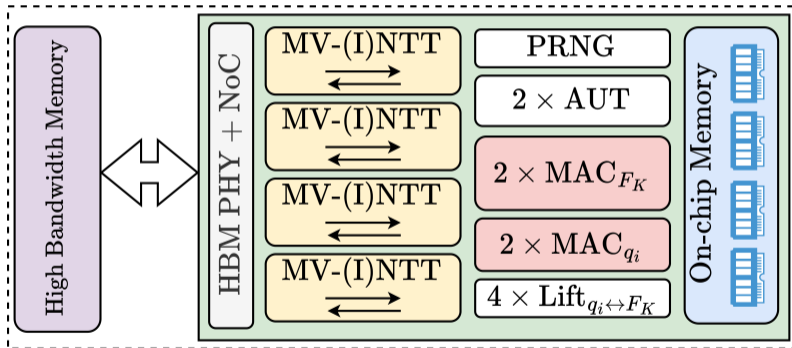


The Implementation Methodology

The Communication vs Computation trade-off: Data in FNTT form is $8\times$ larger than data in prior NTT forms.



To support increased MV-NTT/MV-INTT overhead, we instantiate multiple units.



Implementation Results

The results are for the 28nm ASIC technology, and our design runs at 1.2GHz.

Building Block	Area (mm ²)
MV-NTT/INTT	4×21.5 (vs. 38.2mm ² [REED])
MV-NTT/INTT (mem.)	4×17.4
MAC (for F_K and q_i)	2×2.8
Automorphism	2×0.14
On-chip memory	52
HBM3 PHY+NoC	33.8
Total	250.94

We achieve 1,200× speed-up compared to software implementation.

Future Scope

- **Communication Overhead:** Multiplierless NTTs are cheap, but having to instantiate multiple such units mitigates this advantage.
 - The proposed MV-NTT finds potential in the emerging PIM architectures.

- **Communication Overhead:** Multiplierless NTTs are cheap, but having to instantiate multiple such units mitigates this advantage.
 - The proposed MV-NTT finds potential in the emerging PIM architectures.
- **Computation Overhead:** The amount of arithmetic computation required grows proportional to the reduced multiplications.
 - Implementation of Schönhage-Strassen or Nussbaumer Approach should be explored as they theoretically reduce the operation count.

Exploring the Advantages and Challenges of Fermat NTT in FHE Acceleration

CRYPTO 2024

Andrey Kim^{*†}, Ahmet Can Mert, Anisha Mukherjee, Aikata Aikata, Maxim Deryabin^{*}, Sunmin Kwon^{*}, Hyung Chul Kang^{*}, Sujoy Sinha Roy

IAIK – Graz University of Technology, Austria

^{*} Samsung Advanced Institute of Technology, Suwon, Republic of Korea; [†] Altbridge, Inc.

