# OpenNTT

An Automated Toolchain for Compiling
High-Performance NTT Accelerators in FHE

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy
October 28, 2024

SCIENCE
PASSION
TECHNOLOGY

# Outline

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

**IAIK**
3

# Table of Contents

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

**IAIK**
4

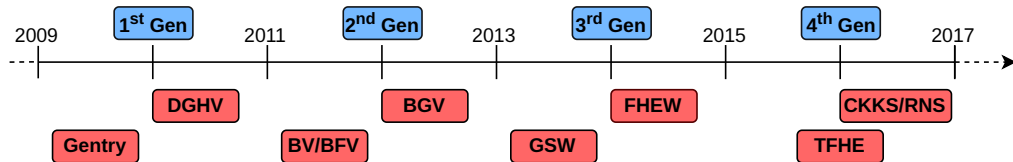# Fully Homomorphic Encryption (FHE)

- Computations over encrypted data
- No information is revealed
- Many different schemes

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

**IAIK**
4

# Fully Homomorphic Encryption (FHE)

- Computations over encrypted data

- No information is revealed

- Many different schemes

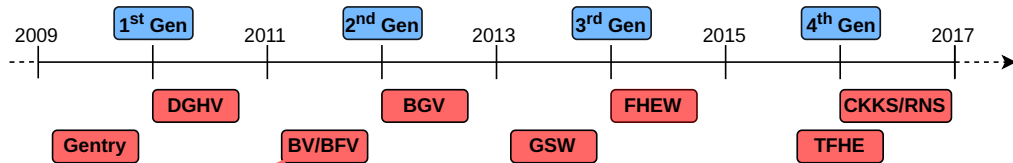Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Fully Homomorphic Encryption (FHE)

- Computations over encrypted data

- No information is revealed

- Many different schemes

# Different FHE Schemes
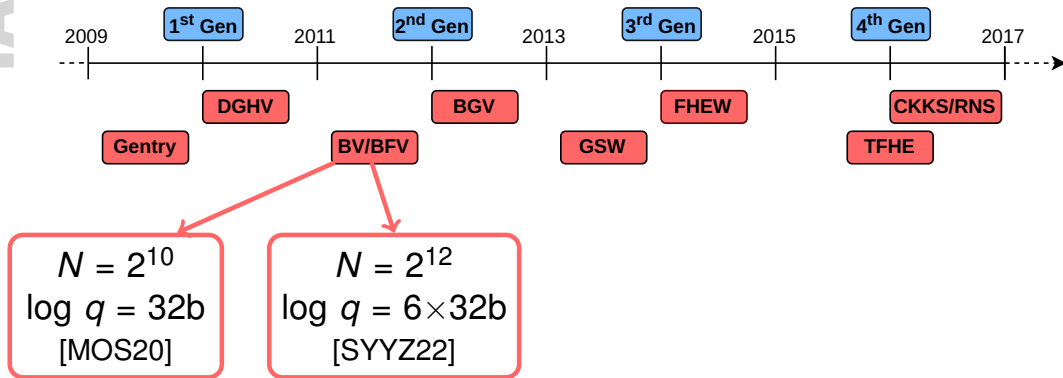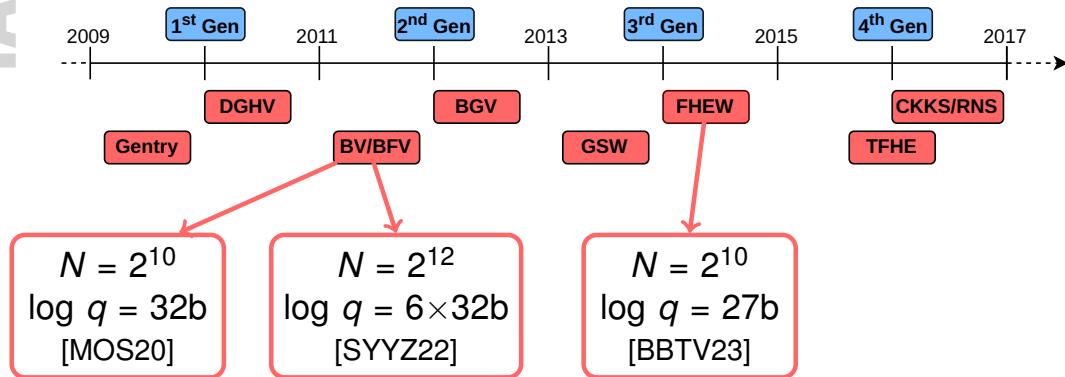
# Different FHE Schemes



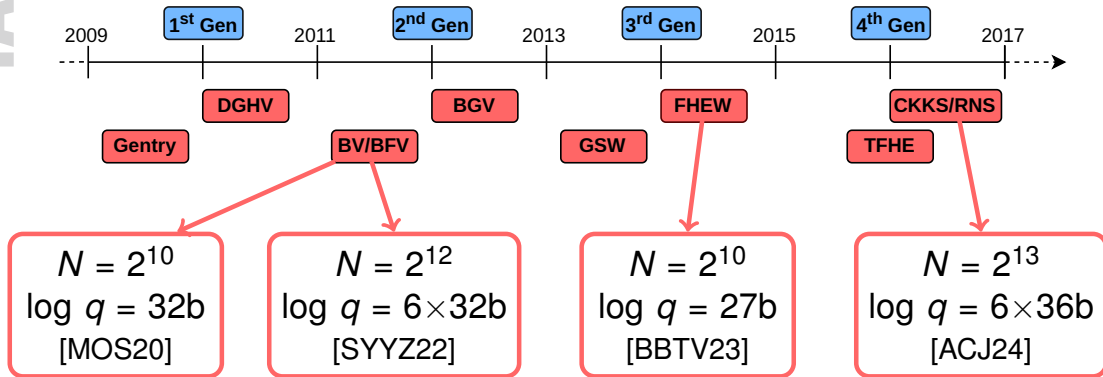$N = 2^{10}$
$\log q = 32b$
[MOS20]

# Different FHE Schemes



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different FHE Schemes



2009    **1st Gen**    2011    **2nd Gen**    2013    **3rd Gen**    2015    **4th Gen**    2017

**DGHV**    **BGV**    **FHEW**    **CKKS/RNS**

**Gentry**    **BV/BFV**    **GSW**    **TFHE**

$N = 2^{10}$
$\log q = 32b$
[MOS20]

$N = 2^{12}$
$\log q = 6 \times 32b$
[SYYZ22]

$N = 2^{10}$
$\log q = 27b$
[BBTV23]

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different FHE Schemes



2009    **1st Gen**    2011    **2nd Gen**    2013    **3rd Gen**    2015    **4th Gen**    2017

DGHV      BGV      FHEW      CKKS/RNS

Gentry      BV/BFV      GSW      TFHE

$N = 2^{10}$
$\log q = 32b$
[MOS20]

$N = 2^{12}$
$\log q = 6 \times 32b$
[SYYZ22]

$N = 2^{10}$
$\log q = 27b$
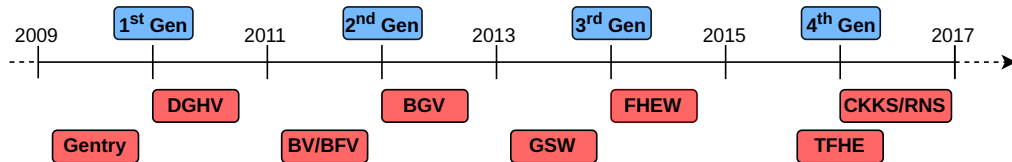[BBTV23]

$N = 2^{13}$
$\log q = 6 \times 36b$
[ACJ24]

**IAIK**
6

# Different FHE Schemes



- Parameters also influenced by:

  - Application scenario
  - Bootstrapping support
  - Target platform

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different FHE Schemes



- **Parameters also influenced by:**

  - Application scenario

  - Bootstrapping support

  - Target platform

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different FHE Schemes



- ■ Parameters also influenced by:

  - ■ Application scenario

  - ■ Bootstrapping support

  - ■ Target platform

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# FHE: Large Variety in Parameters

- Different schemes $\rightarrow$ different parameters
- Influences building-block level:
  - Polynomial multiplication
  - Number Theoretic Transformation

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# FHE: Large Variety in Parameters

- Different schemes $\rightarrow$ different parameters

- Influences building-block level:

  - Polynomial multiplication
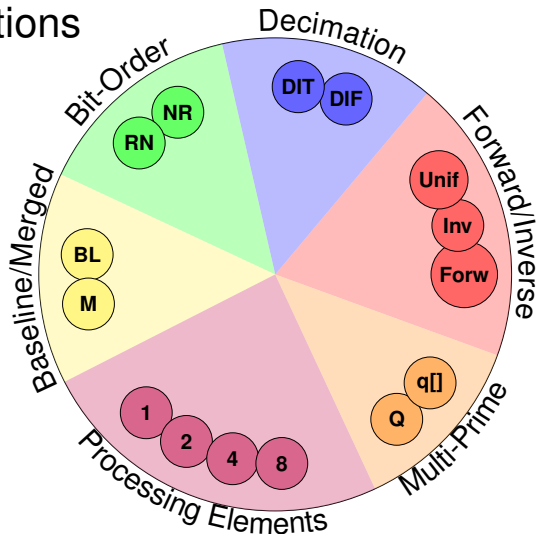  - Number Theoretic Transformation

**IAIK**

8

# The Number Theoretic Transformation

- Computational bottleneck: NTT

  → >70% of computation time [KKK+22]

  → Hardware acceleration

## Many configurations: Challenges

? How can we reduce the hardware design effort?

? How can we enhance flexibility?

♀ Hardware design tools!

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# The Number Theoretic Transformation

- Computational bottleneck: NTT

  → >70% of computation time [KKK+22]
  → Hardware acceleration

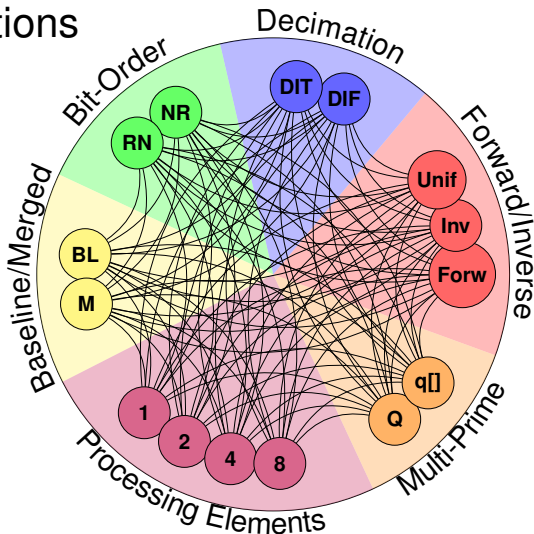## Many configurations: Challenges

? How can we reduce the hardware design effort?

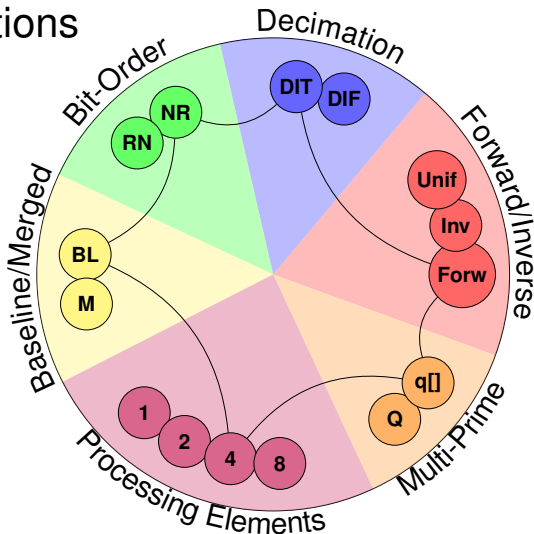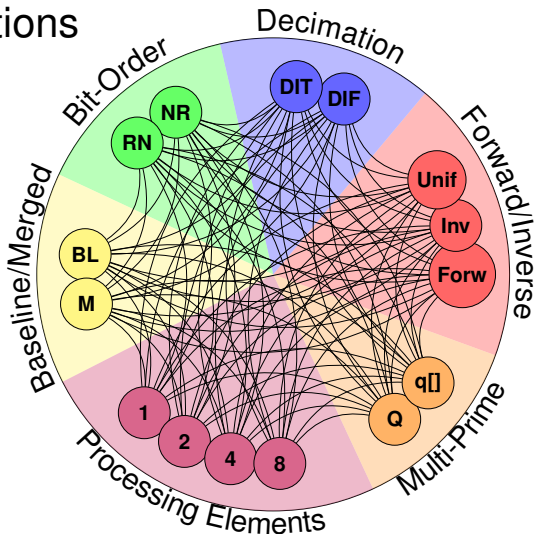? How can we enhance flexibility?

💡 Hardware design tools!

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# NTT Configurations



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# NTT Configurations



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# NTT Configurations



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# NTT Configurations



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# NTT Configurations



**NTT design tools
are challenging**

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# NTT: Versatile Applications

- Not only for FHE
- Post-Quantum Cryptography
- Zero-Knowledge Proof Systems



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

**IAIK**

14

# NTT: Versatile Applications

- Not only for FHE

- Post-Quantum Cryptography

- Zero-Knowledge Proof Systems



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# NTT: Versatile Applications

- Not only for FHE

- Post-Quantum Cryptography

- Zero-Knowledge Proof Systems



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Table of Contents

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different NTT Approaches



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different NTT Approaches

- Butterfly operation

# Different NTT Approaches

- Butterfly operation



- Multiple processing elements (PE)

  → Higher performance

# Different NTT Approaches

- **Iterative NTT**
  - Multi-PE
  - High-radix PE



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different NTT Approaches

- Iterative NTT

    - Multi-PE

    - High-radix PE



Input

Output

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different NTT Approaches

- **Iterative NTT**

  - **Multi-PE**

  - High-radix PE



Input

Output

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different NTT Approaches

- Iterative NTT

    - Multi-PE

    - High-radix PE



Input

Output

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different NTT Approaches

- **Iterative NTT**

    - Multi-PE

    - High-radix PE



Input

Output

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different NTT Approaches

- **Iterative NTT**
  - Multi-PE
  - High-radix PE



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different NTT Approaches

- Iterative NTT

  - Multi-PE

  - High-radix PE

- Pipelined NTT

  - Single-path delay feedback (SDF)

  - Multi-path delay commutator (MDC)

# Different NTT Approaches

- Iterative NTT

  - Multi-PE

  - High-radix PE

- Pipelined NTT

  - Single-path delay feedback (SDF)

  - Multi-path delay commutator (MDC)



[HMR24]

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Different NTT Approaches

- **Iterative NTT**

  - **Multi-PE**

  - High-radix PE

- **Pipelined NTT**

  - Single-path delay feedback (SDF)

  - Multi-path delay commutator (MDC)

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Iterative Multi-PE NTT

- Twiddle factor management (TF)

| Stored TF | On-the-fly generated TF |
|---|---|
| + Less logic | + Less memory |
| − Large memories | − More logic |
| → Simpler to implement | → Better for FHE (Multiple primes, large polys) |

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Iterative Multi-PE NTT

- Twiddle factor management (TF)

| Stored TF | On-the-fly generated TF |
|---|---|
| ➕ Less logic | ➕ Less memory |
| ➖ Large memories | ➖ More logic |
| → Simpler to implement | → Better for FHE (Multiple primes, large polys) |

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Iterative Multi-PE NTT

- Twiddle factor management (TF)

| Stored TF | On-the-fly generated TF |
|---|---|
| ➕ Less logic | ➕ Less memory |
| ➖ Large memories | ➖ More logic |
| ➔ Simpler to implement | ➔ Better for FHE (Multiple primes, large polys) |

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Iterative Multi-PE NTT

🧩 On-the-fly twiddle factor generation

→ Linear twiddle factor order: $\omega^1, \omega^2, \omega^3, \ldots$

🧩 Conflict-free memory accesses

→ Dedicated execution flow

## OpenNTT

- Combines 🧩 and 🧩 generically
- For all relevant configurations

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Iterative Multi-PE NTT

- On-the-fly twiddle factor generation

    → Linear twiddle factor order: $\omega^1, \omega^2, \omega^3, \ldots$

- Conflict-free memory accesses

    → Dedicated execution flow

**OpenNTT**

- Combines 🧩 and 🧩 generically
- For all relevant configurations

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Iterative Multi-PE NTT

🧩 On-the-fly twiddle factor generation

→ Linear twiddle factor order: $\omega^1, \omega^2, \omega^3, \ldots$

🧩 Conflict-free memory accesses

→ Dedicated execution flow

## OpenNTT

- Combines 🧩 and 🧩 generically
- For all relevant configurations

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Table of Contents

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Generic Processing Orders

**IAIK**

- **Derive universal memory access constraints**

  - For each stage
  - For every NTT type
  - For different number of PE

- Combine with efficient twiddle factor order

  - Ensure *somewhat* linear generation order

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Generic Processing Orders

- Derive universal memory access constraints

    - For each stage

    - For every NTT type

    - For different number of PE

- Combine with efficient twiddle factor order

    - Ensure *somewhat* linear generation order

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Generic Processing Orders

- Derive universal memory access constraints

    - For each stage

    - For every NTT type

    - For different number of PE

- Combine with efficient twiddle factor order

    - Ensure *somewhat* linear generation order

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Somewhat Linear Generation Order

- Some NTT configurations need somewhat linear generation:

  → Due to memory constraints

  $\omega^0$, $\omega^8$, $\omega^1$, $\omega^9$, $\omega^2$, $\omega^{10}$, ...

  $\omega^i$, $\omega^j$, $\omega^{i+c}$, $\omega^{j+c}$, ...

- Changes across stages

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Somewhat Linear Generation Order

- Some NTT configurations need somewhat linear generation:

  → Due to memory constraints

  $\omega^0$, $\omega^8$, $\omega^1$, $\omega^9$, $\omega^2$, $\omega^{10}$, ...

  $\omega^i$, $\omega^j$, $\omega^{i+c}$, $\omega^{j+c}$, ...

- Changes across stages

# Somewhat Linear Generation Order

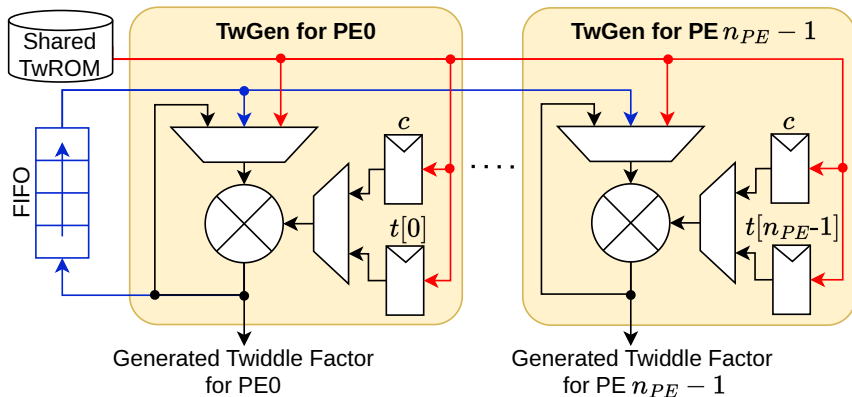- Some NTT configurations need somewhat linear generation:

  → Due to memory constraints

  $\omega^0, \omega^8, \omega^1, \omega^9, \omega^2, \omega^{10}, \ldots$

  $\omega^i, \omega^j, \omega^{i+c}, \omega^{j+c}, \ldots$

- Changes across stages

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Efficient TF Generation Module



Generated Twiddle Factor
for PE0

Generated Twiddle Factor
for PE $n_{PE} - 1$

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Flexible Hardware Architecture



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# The OpenNTT Toolchain

- **Overtakes all steps**

- Provides testing functionality



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# The OpenNTT Toolchain

- Overtakes all steps

- Provides testing functionality



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# The OpenNTT Toolchain

- **Supported parameters:**

  - Polynomial size

  - Prime size

  - Number of primes

  - Baseline or NWC NTT

  - Decimation method

  - Coefficient orders

  - Forward/inverse/unified NTT

  - Number of PE

  - NTT only or with arithmetic

  - Memory depth

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Table of Contents

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Memory vs. Logic Consumption

- Twiddle factor generation trades memory for logic



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Achieved Frequency

- Benchmarked on same FPGA



Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Comparison Benchmarks

- NTT generation tools:

  → Prior tools use stored twiddle factors
  - 1.3× to 2.7× speedup
  - 1.8× improved ATP

- Dedicated designs:

  - Comparable performance

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Comparison Benchmarks

- NTT generation tools:

  → Prior tools use stored twiddle factors
  - $1.3\times$ to $2.7\times$ speedup
  - $1.8\times$ improved ATP
- Dedicated designs:
  - Comparable performance

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Comparison Benchmarks

- NTT generation tools:

  → Prior tools use stored twiddle factors
  - $1.3\times$ to $2.7\times$ speedup
  - $1.8\times$ improved ATP

- Dedicated designs:

  - Comparable performance

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Comparison Benchmarks

- NTT generation tools:

  → Prior tools use stored twiddle factors
  - $1.3\times$ to $2.7\times$ speedup
  - $1.8\times$ improved ATP

- Dedicated designs:

  - Comparable performance

# Table of Contents

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# Conclusion and Takeaways

- **Generic tool for NTT hardware accelerator design**
- Open-sourced our tool
  - → Contribute to faster NTT design times
  - → Supports research in the field
  - → Relevant for industrial and academic domains

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

**IAIK**
33

# Conclusion and Takeaways

- Generic tool for NTT hardware accelerator design

- Open-sourced our tool

  → Contribute to faster NTT design times

  → Supports research in the field

  → Relevant for industrial and academic domains

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

# OpenNTT

An Automated Toolchain for Compiling
High-Performance NTT Accelerators in FHE

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy

October 28, 2024

References

[MOS20] A. C. Mert *et al.*, "Design and Implementation of Encryption/Decryption Architectures for BFV Homomorphic Encryption Scheme," in IEEE TVLSI, 2020

[SYYZ22] Y. Su *et al.*, "ReMCA: A Reconfigurable Multi-Core Architecture for Full RNS Variant of BFV Homomorphic Evaluation," in TCAS-I, 2022

[BBTV23] J. Bertels *et al.*, "Hardware Acceleration of FHEW," in DDECS, 2023

[ACJ24] R. Agrawal *et al.*, "HEAP: A Fully Homomorphic Encryption Accelerator with Parallelized Bootstrapping," in ISCA, 2024

[HMR24] F. Hirner *et al.*, "Proteus: A Pipelined NTT Architecture Generator," in TVLSI, 2024

[HTLW23] Xiao Hu *et al.*, "Ac-pm: an area-efficient and configurable polynomial multiplier for lattice based cryptography," in TCAS-I, 2023

[LKMS23] Si-Huang Liu *et al.*, "An area-efficient, conflict-free, and configurable architecture for accelerating ntt/intt," in TVLSI, 2023

[MKO+20] A. C. Mert *et al.*, "A flexible and scalable ntt hardware: applications from homomorphically encrypted deep learning to post-quantum cryptography," in DATE, 2020

[MRW+23] J. Mu *et al.*, "Scalable and conflict-free ntt hardware accelerator design: methodology, proof, and implementation," in TCAD, 2023

[DKJL22] P. Duong-Ngoc *et al.*, "Area-efficient number theoretic transform architecture for homomorphic encryption," in TCAS-I, 2022

[KKK+22] S. Kim *et al.*, "BTS: an accelerator for bootstrappable fully homomorphic encryption," in ISCA, 2022

Florian Krieger, Florian Hirner, Ahmet Can Mert, Sujoy Sinha Roy