



清華大學
Tsinghua University

CFNTT: Scalable Radix-2/4 NTT Multiplication Architecture with an Efficient Conflict-free Memory Mapping Scheme

Xiangren Chen, Bohan Yang, Shouyi Yin, Shaojun Wei and Leibo Liu*

School of Integrated Circuits, Tsinghua University, China.

CHES 2022, Issue 1



Outlines

- **Introduction**
- **Optimized Radix-4 NTT/INTT Algorithm**
- **Conflict-free Memory Mapping Scheme**
- **Hardware Architecture of Radix-2/4 NTT**
- **Implementation Result and Comparison**



Outlines

- **Introduction**
- Optimized Radix-4 NTT/INTT Algorithm
- Conflict-free Memory Mapping Scheme
- Hardware Architecture of Radix-2/4 NTT
- Implementation Result and Comparison

NTT Related Cryptographic Scheme

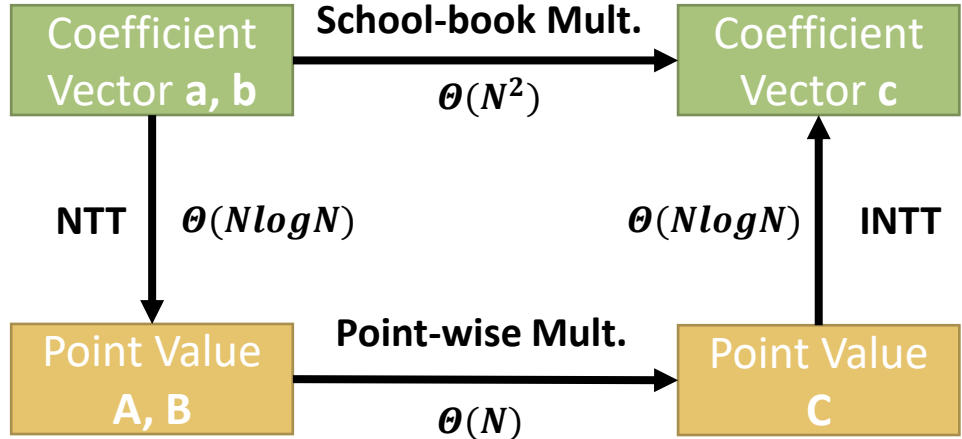
$$\boxed{A} \times \boxed{s} + \boxed{e} = \boxed{b} \quad \text{Key Generation}$$

$$\boxed{A^T} \times \boxed{r} + \boxed{e'} = \boxed{u} \quad \text{Encryption}$$

$$\boxed{b^T} \times \boxed{r} + \boxed{e''} + \boxed{m} = \boxed{v}$$

$$\boxed{v} - \boxed{s^T} \times \boxed{u} \approx \boxed{m} \quad \text{Decryption}$$

School-book Algorithm: $\Theta(N^2)$
 Karatsuba Algorithm: $\Theta(N^{\log_2 3})$
Number Theoretical Transformation: $\Theta(N \log N)$

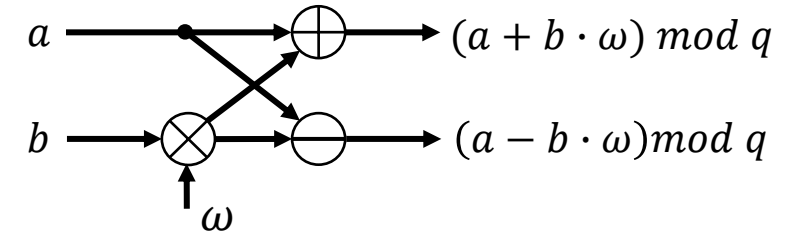


Polynomial computation in MLWE-based scheme e.g. Kyber

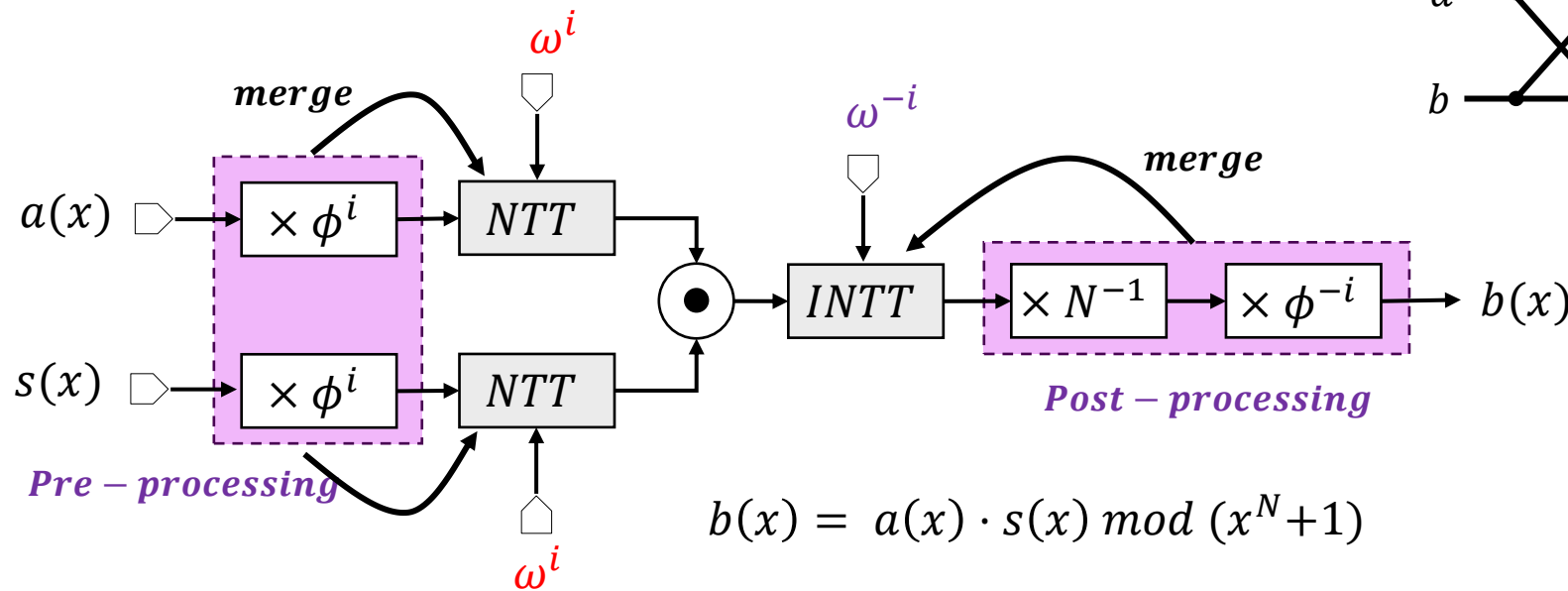
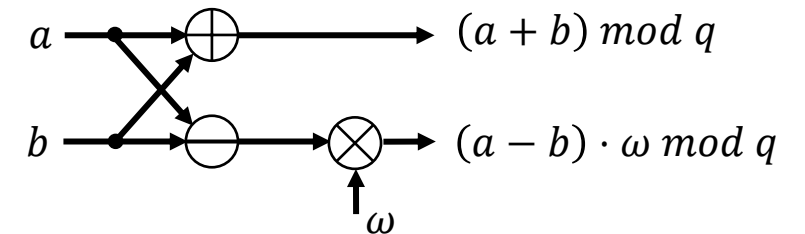
NTT Multiplication over the Ring

- $\mathbb{R}_q = \mathbb{Z}_q[x] / \langle f(x) \rangle \quad q \equiv 1 \pmod{2N}$
- $f(x) = x^N + 1 \Rightarrow$ Negative Wrapped Convolution (NWC)
- $\omega : N$ -th roots of unity
- $\phi : 2N$ -th roots of unity

Radix-2 Cooley-Tukey BFU

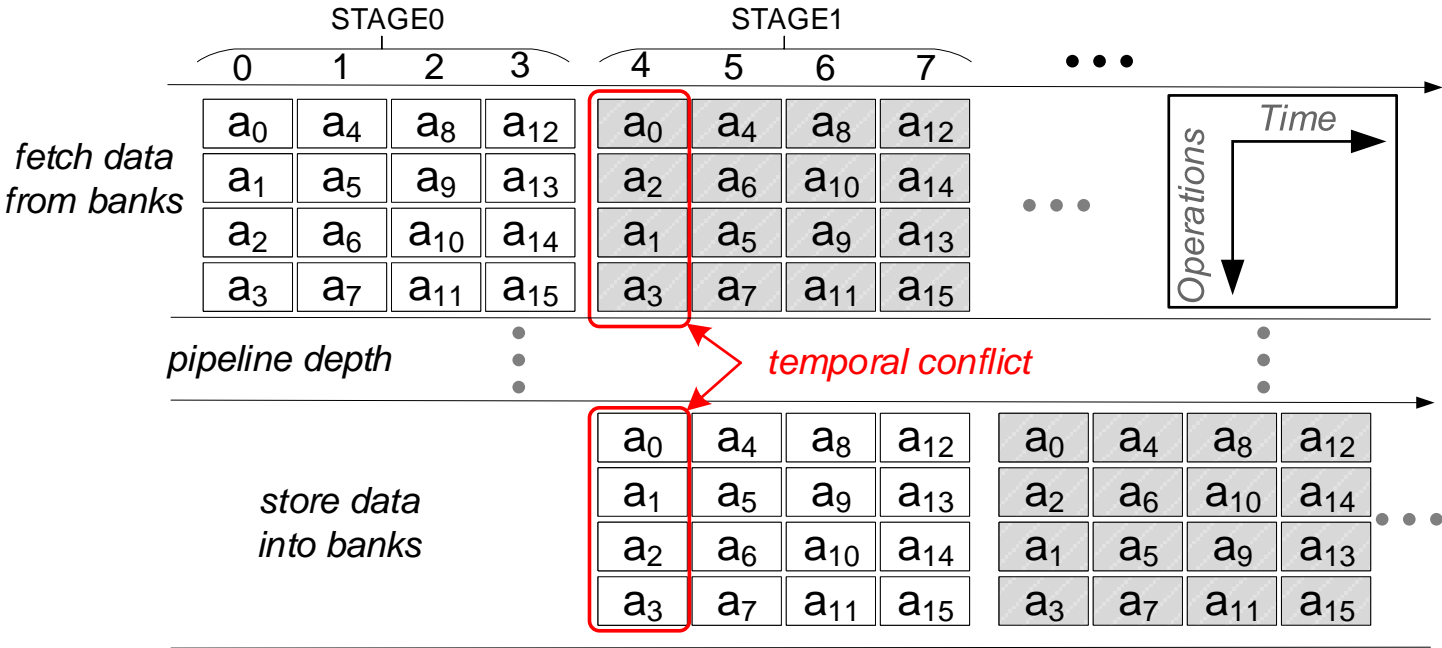


Radix-2 Gentleman-Sande BFU



Temporal Conflict

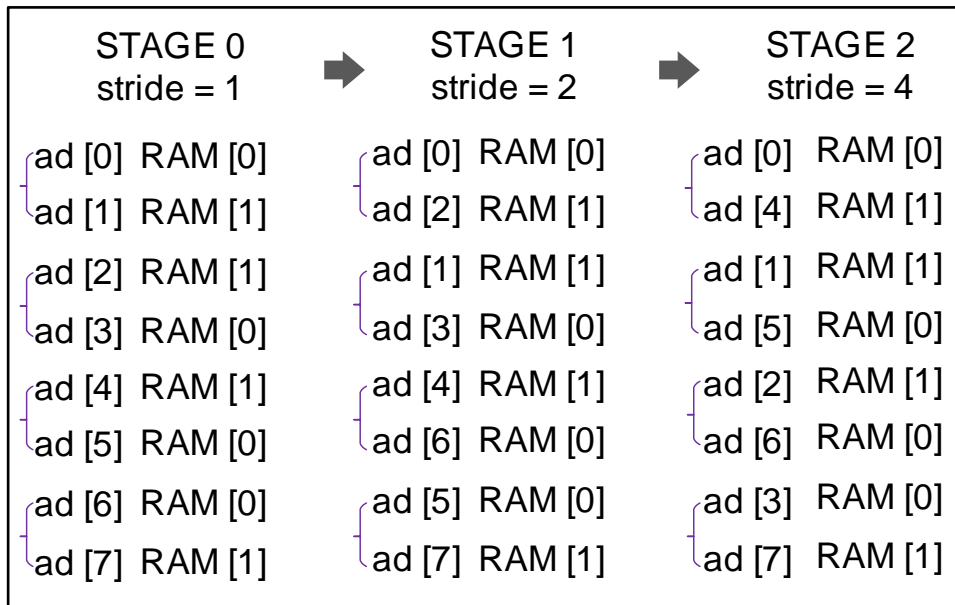
- Read After Write (RAW) pipelined hazard
 - $N/4d \geq \text{pipeline depth}$
 - More stringent when considering the higher radix NTT



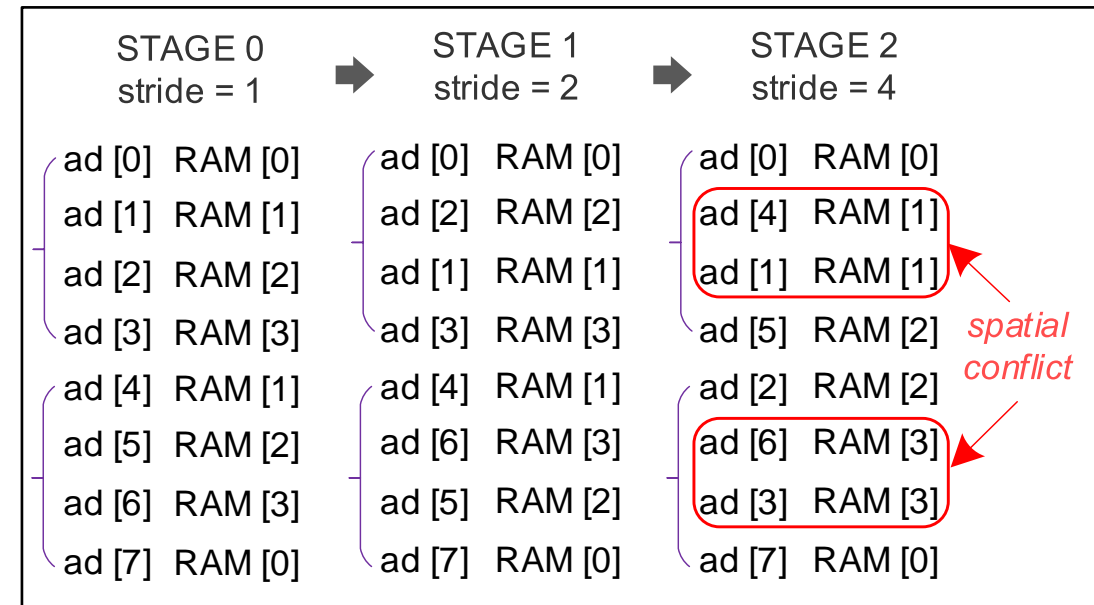
The dataflow of 16-point radix-2 in-place NTT with $d = 2$

Spatial Conflict

[Joh92] ✓ Memory mapping scheme for in-place FFT with arbitrary radix
 ✗ Placing multiple butterfly units leads to access conflict



(a) Radix-2 in-place NTT with $d = 1$.



(b) Radix-2 in-place NTT with $d = 2$.

Motivations

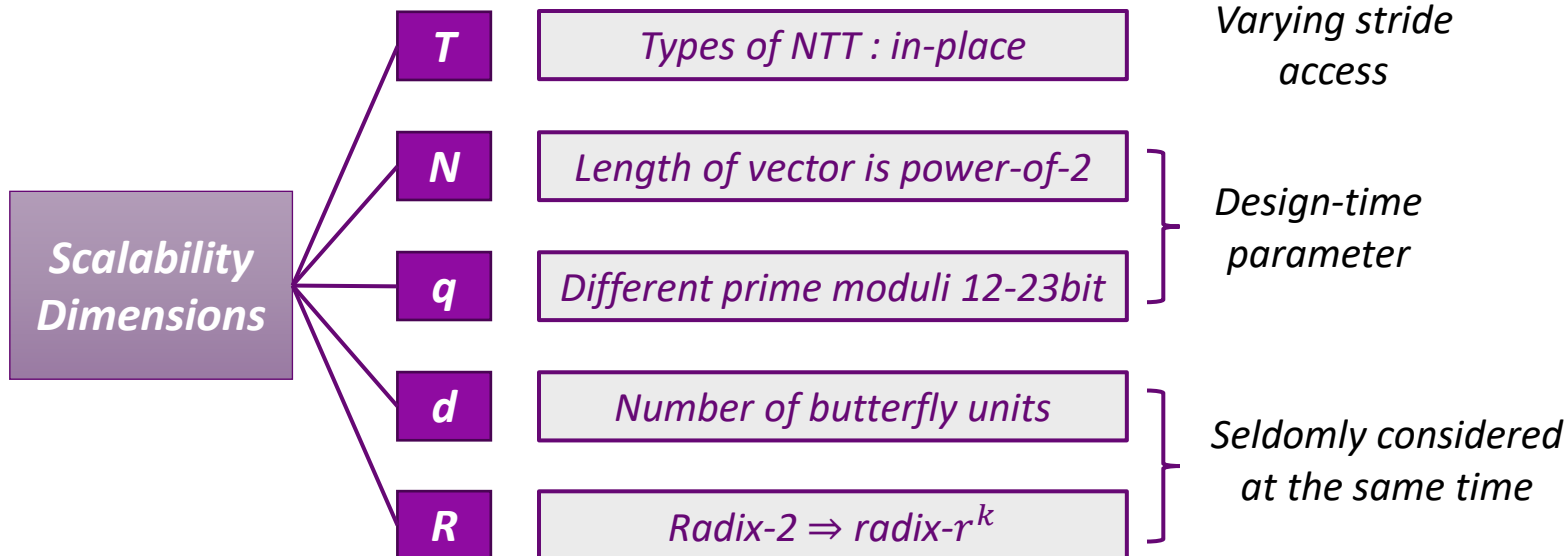
□ 3Ds in lattice-based PQC:

- **D**iverse security parameters of lattice-based scheme
- **D**ifferent resource constraints of computation platform
- **D**ifferent throughput requirements of practical application

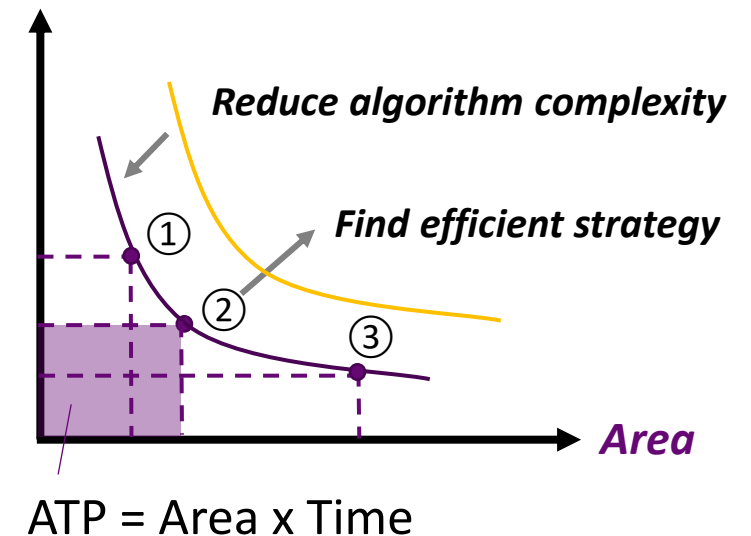
e.g. $\left\{ \begin{array}{l} N = 1024 \quad q = 12289 \quad \text{Falcon} \\ N = 256 \quad q = 8380417 \quad \text{Dilithium} \end{array} \right.$

Embedded device vs. Server

IoT vs. 5G



Time = Cycles / Frequency





Outlines

- Introduction
- **Optimized Radix-4 NTT/INTT Algorithm**
- Conflict-free Memory Mapping Scheme
- Hardware Architecture of Radix-2/4 NTT
- Implementation Result and Comparison

The Derivation of Radix-4 NTT without Preprocessing

$$A_i = \sum_{j=0}^{N-1} a_j \phi_{2N}^j \omega_N^{ij} \pmod q \xrightarrow{\text{Split}} A_i = \sum_{j=0}^{N/4-1} a_{4j} \phi_{2N}^{4j} \omega_N^{i \cdot (4j)} + \sum_{j=0}^{N/4-1} a_{4j+1} \phi_{2N}^{4j+1} \omega_N^{i \cdot (4j+1)} + \sum_{j=0}^{N/4-1} a_{4j+2} \phi_{2N}^{4j+2} \omega_N^{i \cdot (4j+2)} + \sum_{j=0}^{N/4-1} a_{4j+3} \phi_{2N}^{4j+3} \omega_N^{i \cdot (4j+3)} \pmod q$$

$$A_i = \sum_{j=0}^{N/4-1} \boxed{a_{4j} \phi_{N/2}^j \omega_{N/4}^{ij}} + \omega_N^i \cdot \phi_{2N}^1 \cdot \sum_{j=0}^{N/4-1} \boxed{a_{4j+1} \phi_{N/2}^j \omega_{N/4}^{ij}} + \omega_N^{2i} \cdot \phi_{2N}^2 \cdot \sum_{j=0}^{N/4-1} \boxed{a_{4j+2} \phi_{N/2}^j \omega_{N/4}^{ij}} + \omega_N^{3i} \cdot \phi_{2N}^3 \cdot \sum_{j=0}^{N/4-1} \boxed{a_{4j+3} \phi_{N/2}^j \omega_{N/4}^{ij}} \pmod q$$

$\checkmark N \log_4 N + 2N \Rightarrow N \log_4 N$

Periodicity property
 $\omega_N^{k+N} = \omega_N^k$

$i \in [0, N/4 - 1]$ **Four $N/4$ -point operation**

$$\begin{bmatrix} A_i \\ A_{i+N/4} \\ A_{i+2N/4} \\ A_{i+3N/4} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \omega_4^1 & -1 & -\omega_4^1 \\ 1 & -1 & 1 & -1 \\ 1 & -\omega_4^1 & -1 & \omega_4^1 \end{bmatrix} \times \left(\begin{bmatrix} F_0 \\ F_1 \\ F_2 \\ F_3 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ \phi_{2N}^{2i+1} \\ \phi_{2N}^{2(2i+1)} \\ \phi_{2N}^{3(2i+1)} \end{bmatrix} \right)$$

Recursion N/4 4-points

The Derivation of Radix-4 INTT without Postprocessing

$$a_i = N^{-1} \cdot \phi_{2N}^{-i} \cdot \sum_{j=0}^{N-1} A_j \omega_N^{-ij} \pmod q \xrightarrow{\text{Split}} a_i = N^{-1} \cdot (\phi_{2N}^{-i} \cdot \sum_{j=0}^{N/4-1} A_j \phi_{2N}^{-j} \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \sum_{j=N/4}^{2N/4-1} A_j \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \sum_{j=2N/4}^{3N/4-1} A_j \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \sum_{j=3N/4}^{N-1} A_j \omega_N^{-ij}) \pmod q$$

$$a_i = N^{-1} \cdot (\phi_{2N}^{-i} \cdot \sum_{j=0}^{N/4-1} A_j \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-i} \cdot \sum_{j=0}^{N/4-1} A_{j+N/4} \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-2i} \cdot \sum_{j=0}^{N/4-1} A_{j+2N/4} \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-3i} \cdot \sum_{j=0}^{N/4-1} A_{j+3N/4} \omega_N^{-ij}) \pmod q$$

Periodicity and Binary property

$$a_i = N^{-1} \cdot (\phi_{2N}^{-i} \cdot \sum_{j=0}^{N/4-1} A_j \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-i} \cdot \sum_{j=0}^{N/4-1} A_{j+N/4} \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-2i} \cdot \sum_{j=0}^{N/4-1} A_{j+2N/4} \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-3i} \cdot \sum_{j=0}^{N/4-1} A_{j+3N/4} \omega_N^{-ij}) \pmod q$$

Next page

Elimination property

$$\phi_{2N}^{-i} \cdot \omega_4^{-2i} \cdot \sum_{j=0}^{N/4-1} A_{j+2N/4} \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-3i} \cdot \sum_{j=0}^{N/4-1} A_{j+3N/4} \omega_N^{-ij}) \pmod q$$

The Derivation of Radix-4 INTT without Postprocessing

$$a_i = N^{-1} \cdot (\phi_{2N}^{-i} \cdot \sum_{j=0}^{N/4-1} A_j \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-i} \cdot \sum_{j=0}^{N/4-1} A_{j+N/4} \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-2i} \cdot \sum_{j=0}^{N/4-1} A_{j+2N/4} \omega_N^{-ij} + \phi_{2N}^{-i} \cdot \omega_4^{-3i} \cdot \sum_{j=0}^{N/4-1} A_{j+3N/4} \omega_N^{-ij}) \bmod q$$

$$\omega_4^{-i} = \begin{cases} 0 & i = 4r \\ \omega_4^{-1} & i = 4r + 1 \\ -1 & i = 4r + 2 \\ -\omega_4^{-1} & i = 4r + 3 \end{cases}$$

$$\omega_4^{-2i} = \begin{cases} 1 & i = 4r \\ -1 & i = 4r + 1 \\ 1 & i = 4r + 2 \\ -1 & i = 4r + 3 \end{cases}$$

$$\omega_4^{-3i} = \begin{cases} 0 & i = 4r \\ \omega_4^{-1} & i = 4r + 1 \\ -1 & i = 4r + 2 \\ -\omega_4^{-1} & i = 4r + 3 \end{cases}$$

$$a_{4i} = \left(\frac{N}{4}\right)^{-1} \cdot \left(\frac{1}{4} \cdot \phi_{N/2}^{-i} \cdot \sum_{j=0}^{N/4-1} A_j \omega_{N/4}^{-ij} + \frac{1}{4} \cdot \phi_{N/2}^{-i} \cdot \sum_{j=0}^{N/4-1} A_{j+N/4} \omega_{N/4}^{-ij} + \frac{1}{4} \cdot \phi_{N/2}^{-i} \cdot \sum_{j=0}^{N/4-1} A_{j+2N/4} \omega_{N/4}^{-ij} + \frac{1}{4} \cdot \phi_{N/2}^{-i} \cdot \sum_{j=0}^{N/4-1} A_{j+3N/4} \omega_{N/4}^{-ij} \right) \bmod q$$

Sum and Binary property

✓ $N \log_4 N + N \Rightarrow N \log_4 N$

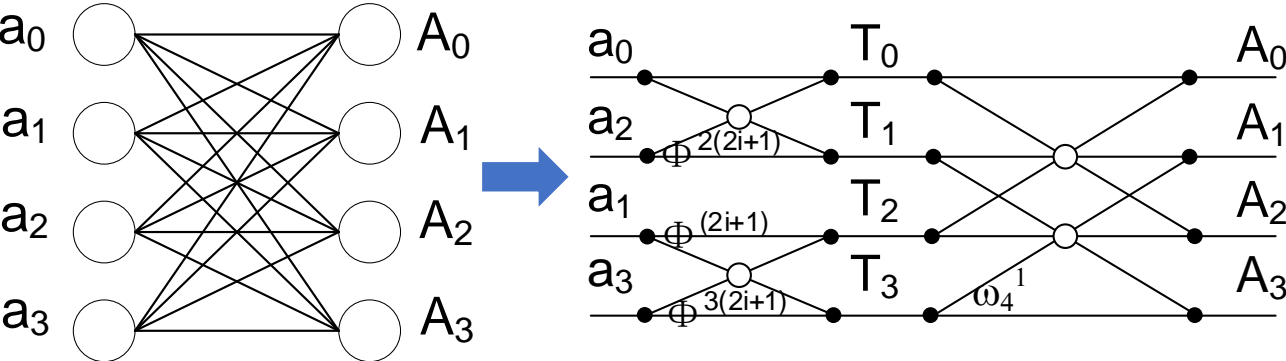
$i \in [0, N/4 - 1]$ Four $N/4$ -point operation

$$\begin{bmatrix} a_{4i} \\ a_{4i+1} \\ a_{4i+2} \\ a_{4i+3} \end{bmatrix} = \left(\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \omega_4^{-1} & -1 & -\omega_4^{-1} \\ 1 & -1 & 1 & -1 \\ 1 & -\omega_4^{-1} & -1 & \omega_4^{-1} \end{bmatrix} \times \begin{bmatrix} G_0 \\ G_1 \\ G_2 \\ G_3 \end{bmatrix} \right) \cdot \begin{bmatrix} 1 \\ \phi_{2N}^{-(2j+1)} \\ \phi_{2N}^{-2(2j+1)} \\ \phi_{2N}^{-3(2j+1)} \end{bmatrix}$$

Recursion

N/4 4-points

Divide and Conquer for Butterfly Operation

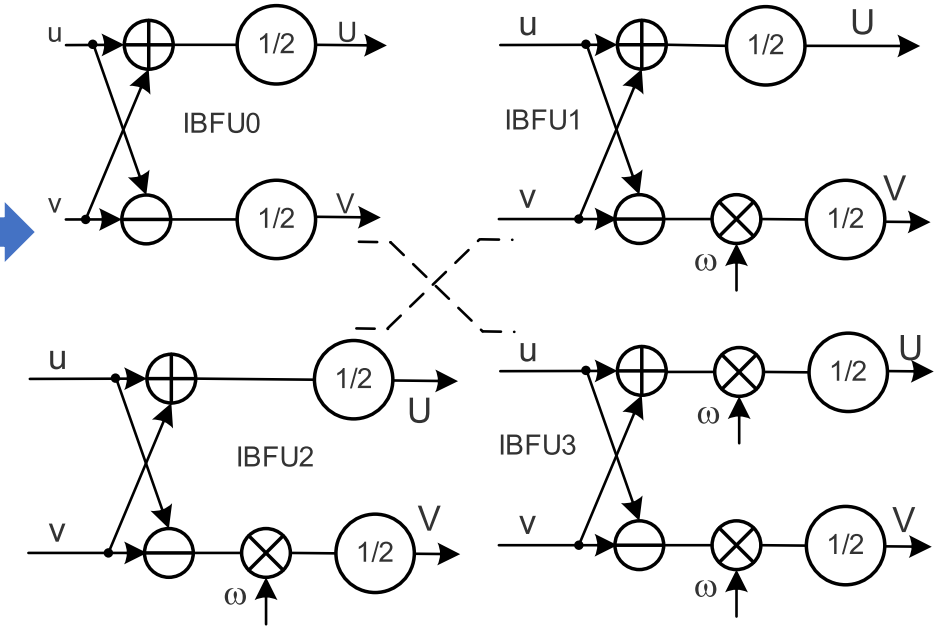
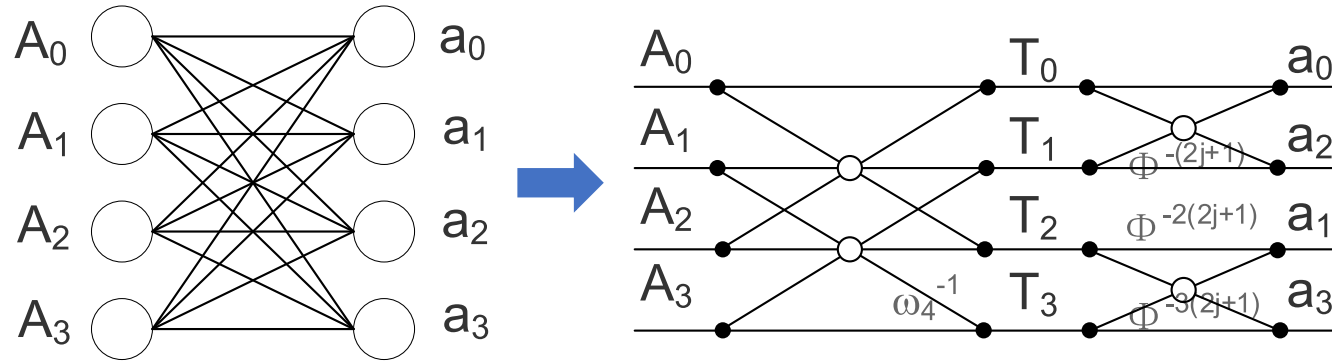


Two-layer radix-4 Cooley-Tukey butterfly unit in forward NTT

$$\begin{aligned}
 T_0 &= (F_0 + F_2 \cdot \phi_{2N}^{2(2i+1)}) & A_i &= T_0 + T_2 \\
 T_1 &= (F_0 - F_2 \cdot \phi_{2N}^{2(2i+1)}) & A_{i+N/4} &= (T_1 + T_3 \cdot \omega_4^1) \\
 T_2 &= (F_1 \cdot \phi_{2N}^{2i+1} + F_3 \cdot \phi_{2N}^{3(2i+1)}) & A_{i+2N/4} &= (T_0 - T_2) \\
 T_3 &= (F_1 \cdot \phi_{2N}^{2i+1} - F_3 \cdot \phi_{2N}^{3(2i+1)}) & A_{i+3N/4} &= (T_1 - T_3 \cdot \omega_4^1)
 \end{aligned}$$

Num. of Op.	MM.	MA.	MS.
Direct	#5	#6	#6
Two-layer	#4	#4	#4
Variation	↓20%	↓33%	↓33%

Divide and Conquer for Butterfly Operation



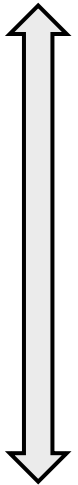
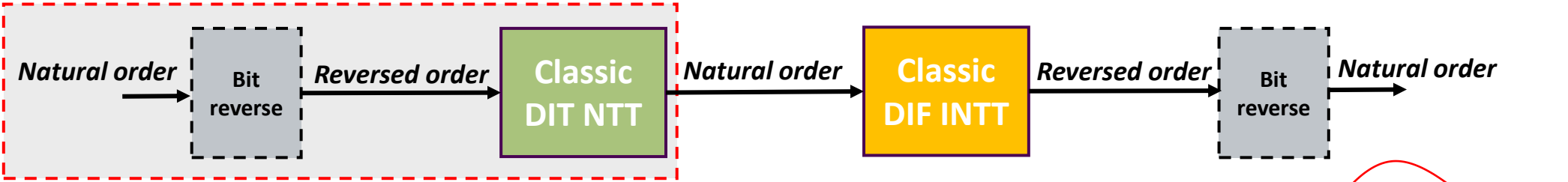
Two-layer radix-4 Gentleman-Sande butterfly unit in inverse NTT

$$\begin{aligned}
 T_0 &= F_0 + F_2 & a_{4i} &= T_0 + T_2 \\
 T_1 &= F_0 - F_2 & a_{4i+1} &= (T_1 + T_3) \cdot \phi_{2N}^{-(2j+1)} \\
 T_2 &= F_1 + F_3 & a_{4i+2} &= (T_0 - T_2) \cdot \phi_{2N}^{-2(2j+1)} \\
 T_3 &= (F_1 - F_3) \cdot \omega_4^{-1} & a_{4i+3} &= (T_1 - T_3) \cdot \phi_{2N}^{-3(2j+1)}
 \end{aligned}$$

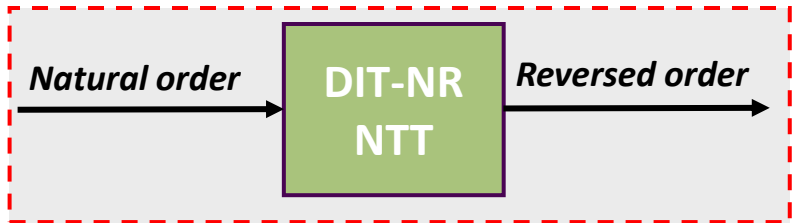
Num. of Op.	MM.	MA.	MS.
Direct	#5	#6	#6
Two-layer	#4	#4	#4
Variation	↓20%	↓33%	↓33%

Proposed DIT-NR Radix-4 NTT

❑ The bit-reversed operation is needed in classic NTT and INTT operation.



A generic method to avoid bit reversed issue in radix-R NTT



① Reverse the first loop p for 0 to $\log_4 N - 1$

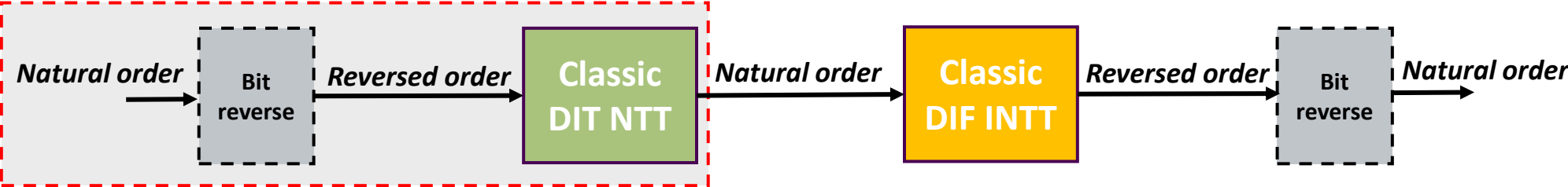
② Adjust the generation of twiddle factors

$$\begin{cases} \phi_{2N}^{N/(4J)} \Rightarrow \phi_{2N}^J \\ \omega_m^{2j+1} \Rightarrow \omega_m^{2 \cdot \text{reversed}(k)+1} \\ \omega_m^{2(2j+1)} \Rightarrow \omega_m^{2 \cdot (2 \cdot \text{reversed}(k)+1)} \\ \omega_m^{3(2j+1)} \Rightarrow \omega_m^{3 \cdot (2 \cdot \text{reversed}(k)+1)} \end{cases}$$

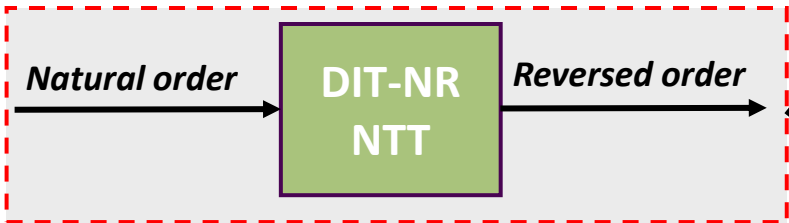
③ Move the place of generating twiddle factors from loop j to loop k .

Proposed DIT-NR Radix-4 NTT

□ The bit-reversed operation is needed in classic NTT and INTT operation.



A generic method to avoid bit reversed issue in radix-R NTT

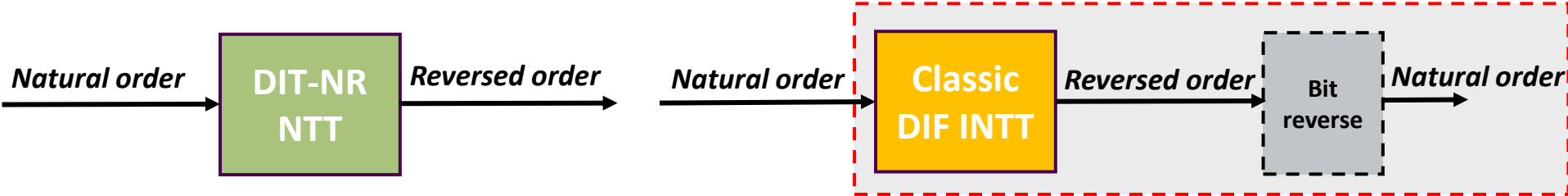


```

for p = log4 N - 1 to 0 do
  J ← 4p
  r ← 0
  for k = 0 to N/(4J) - 1 do
    ω1 ← ωa1_ROM [r]
    ω2 ← ωa2_ROM [r]
    ω3 ← ωa3_ROM [r]
    r ← r + 1
    for j = 0 to J - 1 do
      T0 ← (a4kJ+j + a4kJ+j+2J · ω2) mod q
      T1 ← (a4kJ+j - a4kJ+j+2J · ω2) mod q
      T2 ← (a4kJ+j+J · ω1 + a4kJ+j+3J · ω3) mod q
      T3 ← (a4kJ+j+J · ω1 - a4kJ+j+3J · ω3) mod q
      A4kJ+j ← (T0 + T2) mod q
      A4kJ+j+J ← (T1 + T3 · ω41) mod q
      A4kJ+j+2J ← (T0 - T2) mod q
      A4kJ+j+3J ← (T1 - T3 · ω41) mod q
    end for
  end for
end for
end for
  
```


Proposed DIF-RN Radix-4 INTT

❑ The bit-reversed operation is needed in classic NTT and INTT operation.



for $\log_4 N - 1$ to 0

① reverse the first loop

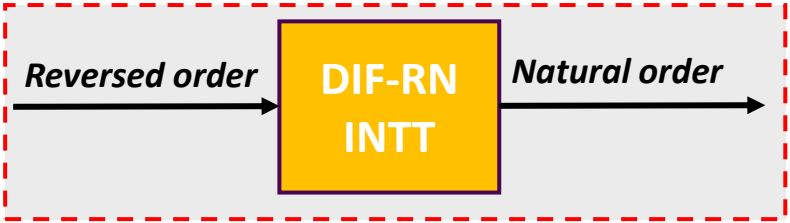
② replace the twiddle factor ω_m^x of NR - NTT with its inverse elements ω_m^{-x}

A generic method to avoid bit reversed issue in radix-R INTT

Double memory footprint is needed for storage of twiddle factors !

Opt. tech.

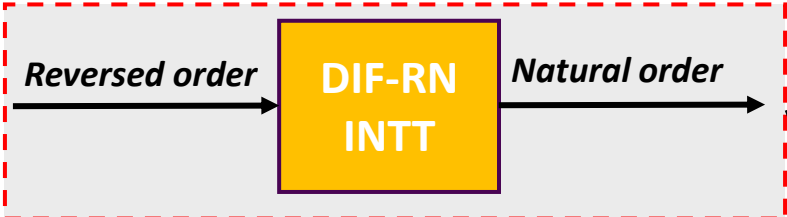
Try to reuse the twiddle factors of forward NTT to reduce memory footprint



Proposed DIF-RN Radix-4 INTT

Opt. tech. 1: Derive three new tricks as following

$$\left\{ \begin{aligned} \phi_{2N}^{-i} \bmod q &= \phi_{2N}^{-N/2} \cdot \phi_{2N}^{N/2-i} = \omega_4^{-1} \cdot \phi_{2N}^{N/2-i} = -\omega_4^1 \cdot \phi_{2N}^{N/2-i} \bmod q \\ \phi_{2N}^{-i} \bmod q &= -\phi_{2N}^N \cdot \phi_{2N}^{N-i} = -\phi_{2N}^{N-i} \bmod q \\ \phi_{2N}^{-i} \bmod q &= \phi_{2N}^{-3N/2} \cdot \phi_{2N}^{3N/2-i} = \omega_4^1 \cdot \phi_{2N}^{3N/2-i} \bmod q \end{aligned} \right.$$



Opt. tech. 2: Apply the derived tricks to obtain the ω_m^{-x} by resuing ω_m^x

$$\left\{ \begin{aligned} \omega_m^{-(2j+1)} &\Rightarrow -\omega_4^1 \cdot \omega_m^{N/2-[2 \cdot \text{reversed}(k)+1]} \\ \omega_m^{-2(2j+1)} &\Rightarrow -\omega_m^{N-[2 \cdot (2 \cdot \text{reversed}(k)+1)]} \\ \omega_m^{-3(2j+1)} &\Rightarrow \omega_4^1 \cdot \omega_m^{3N/2-[3 \cdot (2 \cdot \text{reversed}(k)+1)]} \end{aligned} \right.$$

Reuse the twiddle factor of forward NTT

Modify the two-layer radix-4 GS-BFU

```

for p = log4 N - 1 to 0 do
  J ← 4p
  r ← 0
  for k = 0 to N/(4J) - 1 do
    ω1 ← ωa1_ROM [r]
    ω2 ← ωa2_ROM [r]
    ω3 ← ωa3_ROM [r]
    r ← r + 1
    for j = 0 to J - 1 do
      T0 ← (a4kJ+j + a4kJ+j+2J · ω2) mod q
      T1 ← (a4kJ+j - a4kJ+j+2J · ω2) mod q
      T2 ← (a4kJ+j+J · ω1 + a4kJ+j+3J · ω3) mod q
      T3 ← (a4kJ+j+J · ω1 - a4kJ+j+3J · ω3) mod q
      A4kJ+j ← (T0 + T2) mod q
      A4kJ+j+J ← (T1 + T3 · ω41) mod q
      A4kJ+j+2J ← (T0 - T2) mod q
      A4kJ+j+3J ← (T1 - T3 · ω41) mod q
    end for
  end for
end for
  
```

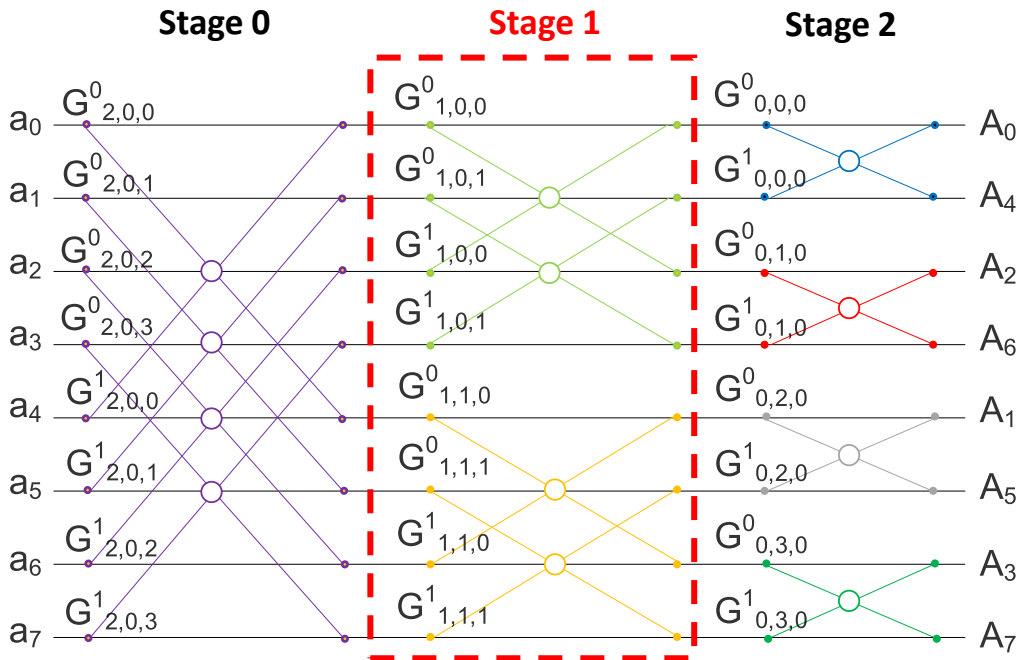
↓ 50% memory footprint of twiddle factors



Outlines

- Introduction
- Optimized Radix-4 NTT/INTT Algorithm
- **Conflict-free Memory Mapping Scheme**
- Hardware Architecture of Radix-2/4 NTT
- Implementation Result and Comparison

Determine the Point-access Order



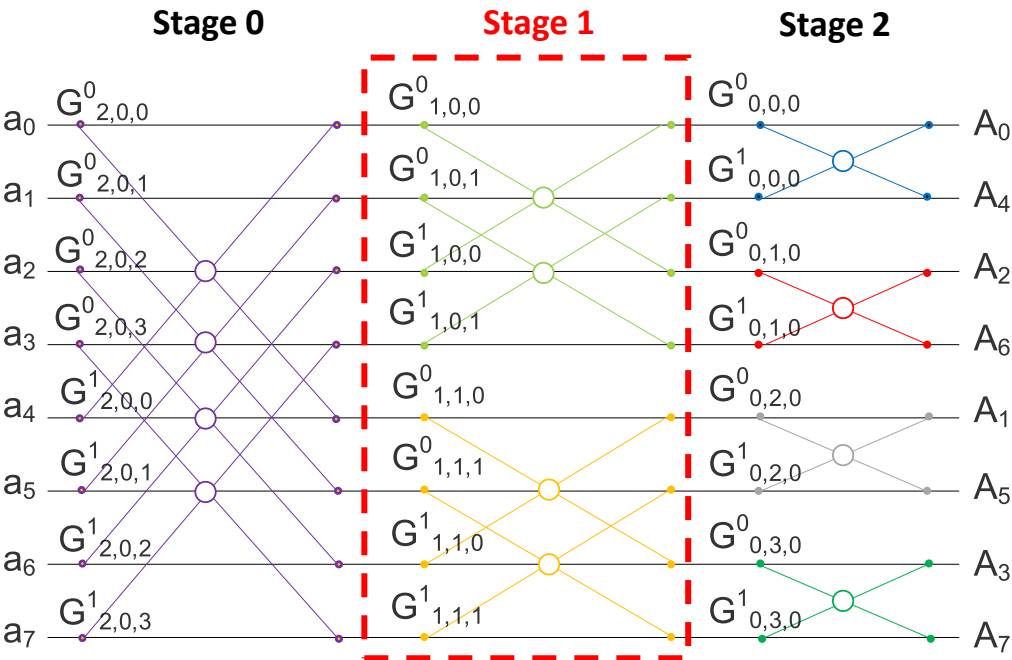
The data flow of 8 points radix-2 DIT-NR NTT.

□ $G^i_{p,k,j}$: The i -th point in stage p , group k and round j
where $i = 0, 1, \dots, R-1$ and R denotes the radix

□ The four data points at **stage 1** can be parallelly accessed in two types of order as below:

- ① Four data points with the same **group index k**
 $\{G^0_{1,0,0}, G^1_{1,0,0}, G^0_{1,0,1}, G^1_{1,0,1}\} \rightarrow \{G^0_{1,1,0}, G^1_{1,1,0}, G^0_{1,1,1}, G^1_{1,1,1}\}$
- ② Four data points with the same **round index j**
 $\{G^0_{1,0,0}, G^1_{1,0,0}, G^0_{1,1,0}, G^1_{1,1,0}\} \rightarrow \{G^0_{1,0,1}, G^1_{1,0,1}, G^0_{1,1,1}, G^1_{1,1,1}\}$

Determine the Point-access Order



The data flow of 8 points radix-2 DIT-NR NTT.

□ $G_{p,k,j}^i$: The i -th point in stage p , group k and round j

where $i = 0, 1, \dots, R-1$ and R denotes the radix

□ The four data points at stage 1 can be parallelly accessed in two types of order as below:

① Four data points with the same group index k

$$\{G_{1,0,0}^0, G_{1,0,0}^1, G_{1,0,1}^0, G_{1,0,1}^1\} \rightarrow \{G_{1,1,0}^0, G_{1,1,0}^1, G_{1,1,1}^0, G_{1,1,1}^1\}$$

② Four data points with the same round index j

$$\{G_{1,0,0}^0, G_{1,0,0}^1, G_{1,1,0}^0, G_{1,1,0}^1\} \rightarrow \{G_{1,0,1}^0, G_{1,0,1}^1, G_{1,1,1}^0, G_{1,1,1}^1\}$$

Proposed Parallel NTT Algorithm with Arbitrary Radix

Algorithm 6 Scalable Iterative NTT Algorithm

Input: a, N, R . Here a denotes a vector of length N . d is the number of butterfly units. R denotes the radix of NTT.

Output: $A = Scalable_NTT(a)$

```

1: for  $p = \log_R N - 1$  to 0 do
2:    $J \leftarrow R^p$ 
3:   if  $J < d$  then
4:     for  $k = 0$  to  $N/(Rd) - 1$  do
5:       for  $i = 0$  to  $d/J - 1$  do
6:         for  $j = 0$  to  $J - 1$  do
7:           
$$\begin{bmatrix} A_0 \\ A_1 \\ \dots \\ A_{R-1} \end{bmatrix} = BF \left( \begin{bmatrix} a_{kRd+iRJ+j} \\ a_{kRd+iRJ+j+J} \\ \dots \\ a_{kRd+iRJ+j+(R-1)J} \end{bmatrix} \right)$$

8:         end for
9:       end for
10:    end for
11:   else
12:     for  $k = 0$  to  $N/(RJ) - 1$  do
13:       for  $i = 0$  to  $J/d - 1$  do
14:         for  $j = 0$  to  $d - 1$  do
15:           
$$\begin{bmatrix} A_0 \\ A_1 \\ \dots \\ A_{R-1} \end{bmatrix} = BF \left( \begin{bmatrix} a_{kRJ+id+j} \\ a_{kRJ+id+j+J} \\ \dots \\ a_{kRJ+id+j+(R-1)J} \end{bmatrix} \right)$$

16:         end for
17:       end for
18:     end for
19:   end if
20: end for
21: return  $A$ 

```

Loop Tile 1

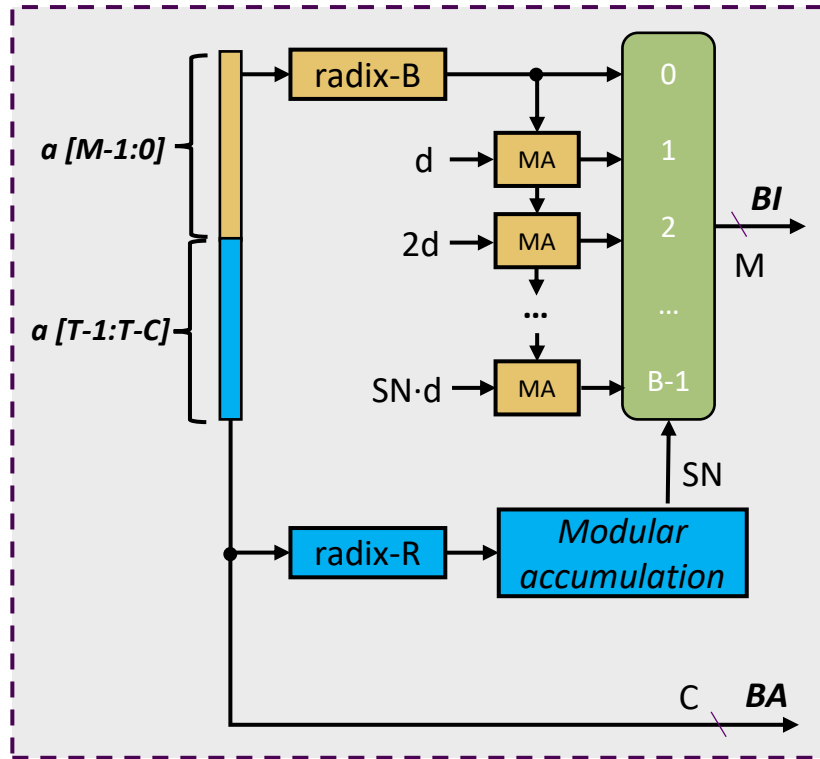
Loop Tile 2

□ The relative size of d and J will influence the number of iterations.

• $J < d$, the parallel data sets will contain several iterations of inner loop.

• $J \geq d$, the inner loop will cover several parallel data sets.

Conflict-free Memory Mapping for Arbitrary Radix NTT



Memory mapper for arbitrary radix r^k

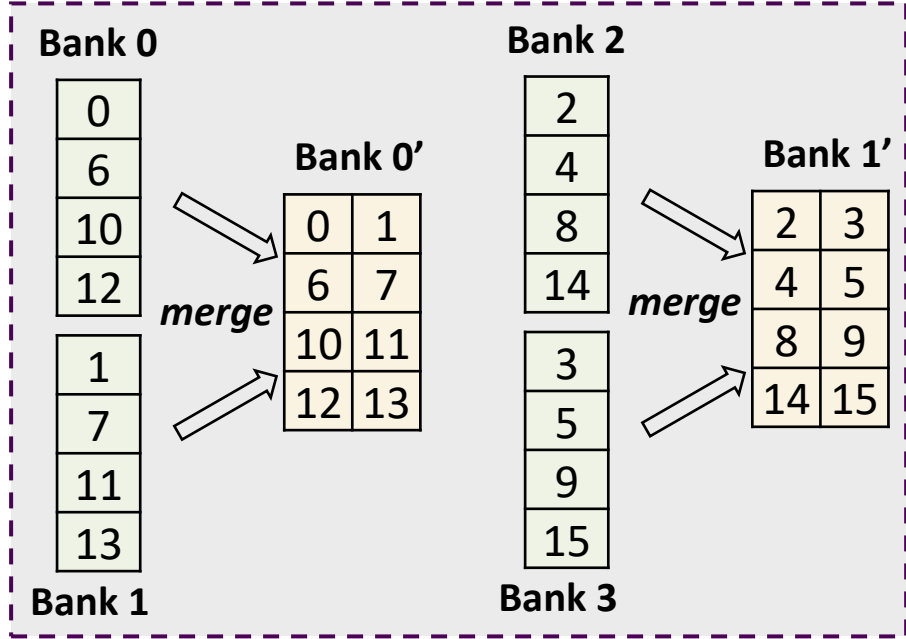
- ❑ d : Number of BFUs;
- ❑ B : Number of banks, $B = R \times d, N \mid B$;
- ❑ $T = \log_r N$; $M = \log_r B$; $C = T - M$;
- ✓ Automated configuration
- ✓ Initialized for different radix $R = r^k$
- ✓ Facilitate other stride access
- ✓ We further extend this method in our paper at DAC 2022*

*Xiangren Chen, Bohan Yang, ..., Leibo Liu. 2022. Efficient access scheme for multi-bank based NTT architecture through conflict graph. (DAC '22).

Case Study of Conflict-free Mapping

Parameter	N = 16	R = 2	d = 2	B = 4
-----------	--------	-------	-------	-------

old address	mix radix R-B	step number	slide distance	bank index	bank address
0	000			$(0+0) \bmod 4 = 0$	00
1	001	$(0+0)$	(0×2)	$(0+1) \bmod 4 = 1$	
2	002	$\bmod 2 = 0$	$\bmod 4 = 0$	$(0+2) \bmod 4 = 2$	
3	003			$(0+3) \bmod 4 = 3$	
4	010			$(2+0) \bmod 4 = 2$	01
5	011	$(0+1)$	(1×2)	$(2+1) \bmod 4 = 3$	
6	012	$\bmod 2 = 1$	$\bmod 4 = 2$	$(2+2) \bmod 4 = 0$	
7	013			$(2+3) \bmod 4 = 1$	
8	100			$(2+0) \bmod 4 = 2$	10
9	101	$(1+0)$	(1×2)	$(2+1) \bmod 4 = 3$	
10	102	$\bmod 2 = 1$	$\bmod 4 = 2$	$(2+2) \bmod 4 = 0$	
11	103			$(2+3) \bmod 4 = 1$	
12	110			$(0+0) \bmod 4 = 0$	11
13	111	$(1+0)$	(0×2)	$(0+1) \bmod 4 = 1$	
14	112	$\bmod 2 = 0$	$\bmod 4 = 0$	$(0+2) \bmod 4 = 2$	
15	113			$(0+3) \bmod 4 = 3$	



Case study of 16-point radix-2 in-place NTT with 2 BFUs

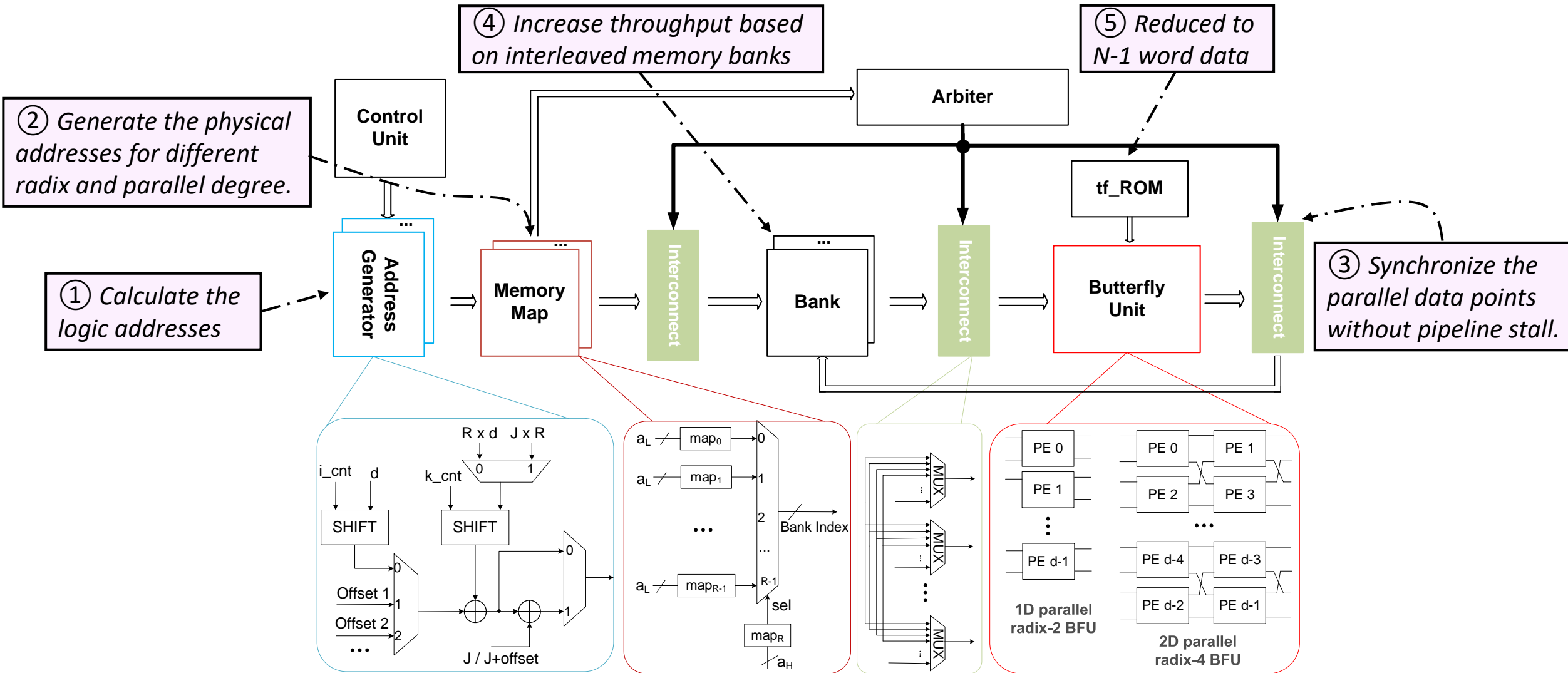
Reduce 4 banks to 2 banks



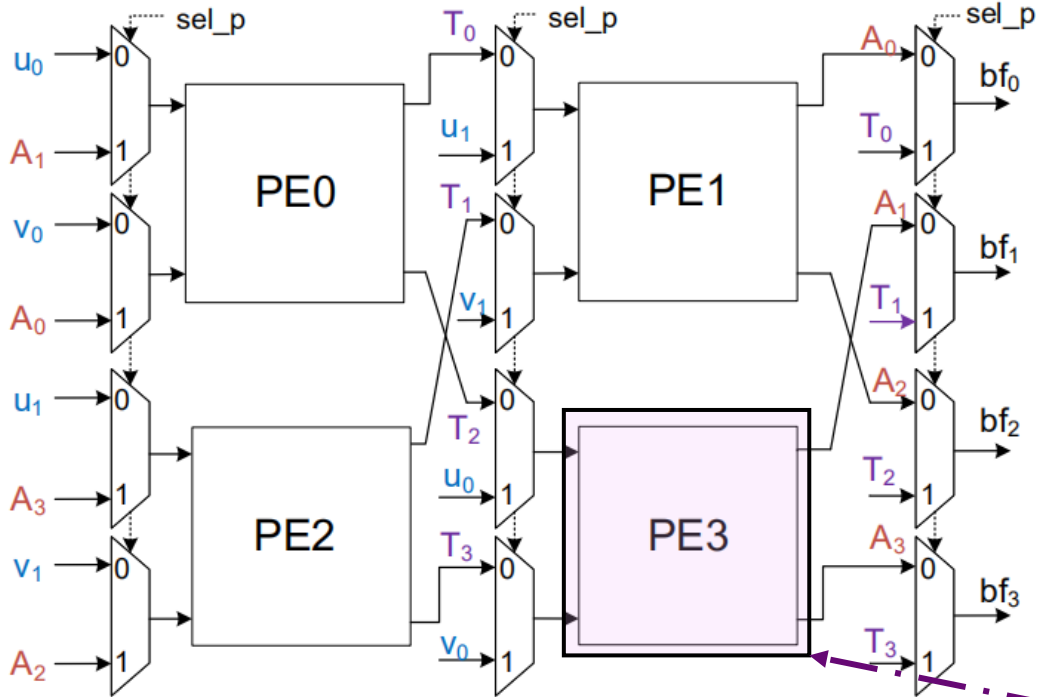
Outlines

- Introduction
- Optimized Radix-4 NTT/INTT Algorithm
- Conflict-free Memory Mapping Scheme
- **Hardware Architecture of Radix-2/4 NTT**
- Implementation Result and Comparison

Hardware Architecture of Radix-2/4 NTT



Unified Design of Radix-4 Butterfly Unit



The routing structure of radix-4 PE array.

PE types	#MM	#MA	#MS	#Half
PE0	1	1/2	1/2	2
PE1	0	1/2	1/2	2
PE2	2	1/2	1/2	2
PE3	1	1/2	1/2	2
Total	4	4/8	4/8	8

- ✓ sel_p = 0/1 perform NTT/INTT
- ✓ Barrett reduction based modular multiplier
- ✓ The multiplication by ω_4^1 is a constant modular multiplier within PE3.



Outlines

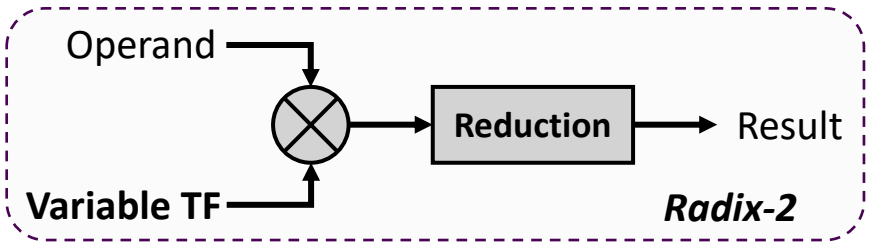
- Introduction
- Optimized Radix-4 NTT/INTT Algorithm
- Conflict-free Memory Mapping Scheme
- Hardware Architecture of Radix-2/4 NTT
- **Implementation Result and Comparison**

Implementation Result and Comparison

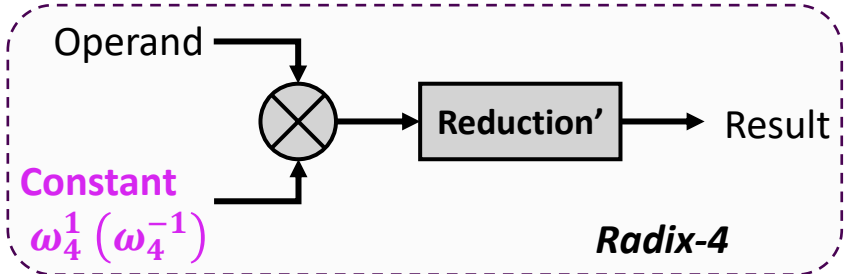
□ Theoretical Evaluation Between Radix-4 and Radix-2 Multi-bank based NTT

Type	AGU Address Generator	MMU Memory Map	AB Arbiter	INN Interconnect Network	Banks	BFU Butterfly Unit	TF ROM Twiddle Factor	NTT/INTT cycles
Radix-2	$2n$	$2n \times 2\text{-to-1 MUXs}$	$2n \times 2\text{-to-1 MUXs}$	$2n \times 2\text{-to-1 MUXs}$	$2n$	$n \times MA/MS/MM$	$N-1$	$(N/2d)\log_2 N$
Radix-4	n	$n \times 4\text{-to-1 MUXs}$	$n \times 4\text{-to-1 MUXs}$	$n \times 4\text{-to-1 MUXs}$	n	$n \times MA/MS/MM^*$	$N-1$	$(N/d)\log_4 N$
Variation	↓ 50%	—	↓ 75%	↓ 75%	↓ 50%	—	—	—

n : number of processing element



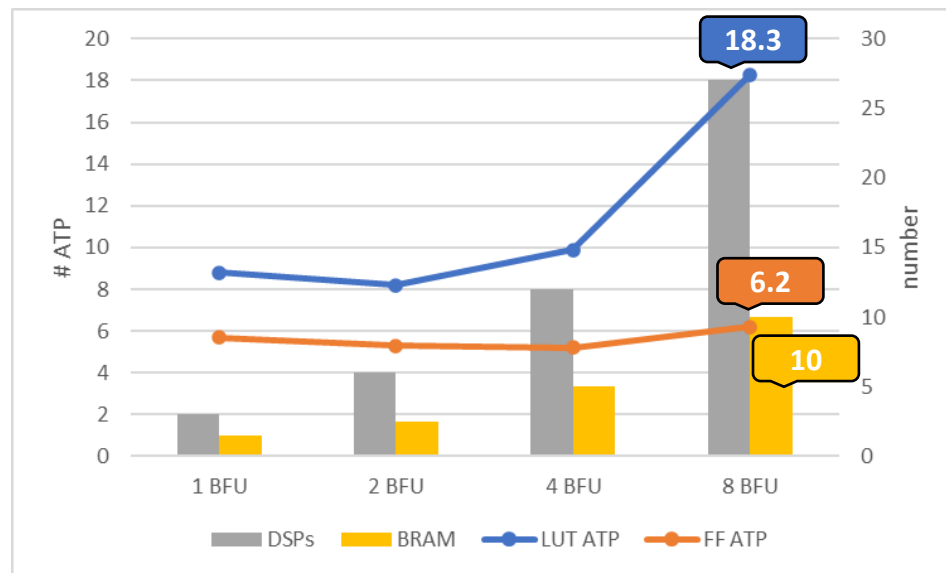
↓ Mult. Complexity
→



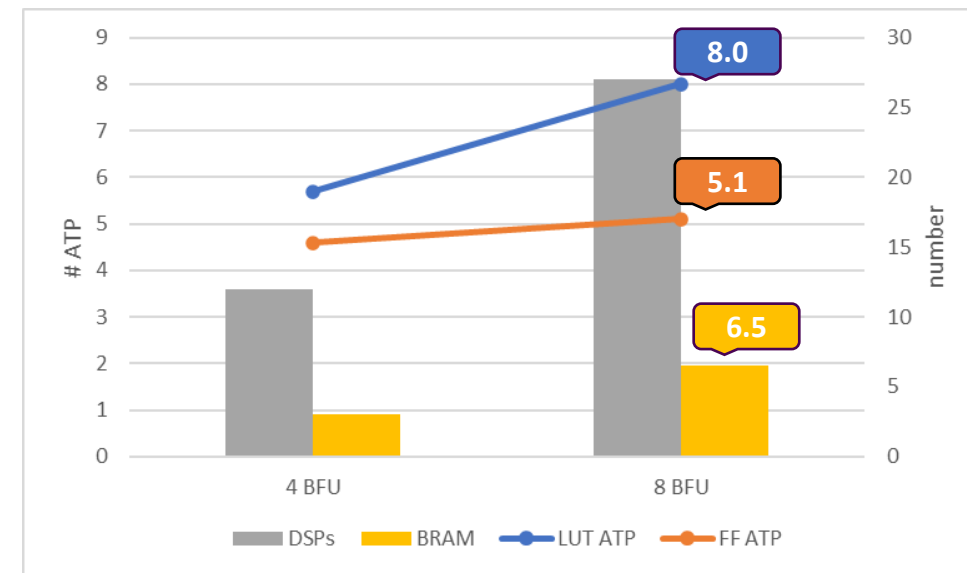
Implementation Result and Comparison

□ Implementation platform: Vivado 2020.2 + Virtex-7 FPGA (xc7vx690tfffg1761-3)

□ Benefits of radix-4 : e.g. #BFU = 8 \Rightarrow LUT ATP $\downarrow 2.2 \times$ FF ATP $\downarrow 1.2 \times$ #BRAMs $\downarrow 1.5$



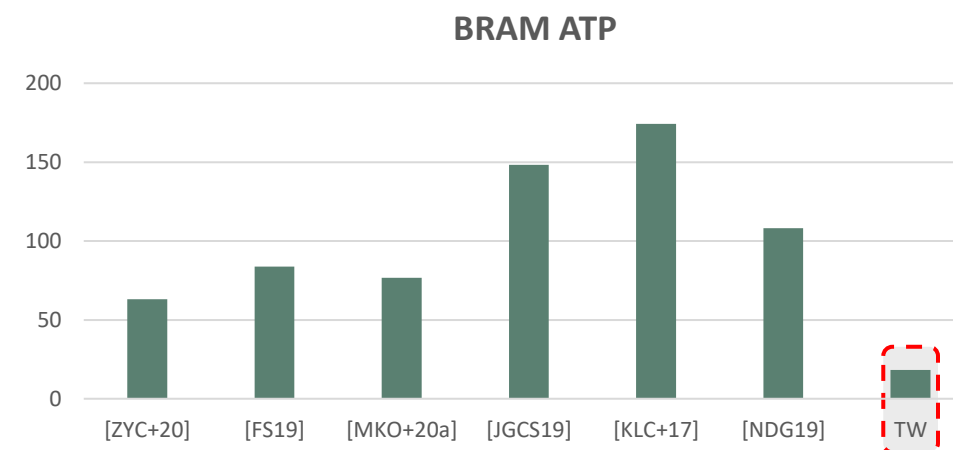
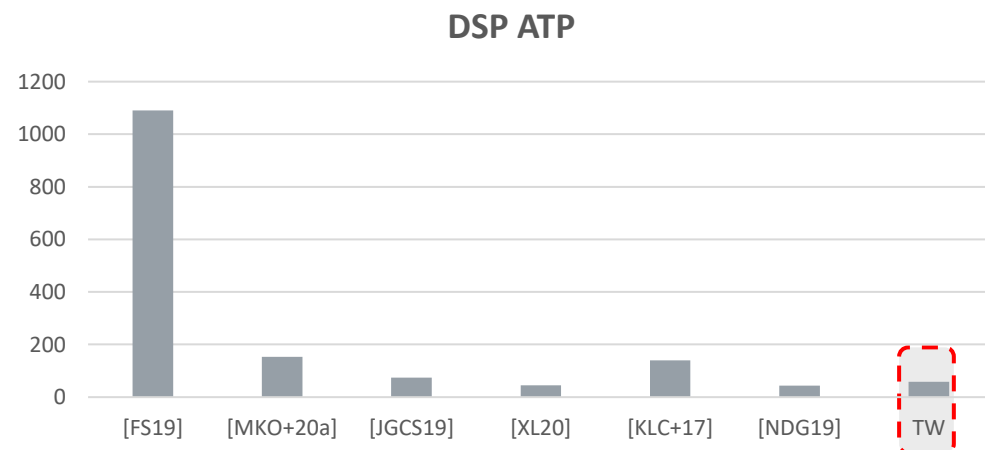
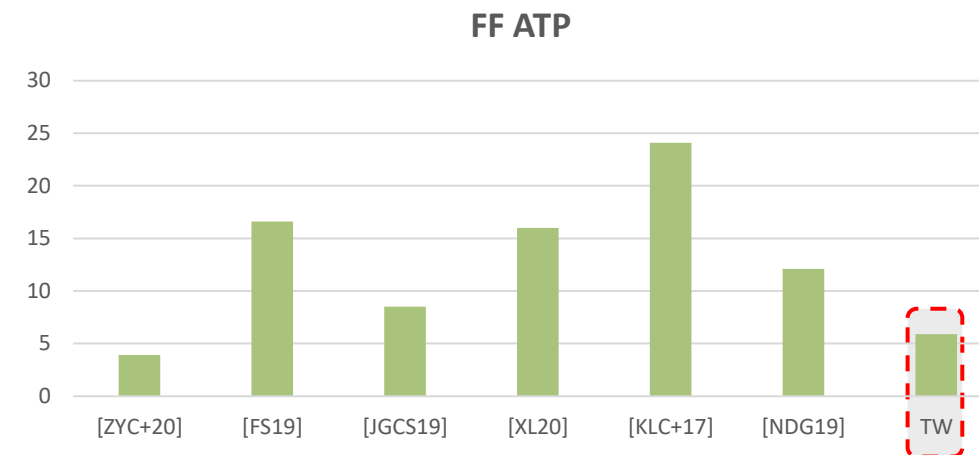
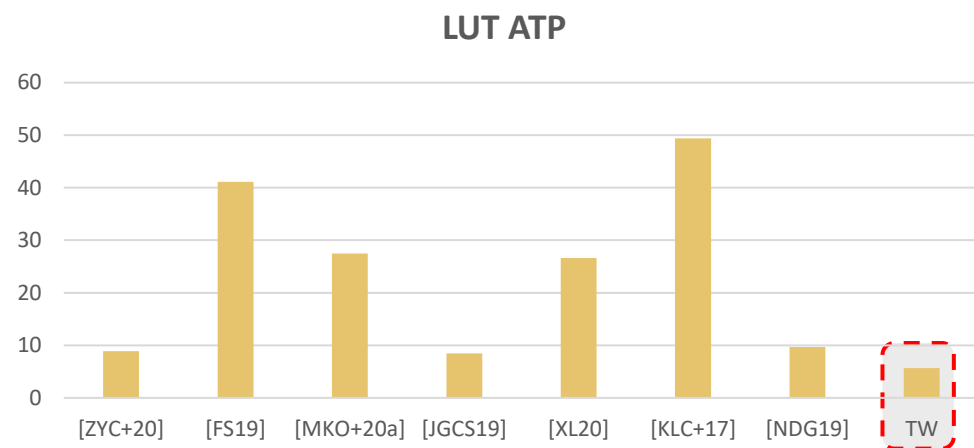
(a) radix-2 NTT



(b) radix-4 NTT

The variation of ATP in radix-2 and radix-4 NTT ($N = 1024, 14\text{-bit } q$)

Implementation Result and Comparison



The comparison between radix-4 NTT core and other works.



Thank You!

Q & A

Reference code: https://github.com/xiang-rc/cfntt_ref