

Engineering Lattice-based Cryptography

Sujoy Sinha Roy

Solving system of linear equations

- System of linear equations with unknown s

$$\begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n,1} & a_{n,2} & \cdots & a_{n,n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{pmatrix} \cdot \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \\ \vdots \\ b_m \end{pmatrix}$$

- Gaussian elimination solves s when number of equations $m \geq n$

System of linear equations with errors

$$\begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n,1} & a_{n,2} & \cdots & a_{n,n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{pmatrix} \cdot \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \\ \vdots \\ e_m \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \\ \vdots \\ b_m \end{pmatrix} \pmod{q}$$

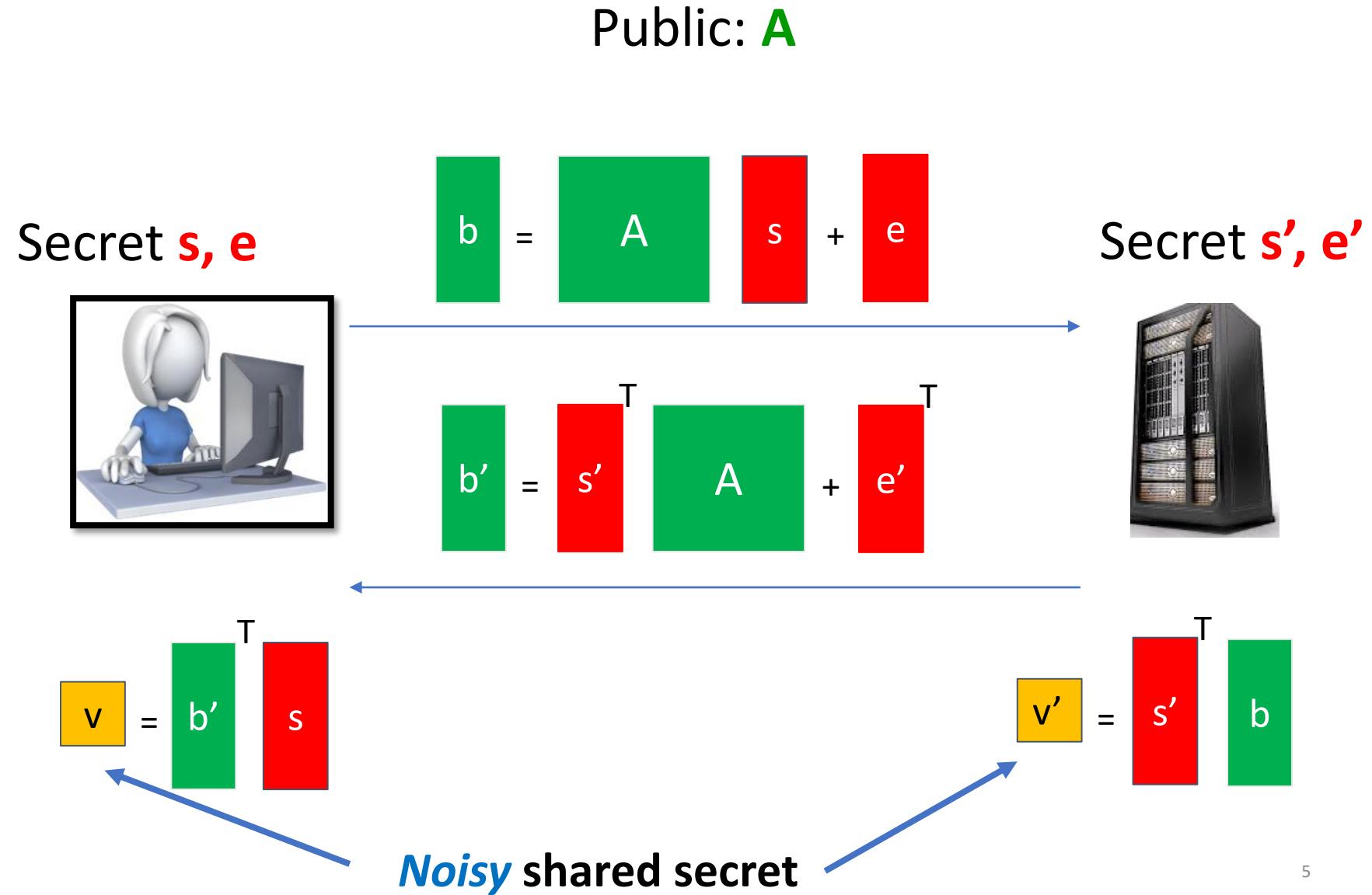
- Search **Learning with Errors (LWE)**-problem:
Given $(\mathbf{A}, \mathbf{b}) \rightarrow$ computationally infeasible to solve (\mathbf{s}, \mathbf{e})
- Decisional LWE-problem:
Given $(\mathbf{A}, \mathbf{b}) \rightarrow$ hard to distinguish from random

Learning with rounding (LWR)

$$\frac{p}{q} \begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n,1} & a_{n,2} & \cdots & a_{n,n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{pmatrix} \cdot \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \\ \vdots \\ b_m \end{pmatrix}$$

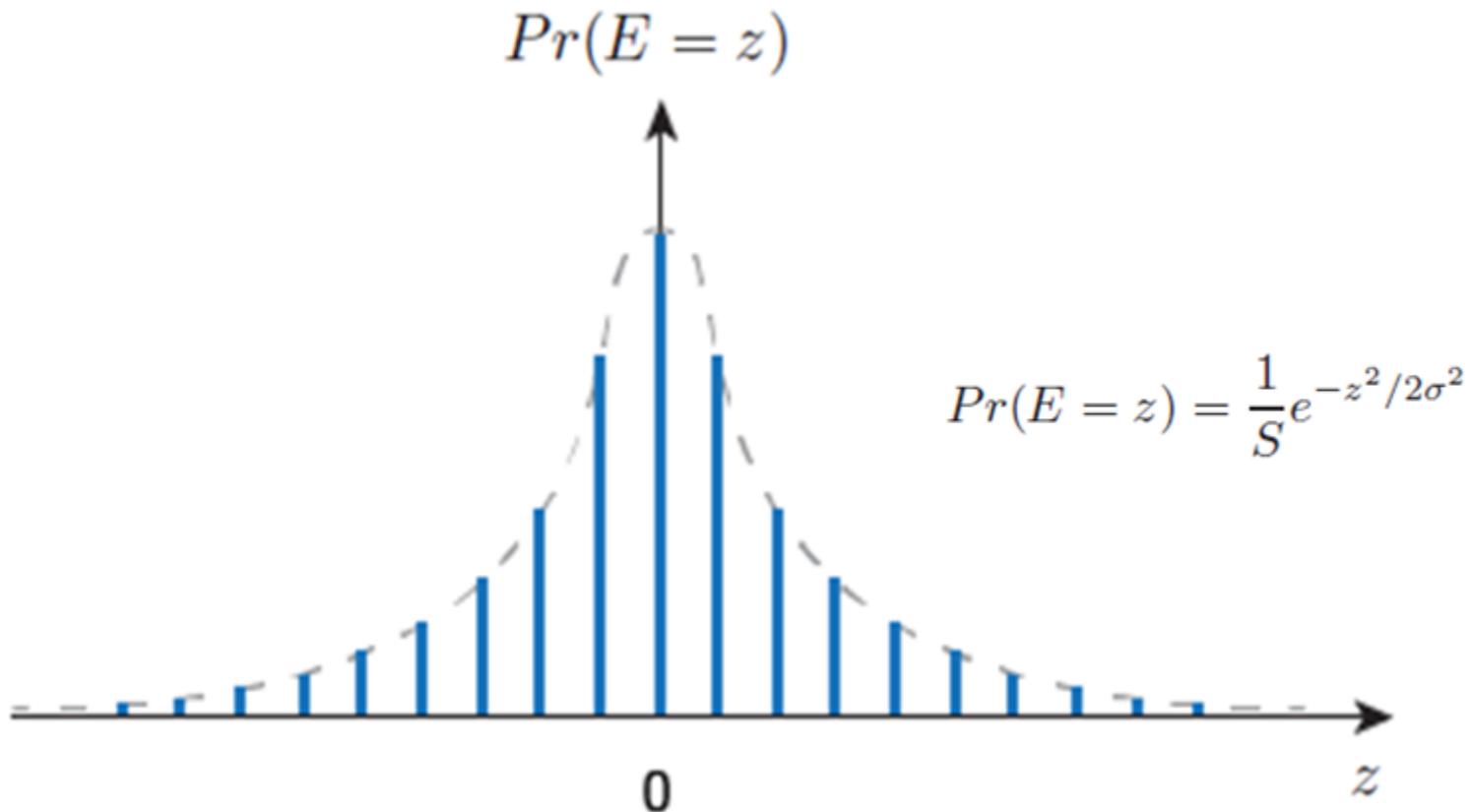
- Search LWE-problem:
Given $(\mathbf{A}, \mathbf{b}) \rightarrow$ computationally infeasible to solve (\mathbf{s}, \mathbf{e})
- Decisional LWE-problem:
Given $(\mathbf{A}, \mathbf{b}) \rightarrow$ hard to distinguish from random

LWE Diffie-Hellman key-exchange



- **Error sampling**
- Matrix-vector multiplication

Sampling from a known discrete distribution



Knuth-Yao random walk

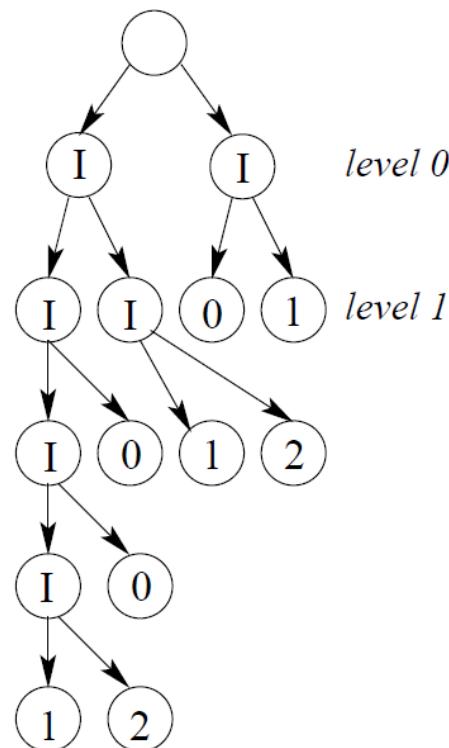
- Example: Let a sample space $S = \{ 0, 1, 2 \}$
 - $p_0 = 0.01110$
 - $p_1 = 0.01101$
 - $p_2 = 0.00101$
- Probability matrix

$$P_{mat} = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$$

row 0 → ↓ *column 0*

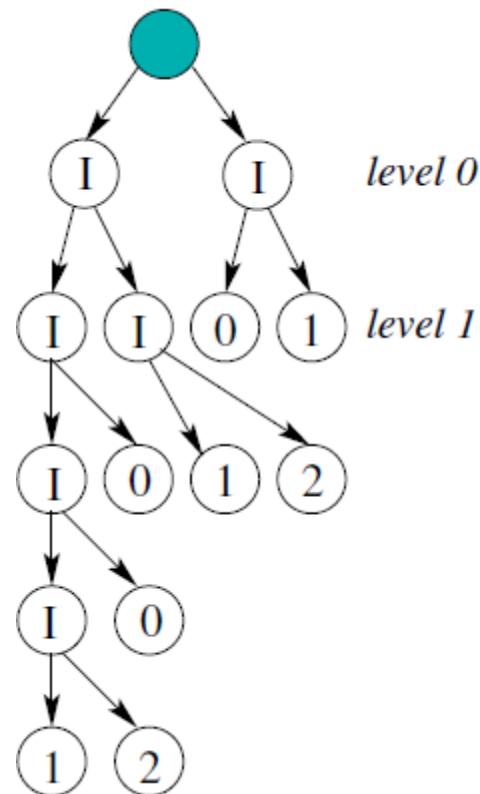
Knuth-Yao random walk: example

Binary tree corresponding to P_{mat}

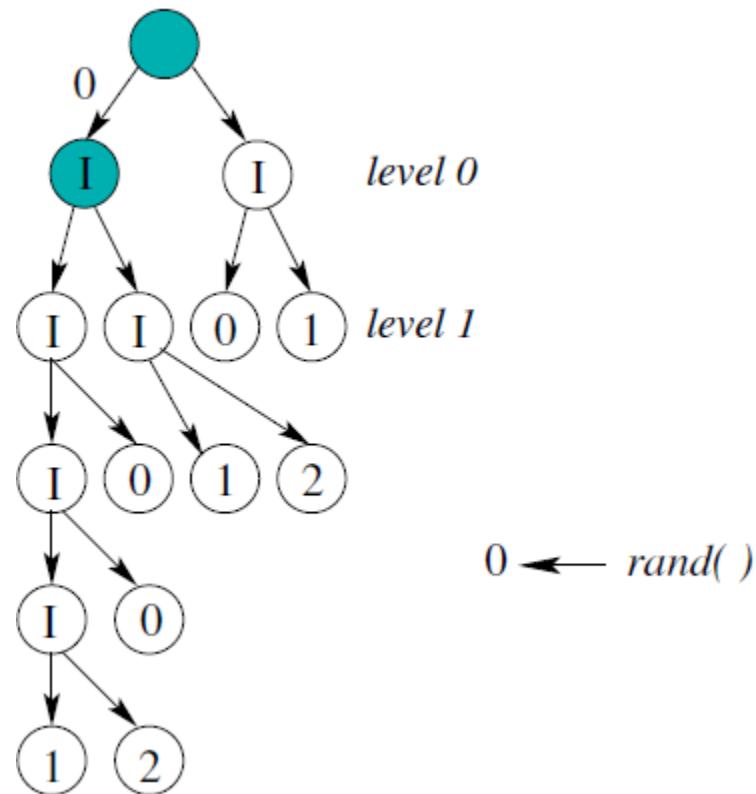


$$P_{mat} = \begin{bmatrix} & & & & \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$$

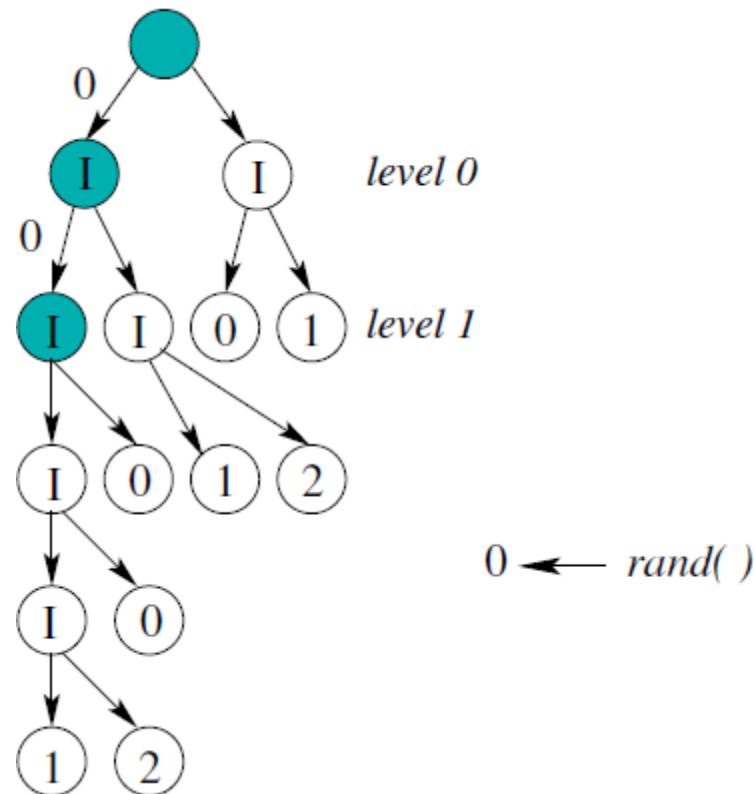
Knuth-Yao random walk: example



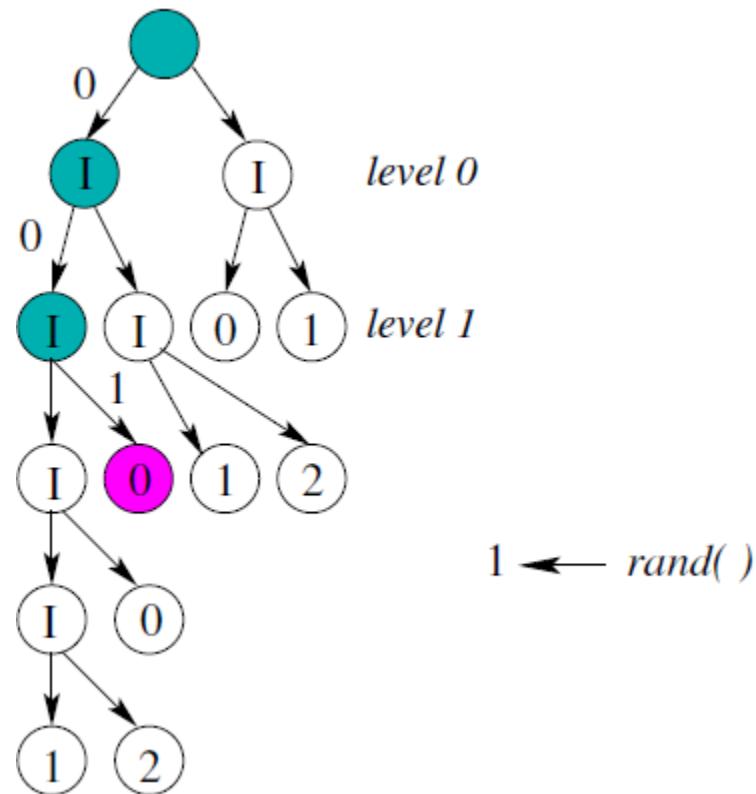
Knuth-Yao random walk: example



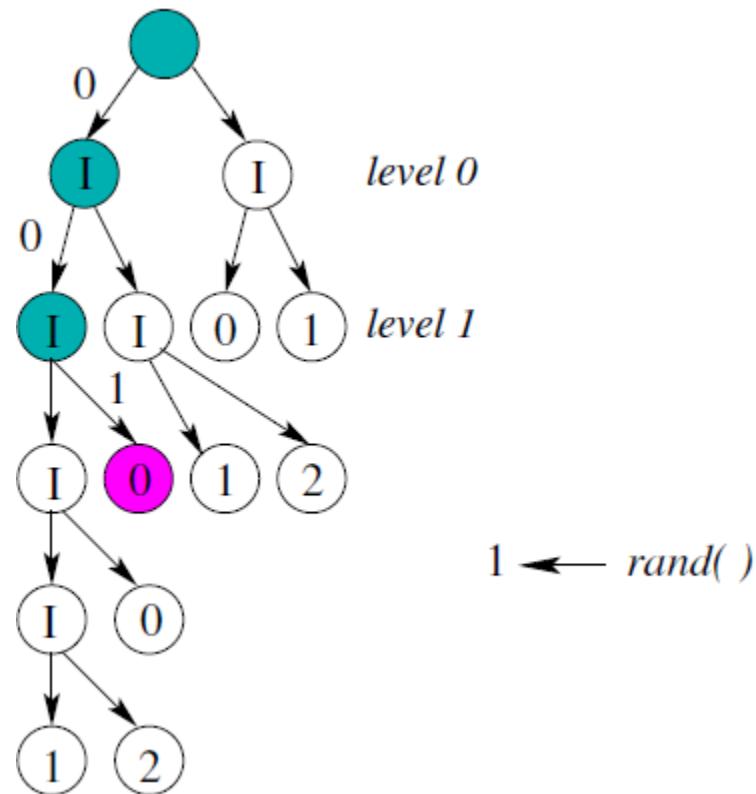
Knuth-Yao random walk: example



Knuth-Yao random walk: example

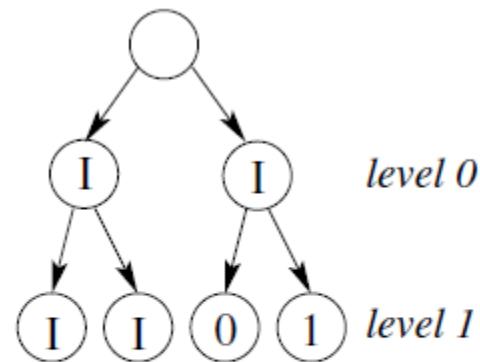


Knuth-Yao random walk: example



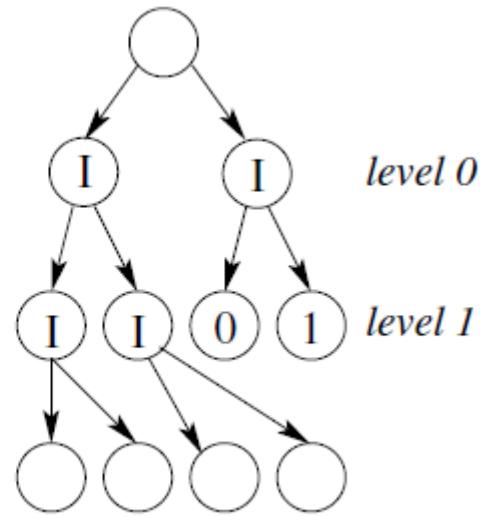
On-the-fly Probability Tree Generation

$$P_{mat} = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$$



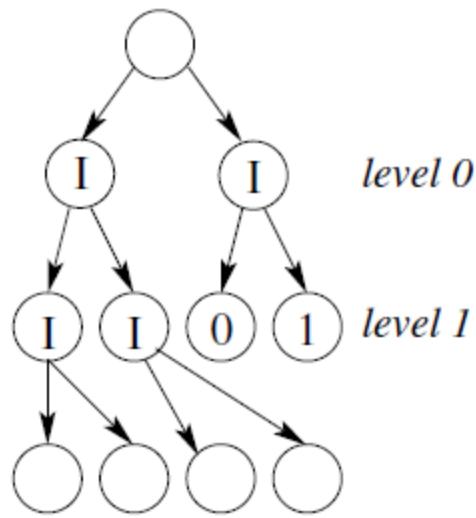
On-the-fly Probability Tree Generation

$$P_{mat} = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$$



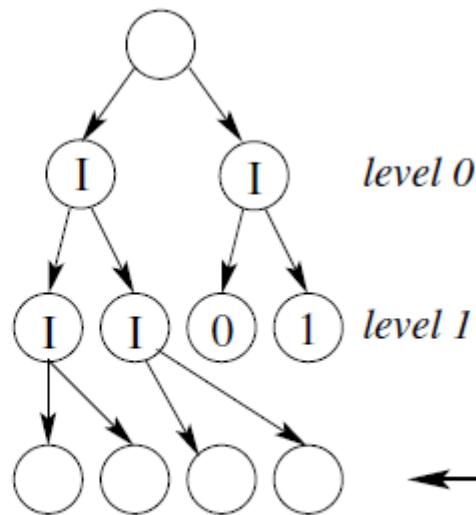
On-the-fly Probability Tree Generation

$$P_{mat} = \begin{bmatrix} 0 & 1 & | & 1 & 1 & 0 \\ 0 & 1 & | & 1 & 0 & 1 \\ 0 & 0 & | & 1 & 0 & 1 \end{bmatrix}$$



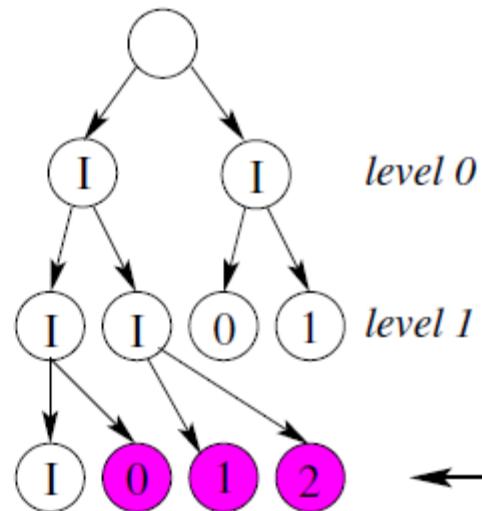
On-the-fly Probability Tree Generation

$$P_{mat} = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$$

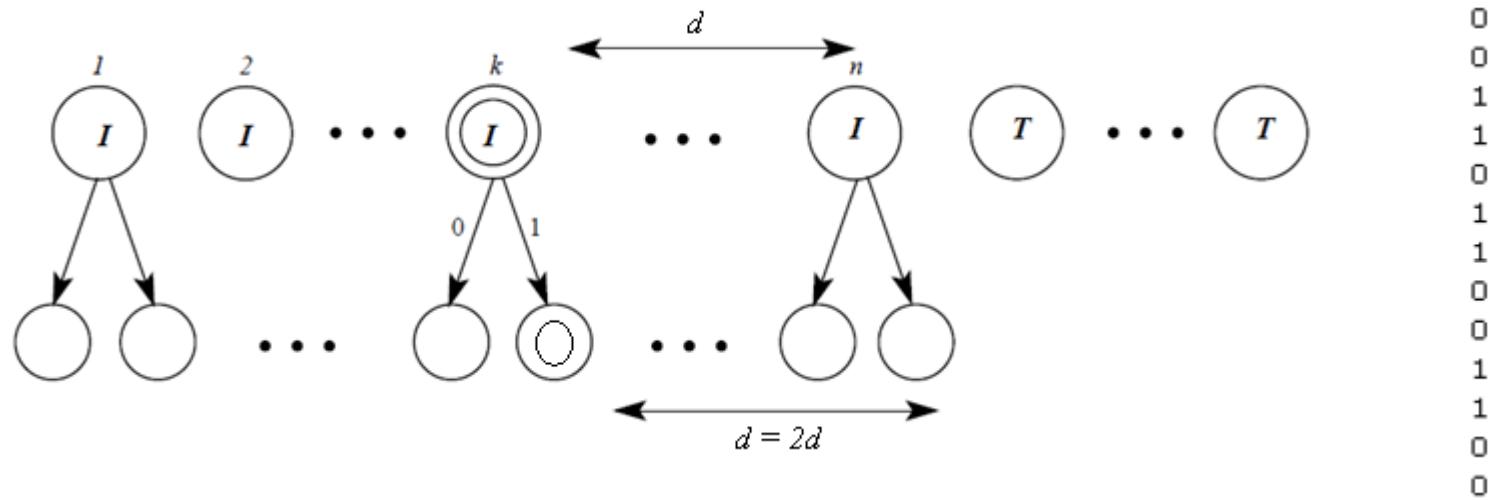
On-the-fly Probability Tree Generation

$$P_{mat} = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$$

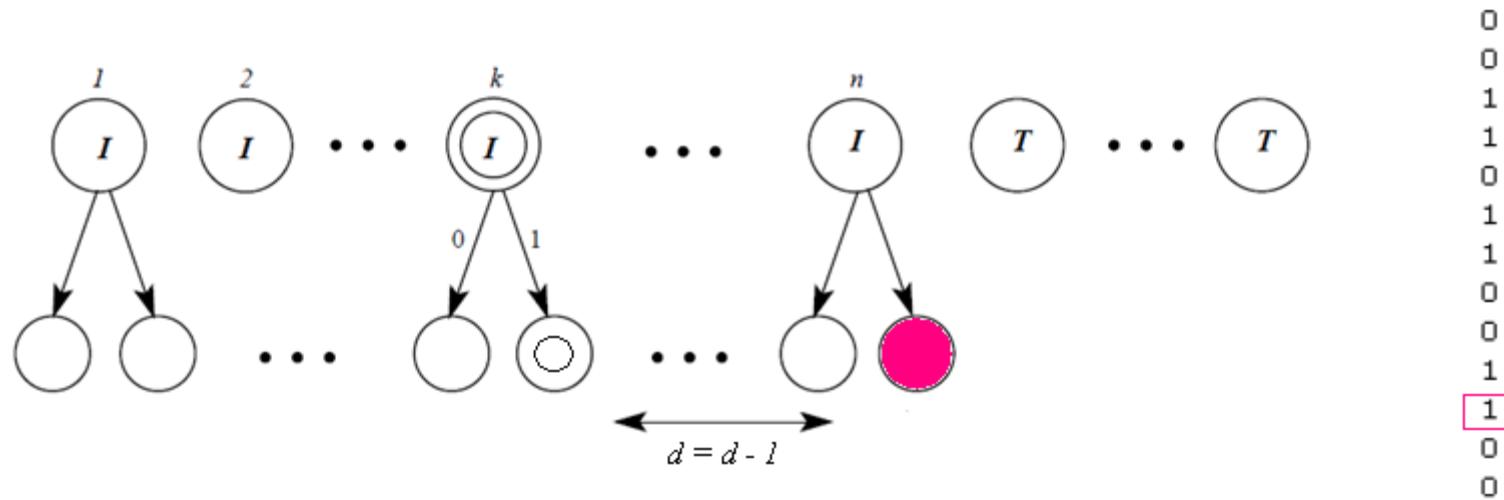
Counter-based algorithm

- Construction of i -th level during sampling : Counter d for distance



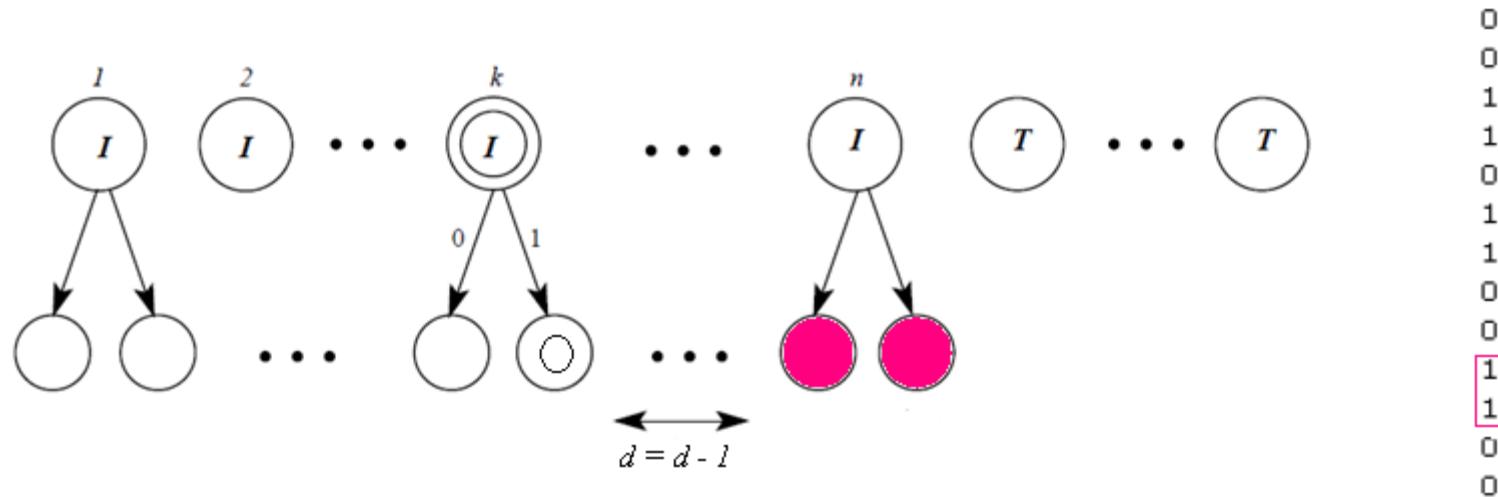
Counter-based algorithm

- Construction of i -th level during sampling : Counter d for distance



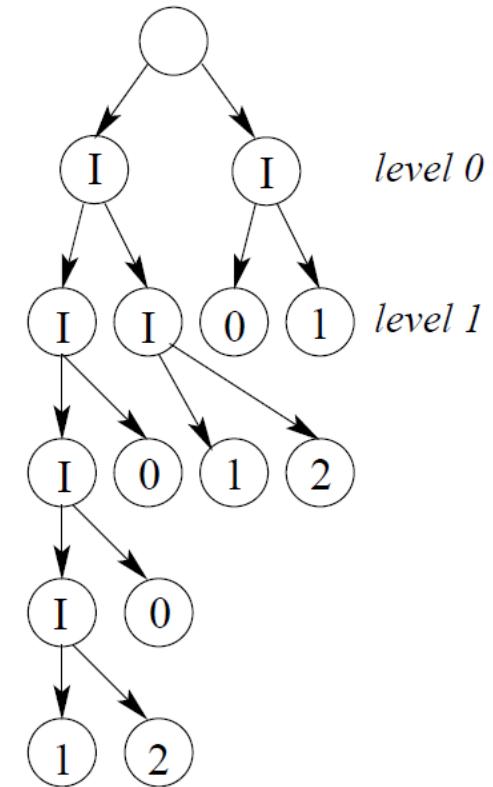
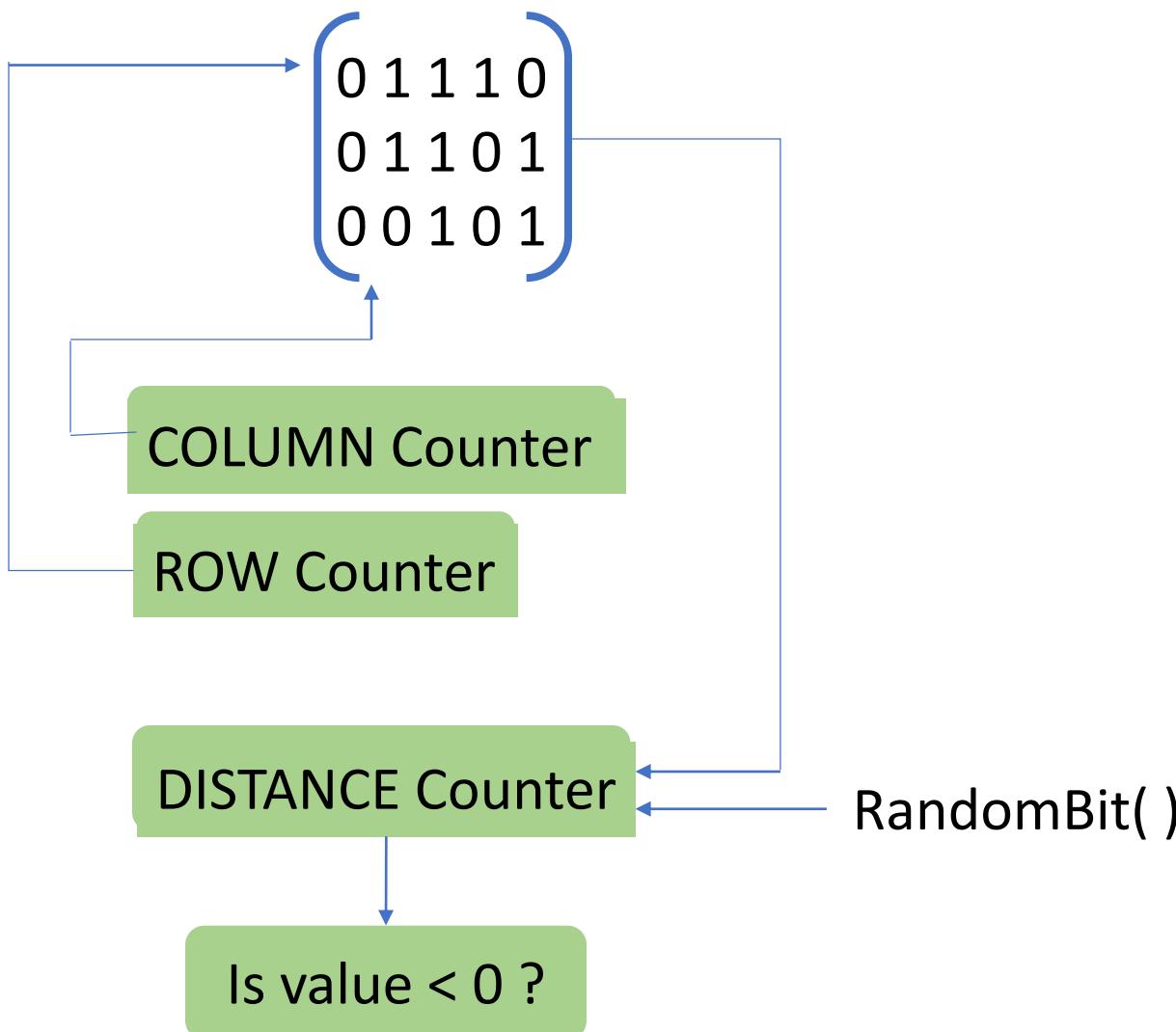
Counter-based algorithm

- Construction of i -th level during sampling : Counter d for distance



- When $d < 0$ for the first time, the visited node is a terminal node
- We need counters for d and row-number

Tree Traversal: Simple Algorithm



Memory optimization: Probability Matrix

Probability Matrix : Column-wise Optimization

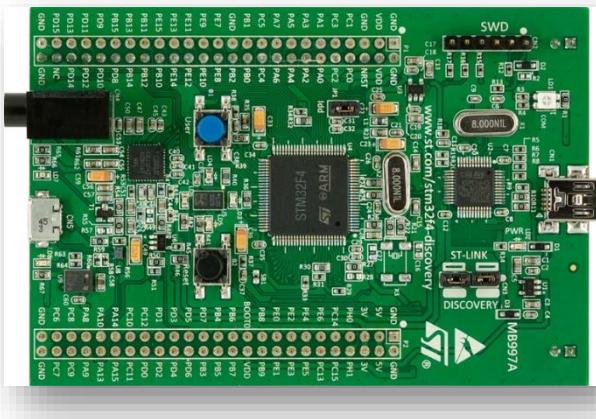
- Observation
 - Difference in length is 1 for most consecutive columns
- One-step difference in column length
 - One bit per differential column-length
 - 1 for increment
 - 0 for no-increment

```
0001111111010111000101110101  
00111001101110110011011001101  
001101001000110011101100011010  
001010010010001110000011001110  
000111010011001101100110100000  
000100101100101100100011010010  
000010101111011110010010001110  
000001011100110110001001011000  
000000101100100010110010101101  
000000010011011000000110100010  
00000000111101001000111111011  
0000000001010111011011001001  
0000000000111000101110001100  
0000000000010000101011010101  
000000000000100011100100010  
000000000000010001001100010001  
000000000000001111000100  
00000000000000010111111111111111
```

Easy to implement in both HW and SW

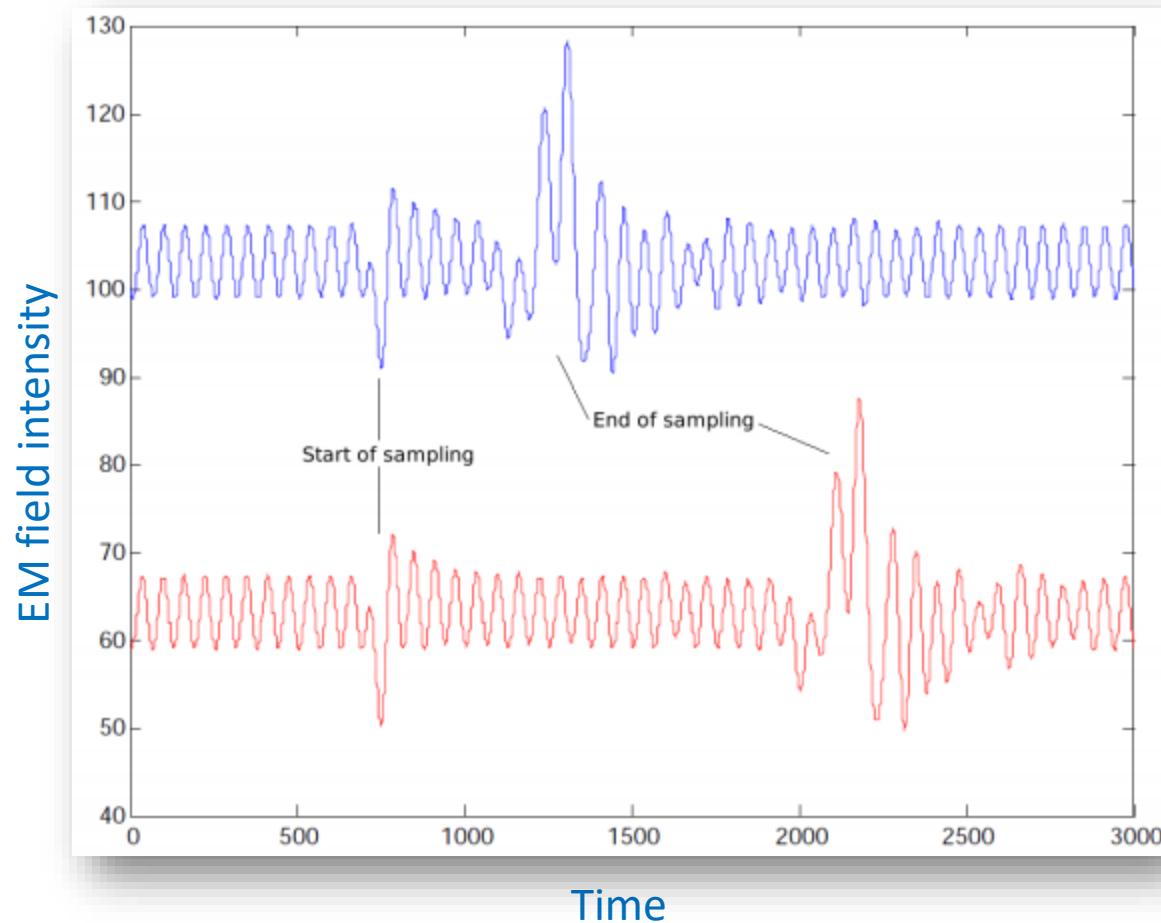
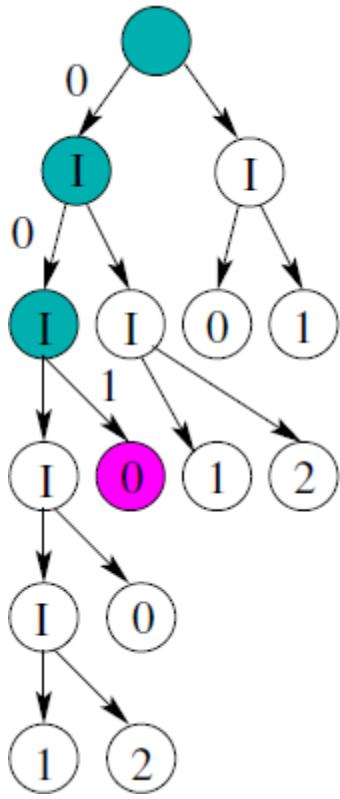


Resources: 1% of smallest FPGA
Cycles per sample (avg.) ~1.5



Memory 0.7 KB
Cycles per sample (avg.) ~28

Variation in execution time

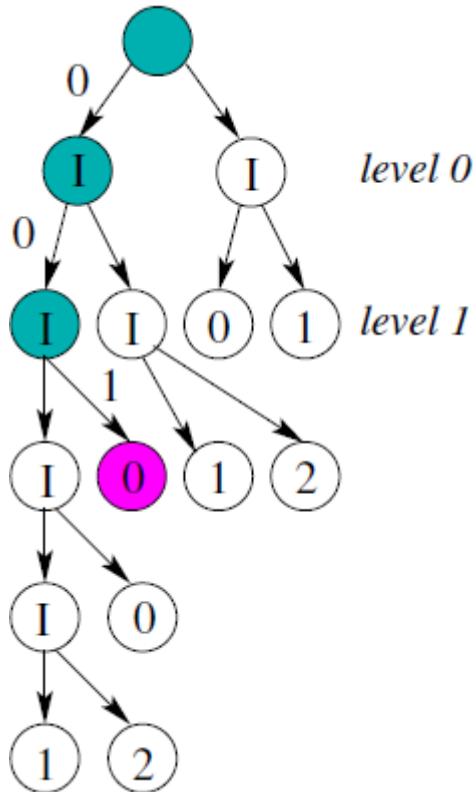


Ideal world: Hard to solve s when $\mathbf{b} = \mathbf{A} \cdot \mathbf{s} + \mathbf{e}$

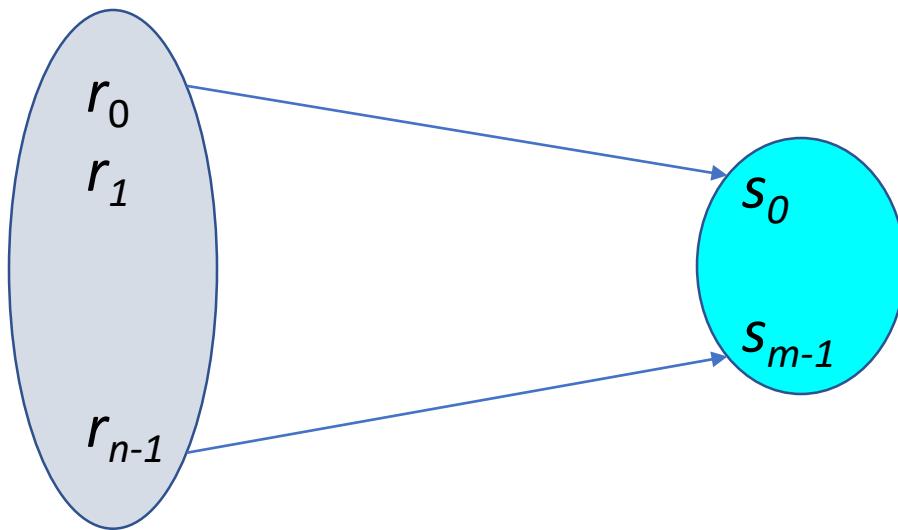
Timing leak \rightarrow Easy to solve s when $\mathbf{b} = \mathbf{A} \cdot \mathbf{s} + \mathbf{e}$

Constant-time Discrete Gaussian sampling

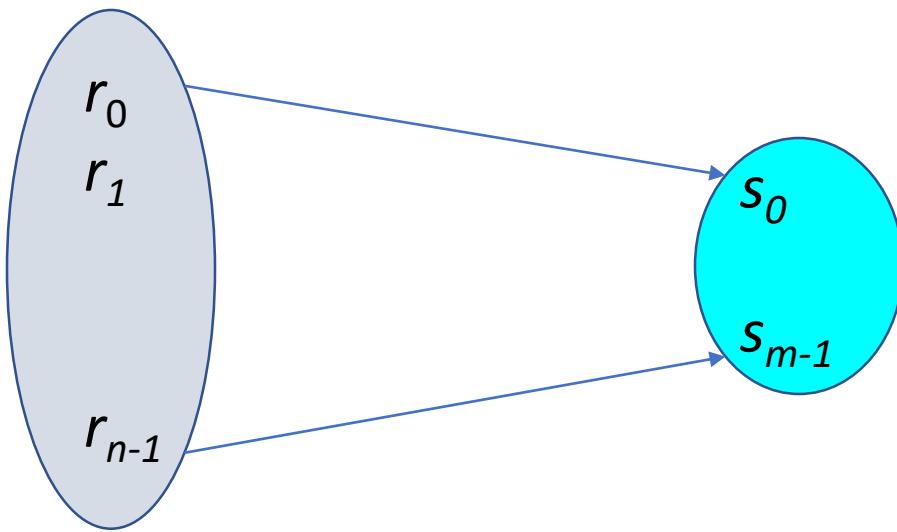
Knuth-Yao random walk: as a mapping



Random string	→	Sample value
001	→	0
010	→	1
011	→	2
...		



Mapping → Boolean equations



$$s_0 = f^0(r_0, r_1, \dots, r_{n-1})$$

...

$$s_{m-1} = f^{m-1}(r_0, r_1, \dots, r_{n-1})$$

}

Needs to be evaluated
in constant time

Constant-time evaluation

E.g. $s_0 = r_0 \wedge (r_1 \vee r_2).$

Result is 0 when $r_0 = 0$ → Non-constant time

Constant-time evaluation

E.g. $s_0 = r_0 \wedge (r_1 \vee r_2).$

Result is 0 when $r_0 = 0$ → Non-constant time

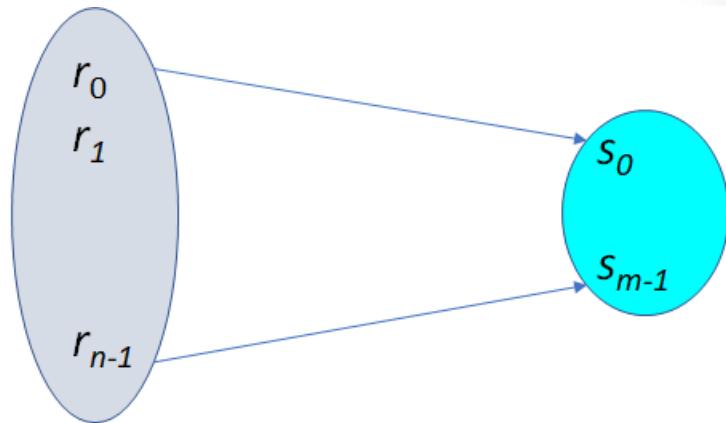
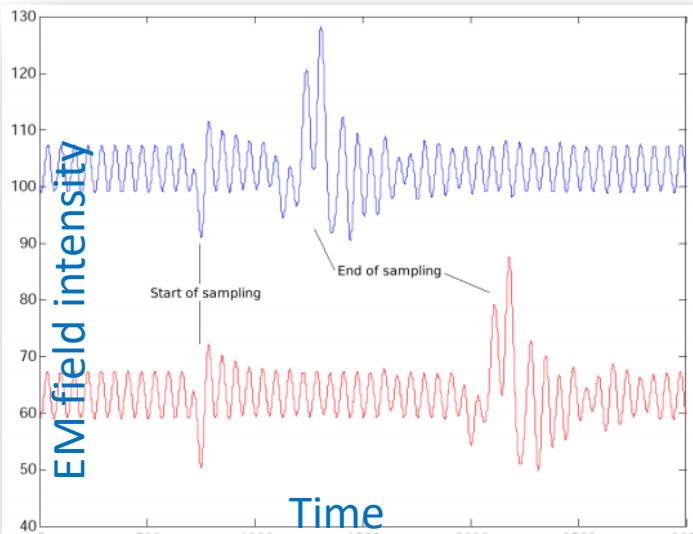
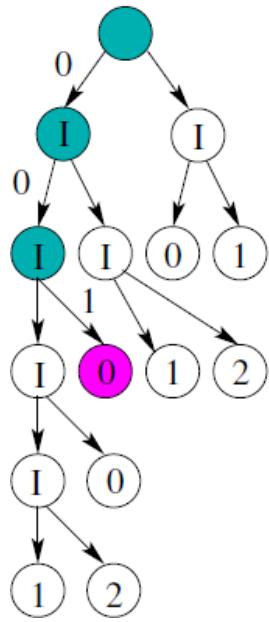
Bit-slicing

- constant time execution and performance
- 64-bit computer architecture
- 64-bit random integers R_0, R_1 , and R_2

$$S_0 = R_0 \wedge (R_1 \vee R_2).$$

Parallel computation of 64 output bits in constant-time

Constant-time Gaussian Sampling

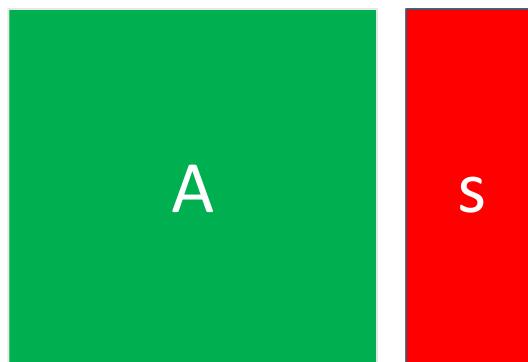


$$s_0 = f^0(r_0, r_1, \dots, r_{n-1})$$
$$\dots$$
$$s_{m-1} = f^{m-1}(r_0, r_1, \dots, r_{n-1})$$

**Convolution +
Boolean Decomp
+ Bitslicing**

Followup → Const. time *Falcon* signature becomes 33% slower

- Error sampling 
- **Matrix-vector multiplication**



Complexity $O(n^2)$

Performance of FRODO KEM

- FRODO KEM uses matrix dimension
 - $n=640$ for 100 bit quantum security
 - $n=976$ for 150 bit quantum security

FRODO KEM ($n=640$) on ARM Cortex M4

Key gen	Encapsulation	Decapsulation
81 M	86 M	87 M

- Roughly 3.3 sec at 24 MHz
- Slow due to expensive matrix-vector multiplications

Standard LWE

$$\begin{pmatrix} a_{0,0} & a_{0,1} & a_{0,2} & a_{0,3} \\ a_{1,0} & a_{1,1} & a_{1,2} & a_{1,3} \\ a_{2,0} & a_{2,1} & a_{2,2} & a_{2,3} \\ a_{3,0} & a_{3,1} & a_{3,2} & a_{3,3} \end{pmatrix} * \begin{pmatrix} s_0 \\ s_1 \\ s_2 \\ s_3 \end{pmatrix} + \begin{pmatrix} e_0 \\ e_1 \\ e_2 \\ e_3 \end{pmatrix} \approx \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \end{pmatrix} \pmod{q}$$

Uniformly random matrix

Special LWE

$$\begin{pmatrix} a_0 & -a_3 & -a_2 & -a_1 \\ a_1 & a_0 & -a_3 & -a_2 \\ a_2 & a_1 & a_0 & -a_3 \\ a_3 & a_2 & a_1 & a_0 \end{pmatrix} * \begin{pmatrix} s_0 \\ s_1 \\ s_2 \\ s_3 \end{pmatrix} + \begin{pmatrix} e_0 \\ e_1 \\ e_2 \\ e_3 \end{pmatrix} \approx \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \end{pmatrix} \pmod{q}$$

Matrix from uniformly random vector

Special LWE: known as Ring-LWE

$$\begin{pmatrix} a_0 & -a_3 & -a_2 & -a_1 \\ a_1 & a_0 & -a_3 & -a_2 \\ a_2 & a_1 & a_0 & -a_3 \\ a_3 & a_2 & a_1 & a_0 \end{pmatrix} * \begin{pmatrix} s_0 \\ s_1 \\ s_2 \\ s_3 \end{pmatrix} + \begin{pmatrix} e_0 \\ e_1 \\ e_2 \\ e_3 \end{pmatrix} \approx \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \end{pmatrix} \pmod{q}$$

Matrix from uniformly random vector



$$a(x) * s(x) + e(x) \approx b(x) \pmod{q} \pmod{x^4 + 1}$$

where

$$a(x) = (a_0 + a_1x + a_2x^2 + a_3x^3)$$

$$s(x) = (s_0 + s_1x + s_2x^2 + s_3x^3)$$

$$e(x) = (e_0 + e_1x + e_2x^2 + e_3x^3)$$

$$b(x) = (b_0 + b_1x + b_2x^2 + b_3x^3)$$

Efficient Polynomial multiplication

Polynomial multiplication: methods

- Schoolbook multiplication: complexity n^2
- Karatsuba multiplication: complexity $n^{1.585}$
- FFT-based multiplication: **complexity $(n \log n)$**

FFT = Fast Fourier Transform

Polynomial multiplication: FFT-based

$$a(x) = a_{n-1}x^{n-1} + \dots + a_1x + a_0$$

a_{n-1}	\dots	a_1	a_0
-----------	---------	-------	-------

0	\dots	0	a_{n-1}	\dots	a_1	a_0
---	---------	---	-----------	---------	-------	-------



2n-FFT



A_{2n-1}		A_n	A_{n-1}		A_1	A_0
------------	--	-------	-----------	--	-------	-------

Polynomial multiplication: FFT-based

$$b(x) = b_{n-1}x^{n-1} + \dots + b_1x + b_0$$

b_{n-1}	...	b_1	b_0
-----------	-----	-------	-------

0	...	0	b_{n-1}	...	b_1	b_0
---	-----	---	-----------	-----	-------	-------

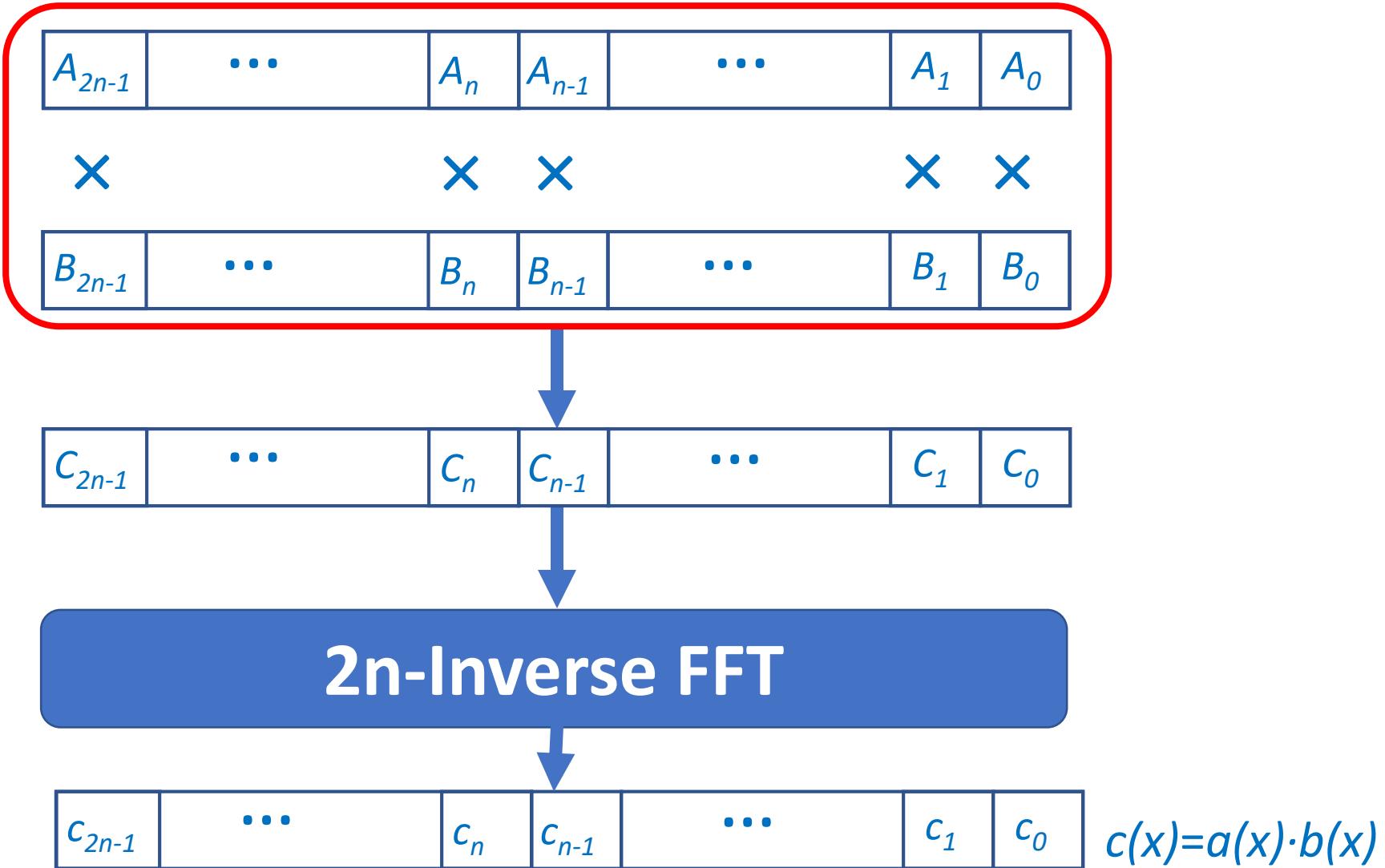


2n-FFT



B_{2n-1}		B_n	B_{n-1}		B_1	B_0
------------	--	-------	-----------	--	-------	-------

Polynomial multiplication: FFT-based



FFT to NTT

- FFT involves real numbers
- Number Theoretic Transform (**NTT**)
 - is a generalization of FFT
 - Only integer arithmetic modulo q

Requirements: $\mathbb{Z}_q[x]/(x^n + 1)$

- $n = 2^k$
- Prime q with $q \equiv 1 \pmod n$

FFT to NTT

- FFT involves real numbers
- Number Theoretic Transform (**NTT**)
 - is a generalization of FFT
 - Only integer arithmetic modulo q

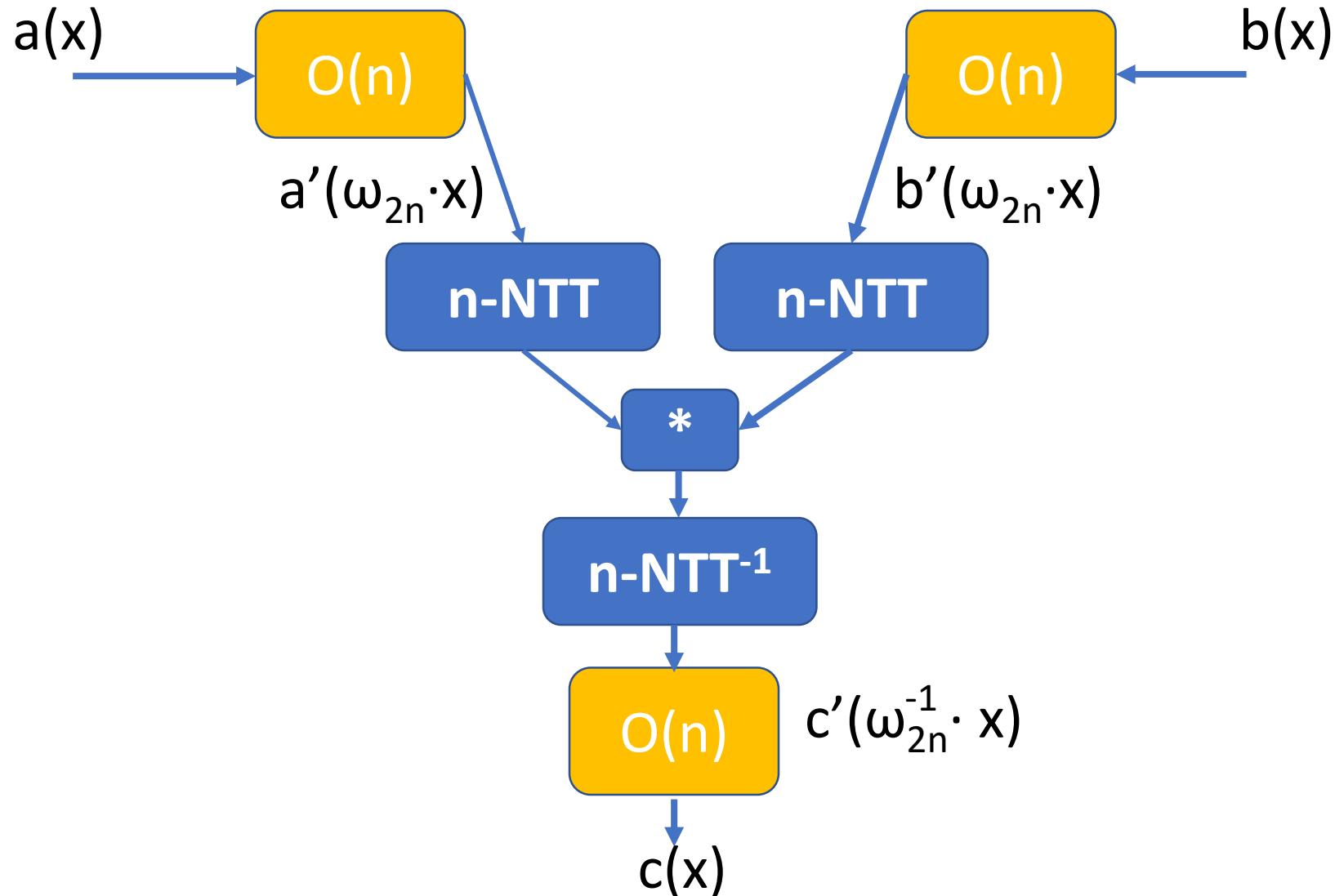
Requirements: $\mathbb{Z}_q[x]/(x^n + 1)$

- $n = 2^k$
- Prime q with $q \equiv 1 \pmod{2n}$

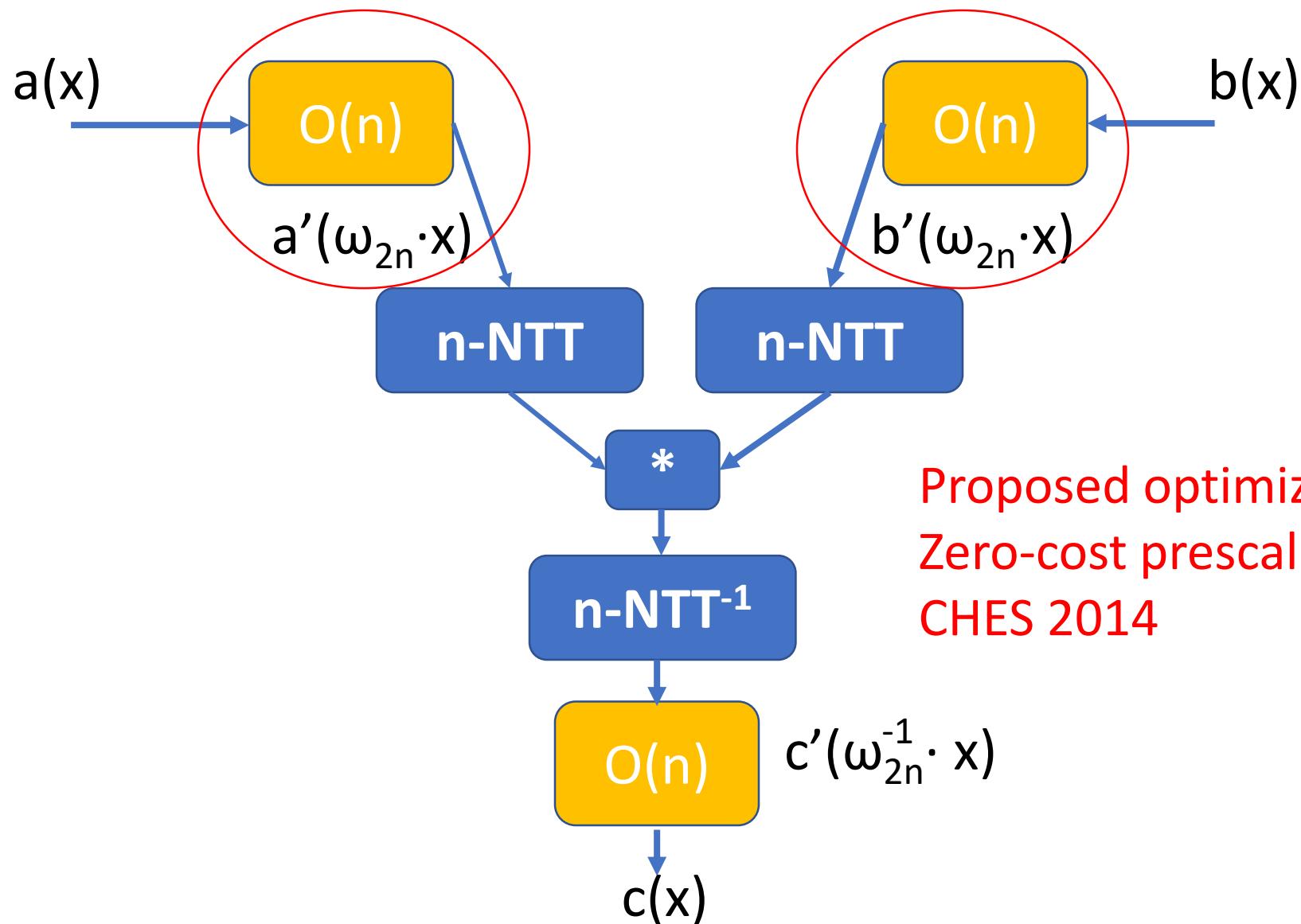
'Negative-wrapped convolution'

Requires n-NTT
instead of 2n-NTT

Special optimization: scaling overhead



Special optimization: scaling overhead



Memory

Simplified NTT loops



```
for (m=2; m<=n; m=2m)
{
    for (j=0; j<=m/2-1; j++)
    {
        for (k=0; j<n; k=k+m)
        {
            index = f(m, j, k);
            Butterfly(A[index], A[index+m/2]);
        }
    }
}
```

Memory

Simplified NTT loops

A[n-1]

A[n-2]

A[3]

A[2]

A[1]

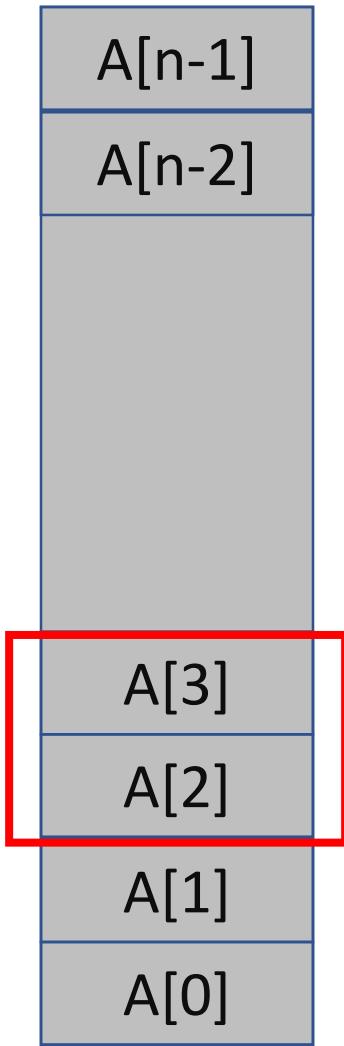
A[0]

```
for (m=2; m<=n; m=2m)
{
    for (j=0; j<=m/2-1; j++)
    {
        for (k=0; j<n; k=k+m)
        {
            index = f(m, j, k);
            Butterfly(A[index], A[index+m/2]);
        }
    }
}
```

NTT starts with $m=2$

Butterfly(A[0], A[1])

Memory

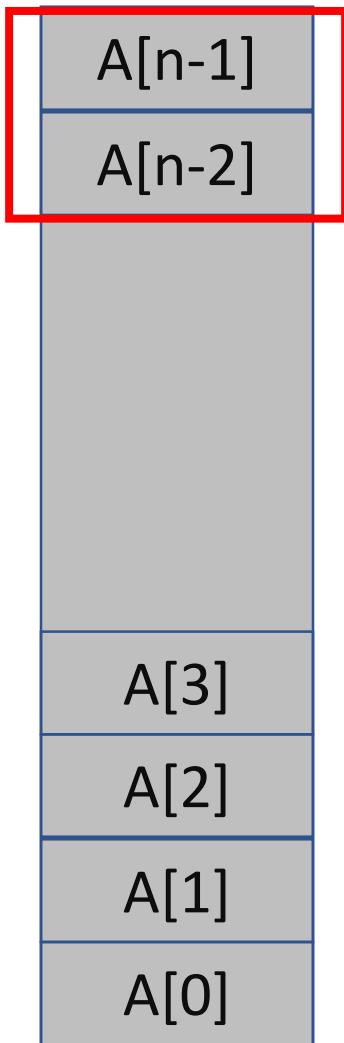


Simplified NTT loops

```
for (m=2; m<=n; m=2m)
{
    for (j=0; j<=m/2-1; j++)
    {
        for (k=0; j<n; k=k+m)
        {
            index = f(m, j, k);
            Butterfly(A[index], A[index+m/2]);
        }
    }
}
```

NTT starts with $m=2$
Butterfly($A[2]$, $A[3]$)

Memory



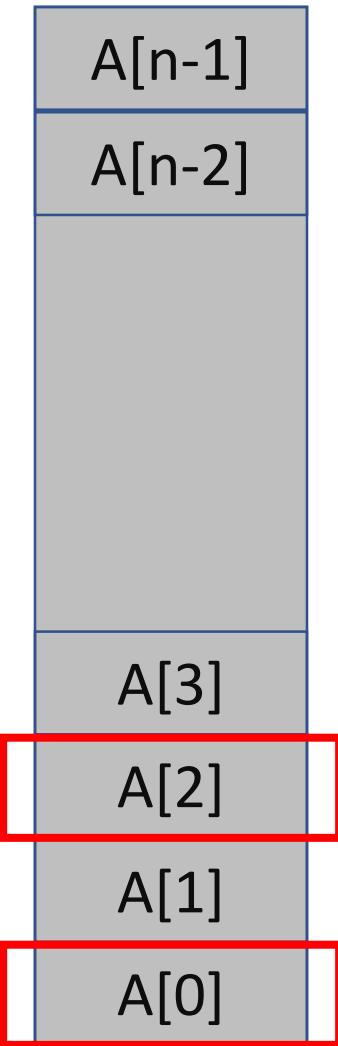
Simplified NTT loops

```
for (m=2; m<=n; m=2m)
{
    for (j=0; j<=m/2-1; j++)
    {
        for (k=0; j<n; k=k+m)
        {
            index = f(m, j, k);
            Butterfly(A[index], A[index+m/2]);
        }
    }
}
```

NTT starts with $m=2$

Butterfly(A[n-2], A[n-1])

Memory



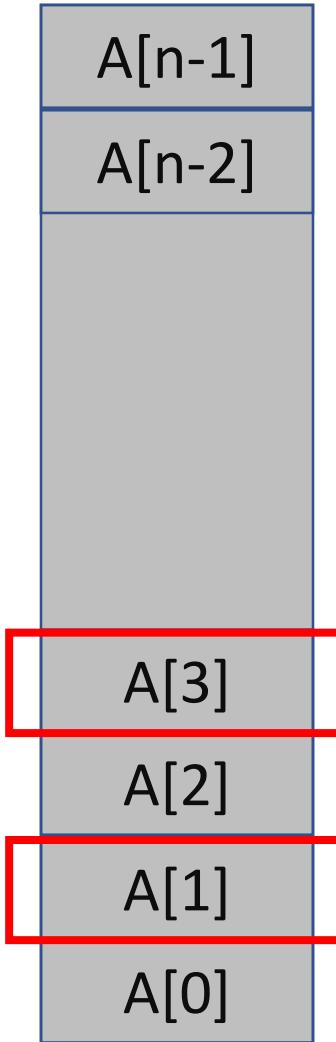
Simplified NTT loops

```
for (m=2; m<=n; m=2m)
{
    for (j=0; j<=m/2-1; j++)
    {
        for (k=0; j<n; k=k+m)
        {
            index = f(m, j, k);
            Butterfly(A[index], A[index+m/2]);
        }
    }
}
```

Next, m increments to **m=4**.
Butterfly(A[0], A[2]), Butterfly(A[4], A[6]) ...

Memory

Simplified NTT loops



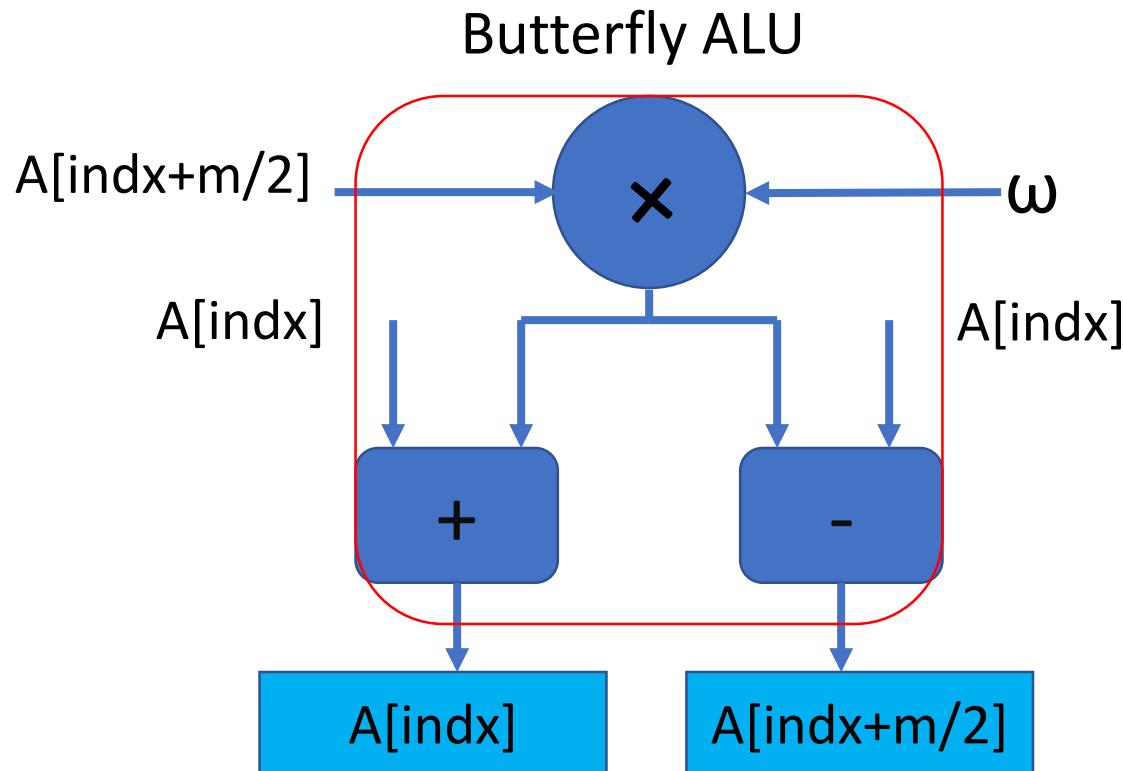
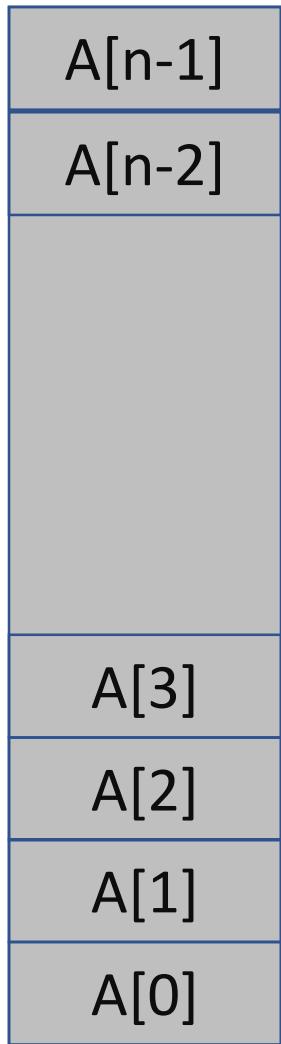
```
for (m=2; m<=n; m=2m)
{
    for (j=0; j<=m/2-1; j++)
    {
        for (k=0; j<n; k=k+m)
        {
            index = f(m, j, k);
            Butterfly(A[index], A[index+m/2]);
        }
    }
}
```

Next, m increments to **m=4**.

Butterfly(A[1], A[3]), Butterfly(A[5], A[7]) ...

Memory access problem

Memory is sequential



1R + 1R + 1C + 1W + 1W



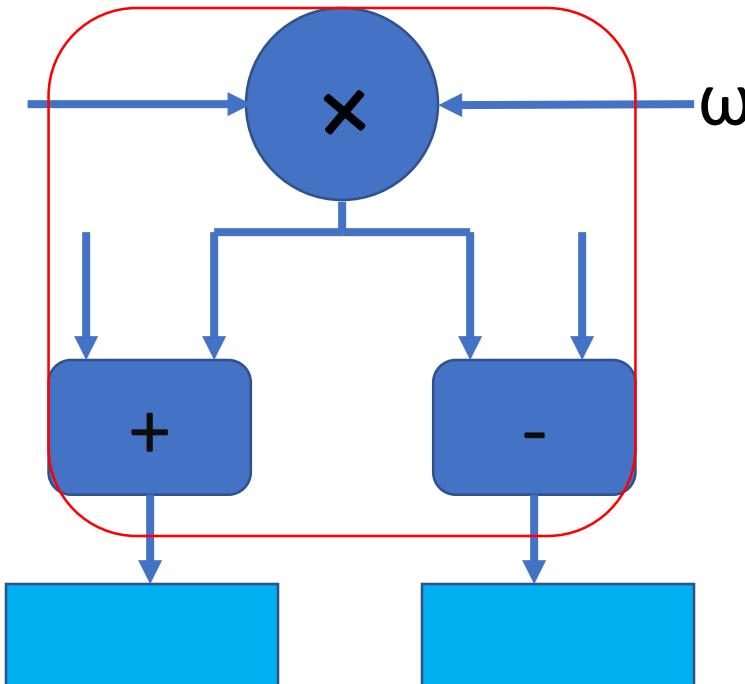
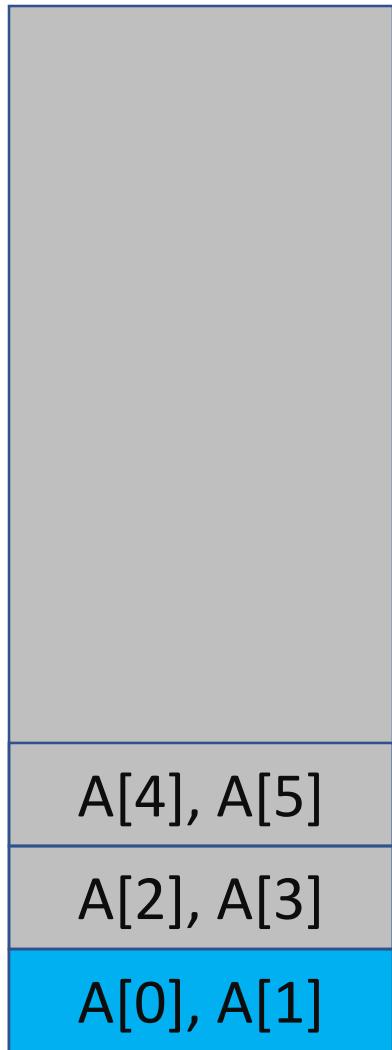
- If we have small coefficients
- Two coefficients in a single word?

A[4], A[5]

A[2], A[3]

A[0], A[1]

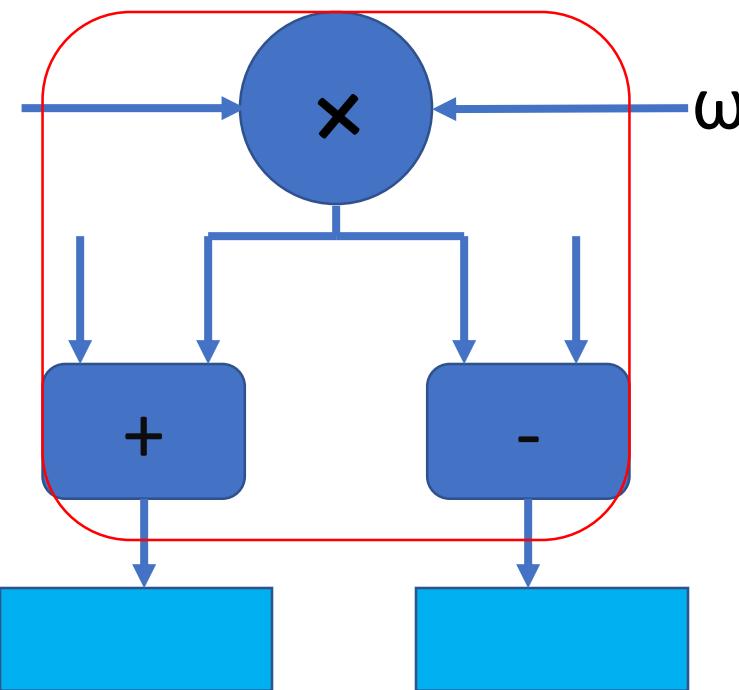
Observation



During NTT loop $m = 2$

1. Read {A[0], A[1]} [1 Cycle]
2. Butterfly(A[0], A[1]) [1 Cycle]
3. Write {A[0], A[1]} [1 Cycle]

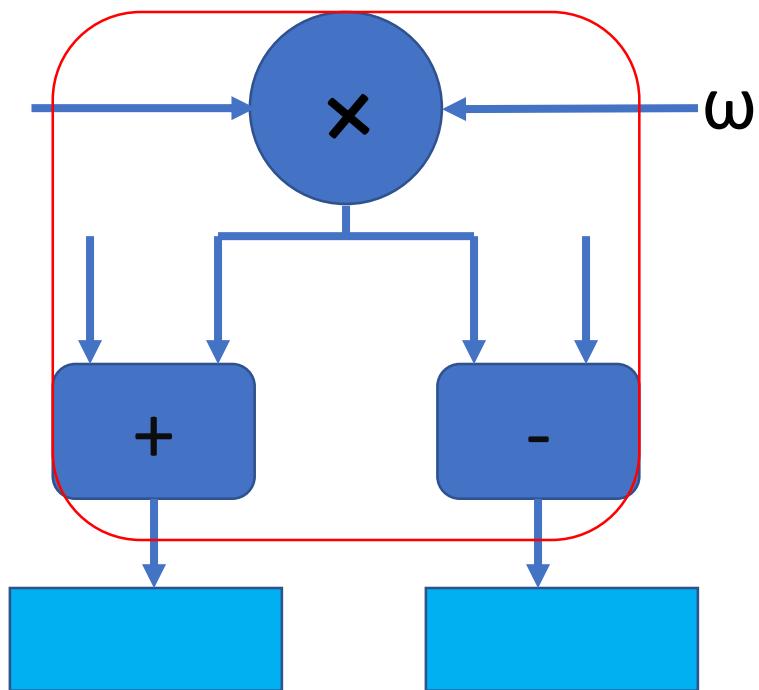
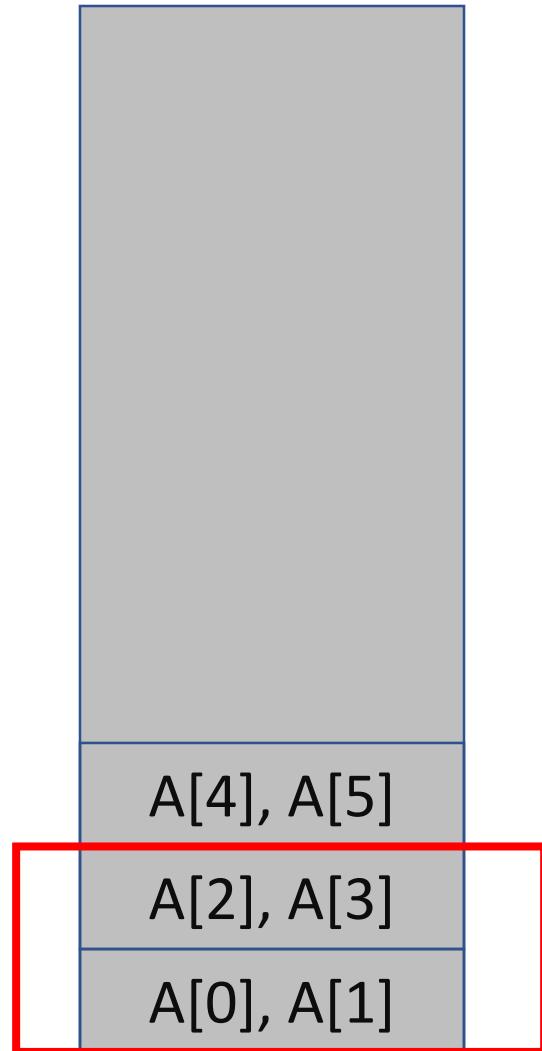
Observation



Problem happens next, when $m=4$.
Butterfly(A[0], A[2]), Butterfly(A[4], A[6]) ...

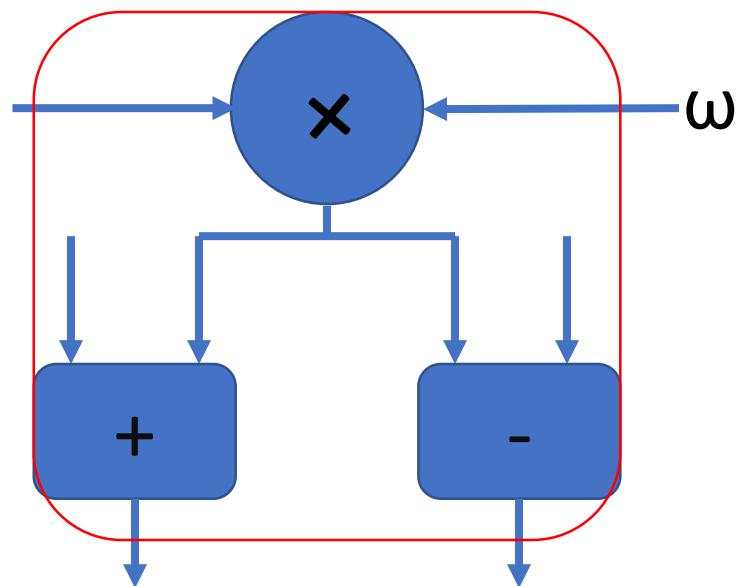
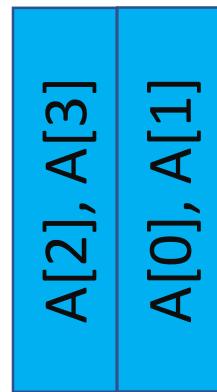
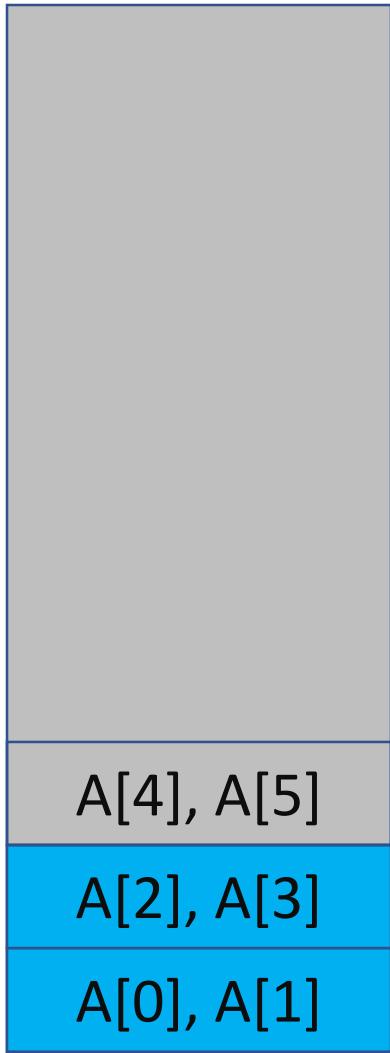
Solution

Process 4 coefficients together



Solution

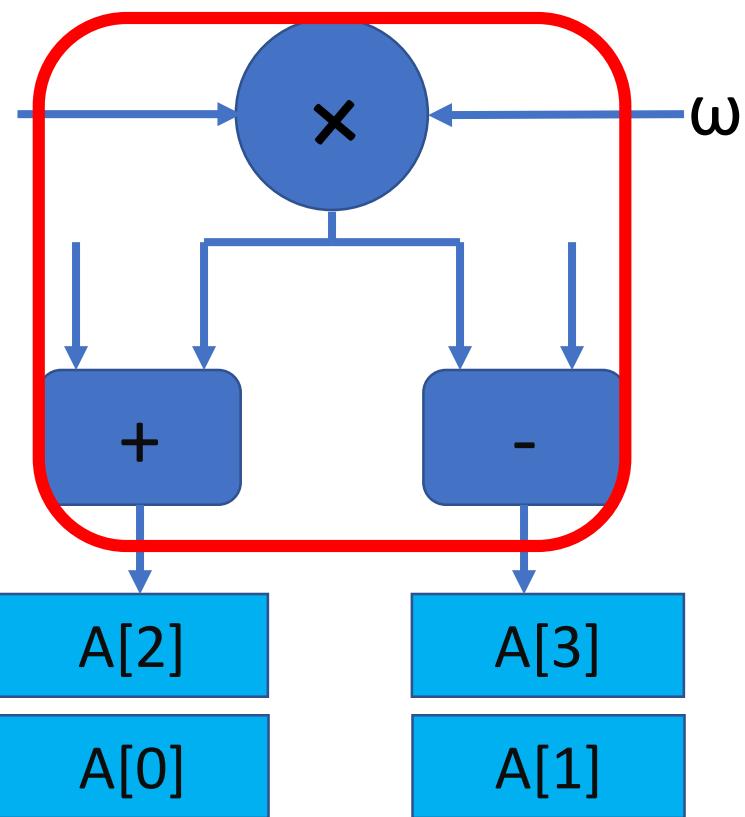
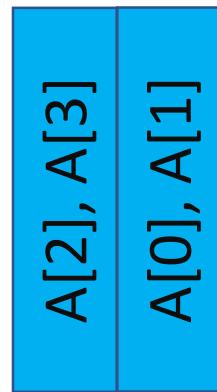
Process 4 coefficients together



2R

Solution

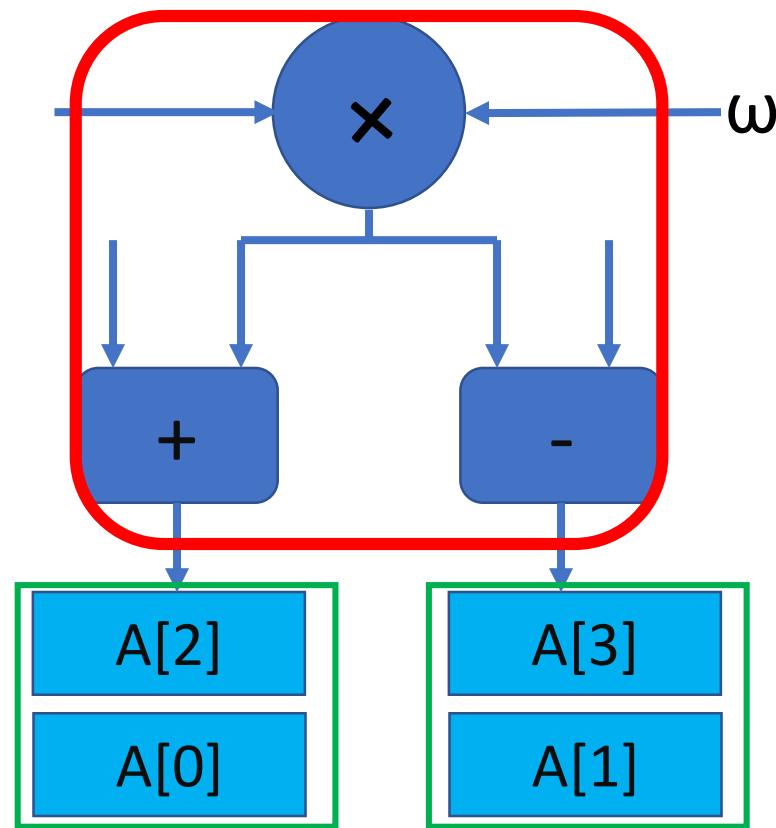
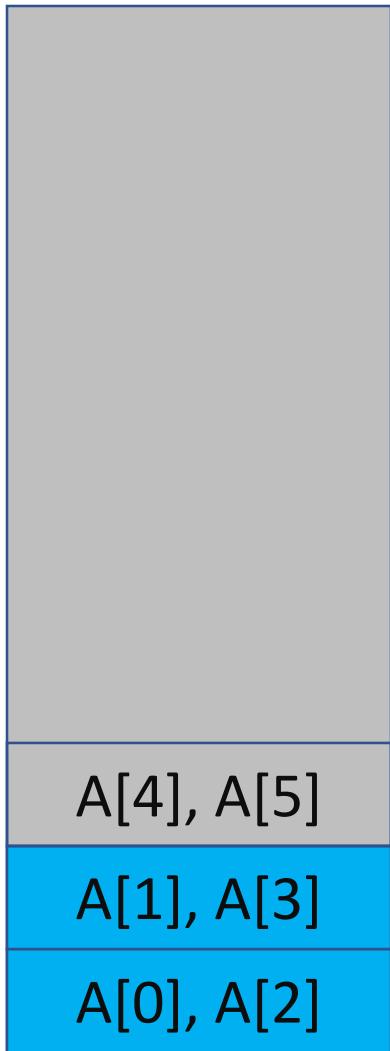
Process 4 coefficients together



$2R + 2C$

Solution

Process 4 coefficients together



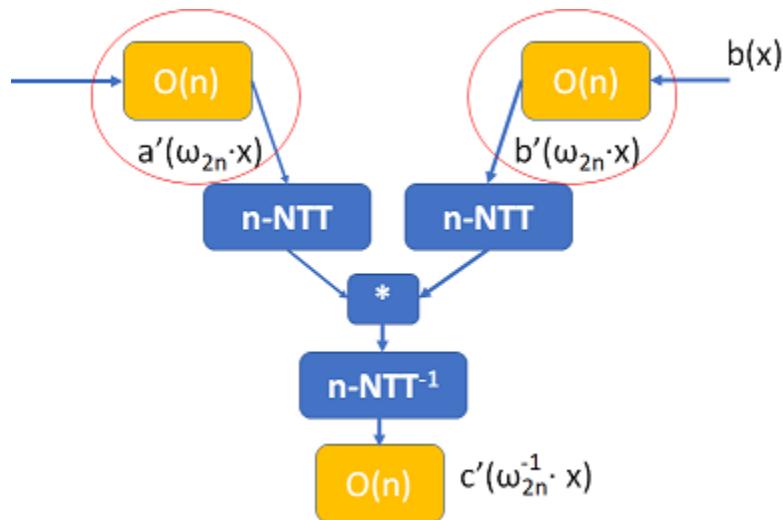
$$2R + 2C + 2W$$

Results

Before $O(n) + O(n) + O(n \log n) + O(n)$

Results

Before $O(n) + O(n) + O(n \log n) + O(n)$

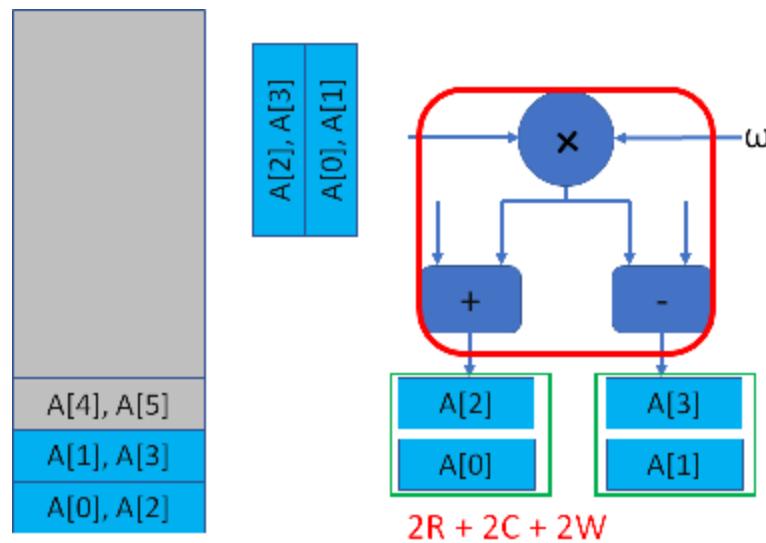


Optimization 1
Zero-cost
prescaling

Now $\cancel{O(n)} + \cancel{O(n)} + O(n \log n) + O(n)$

Results

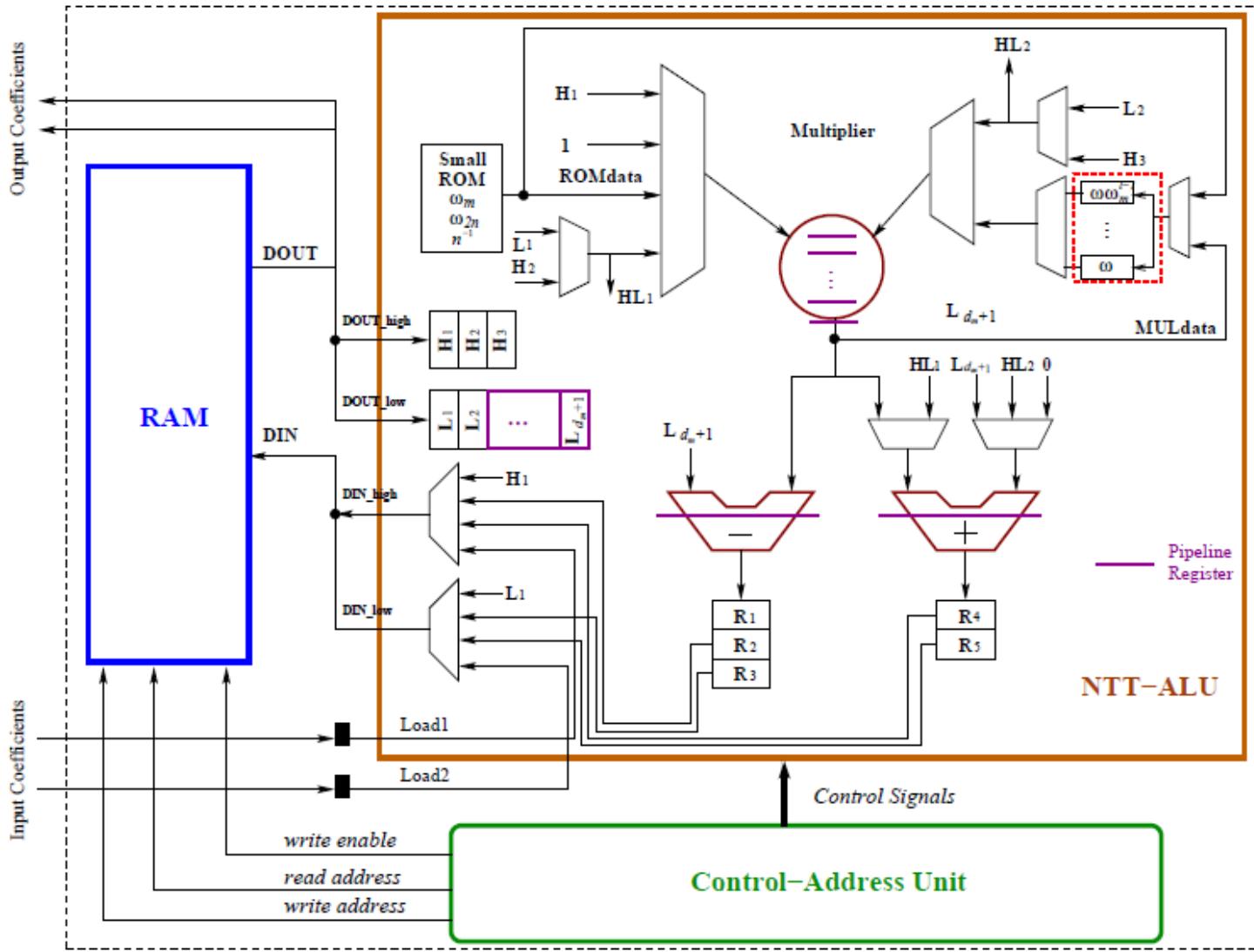
Before $O(n) + O(n) + O(n \log n) + O(n)$



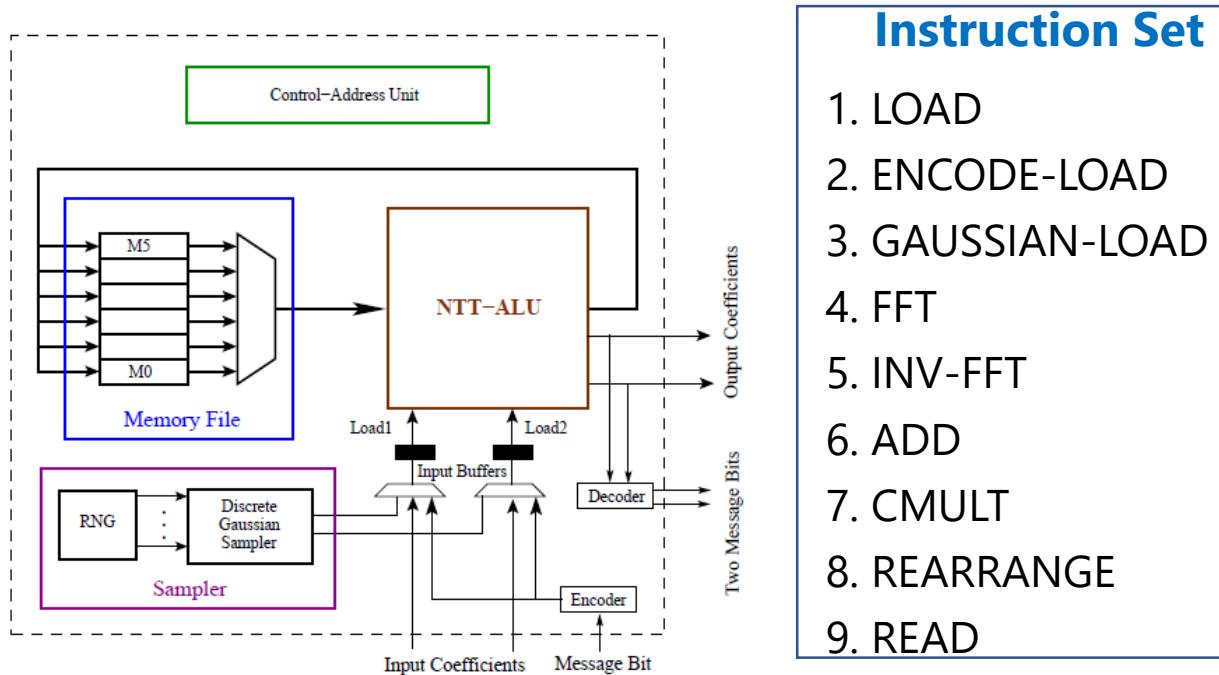
Optimization 2
Memory access
reduction

Now $O(n) + O(n) + \frac{1}{2}O(n \log n) + O(n)$

Architecture of NTT-based polynomial multiplier



Lattice-based public-key instruction-set encryption processor



Publication: CHES 2014

Instruction-set ring-LWE cryptoprocessor

Throughput:

50,000 encryption/sec

100,000 decryption/sec

CHES 2014



Area: 1.3K LUT, 860 FF,
1 DSPMULT, 2 BRAM18

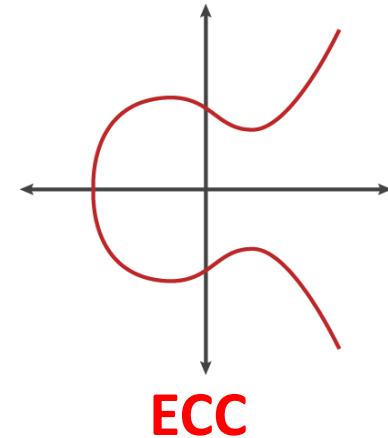
Throughput:

40,000 encryption/sec <

80,000 decryption/sec <

CHES 2012

Area: 18349 LUT, 5644 FF



Ring-LWE encryption: followup works

Software implementations

- On 32-bit ARM

R. de Clercq, S. Sinha Roy, F. Vercauteren, I. Verbauwhede,
"Efficient software implementation of ring-LWE encryption", DATE 2015

Encryptions 121,166 cycles

Decryptions 43,324 cycles

Orders of magnitude
faster than ECC

- On 8-bit AVR

Z. Liu, H. Seo, S. Sinha Roy, J. Großschädl, H. Kim, I. Verbauwhede,
"Efficient Ring-LWE Encryption on 8-Bit AVR Processors", CHES 2015

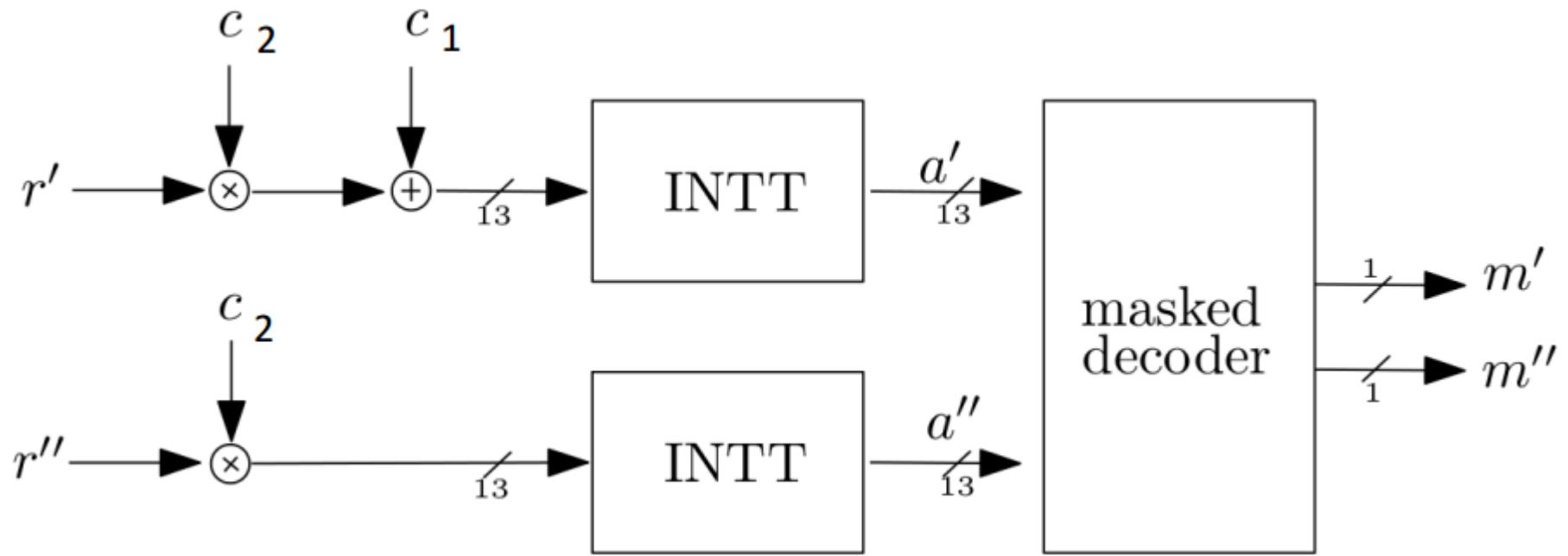
Encryptions 671,628 cycles

Decryptions 275,646 cycles

Ring-LWE encryption: followup works

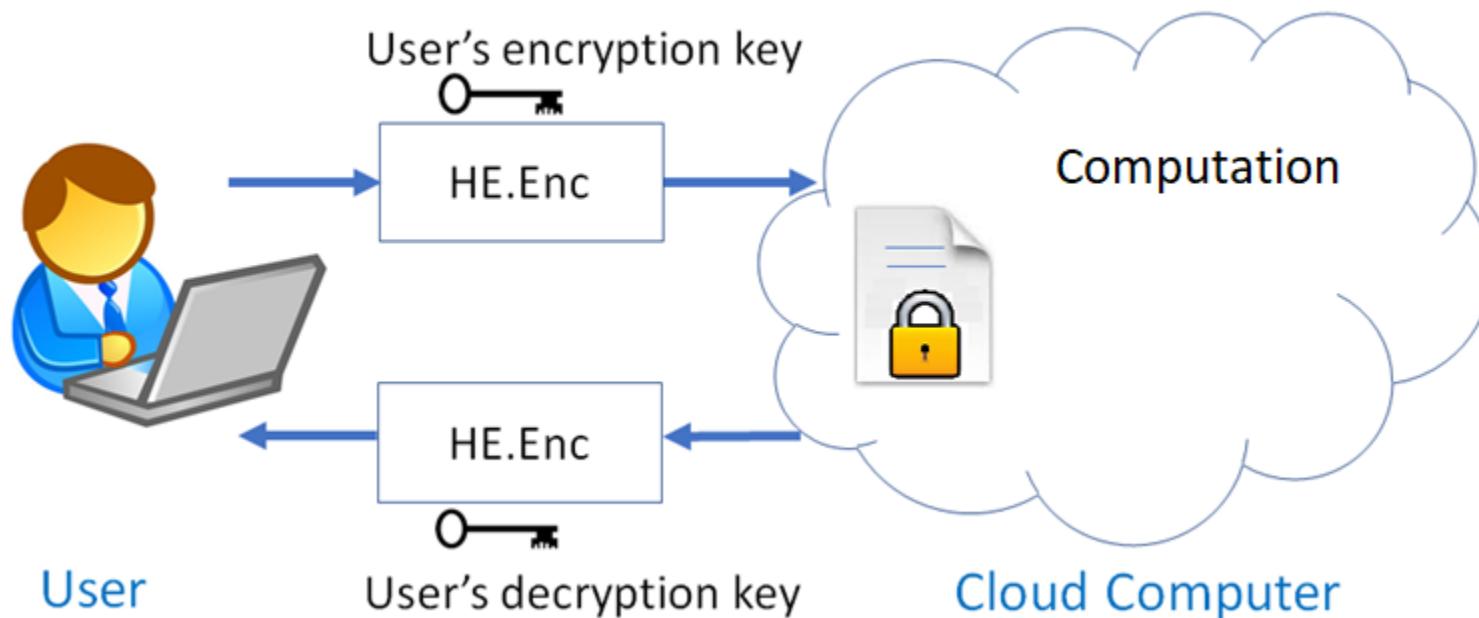
Side channel security: masking scheme

O. Reparaz, S. Sinha Roy, F. Vercauteren, I. Verbauwhede,
"A masked ring-LWE implementation", in CHES 2015



Hardware accelerators for Homomorphic Computation

Homomorphic computation



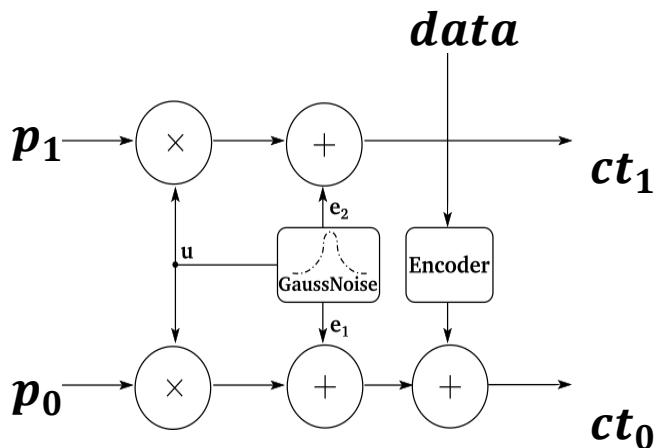
Interesting applications :

- Machine learning on encrypted data
- Prediction from consumption data in smart electricity meters
- Health-care applications
- Encrypted web-search engine

Lattice-based Homomorphic encryption

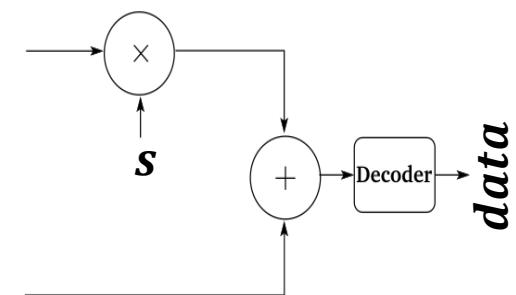
1. Encrypt

public keys (p_0, p_1)



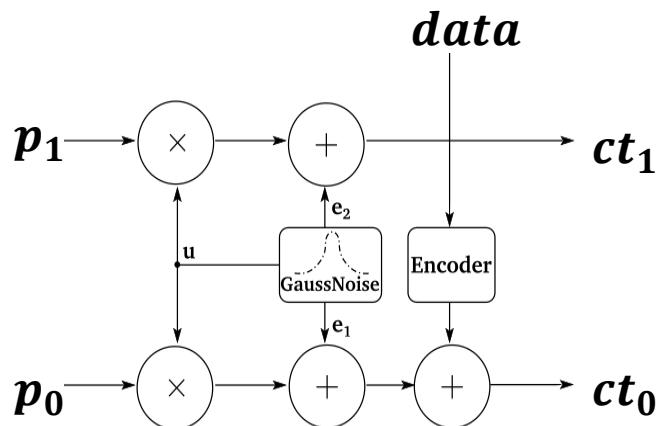
2. Decrypt

private key s

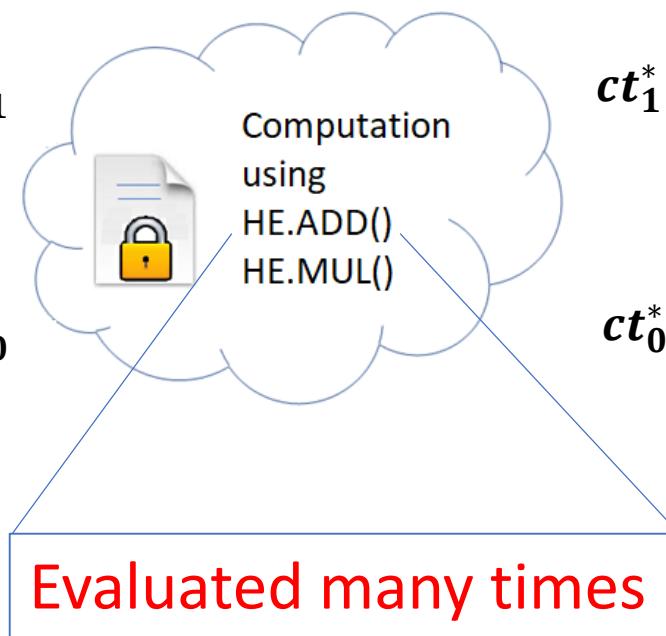


Lattice-based Homomorphic encryption

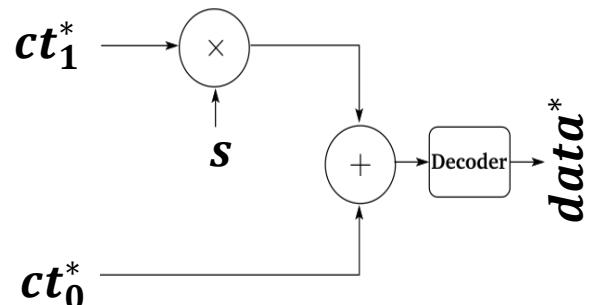
1. Encrypt locally
public keys (p_0, p_1)



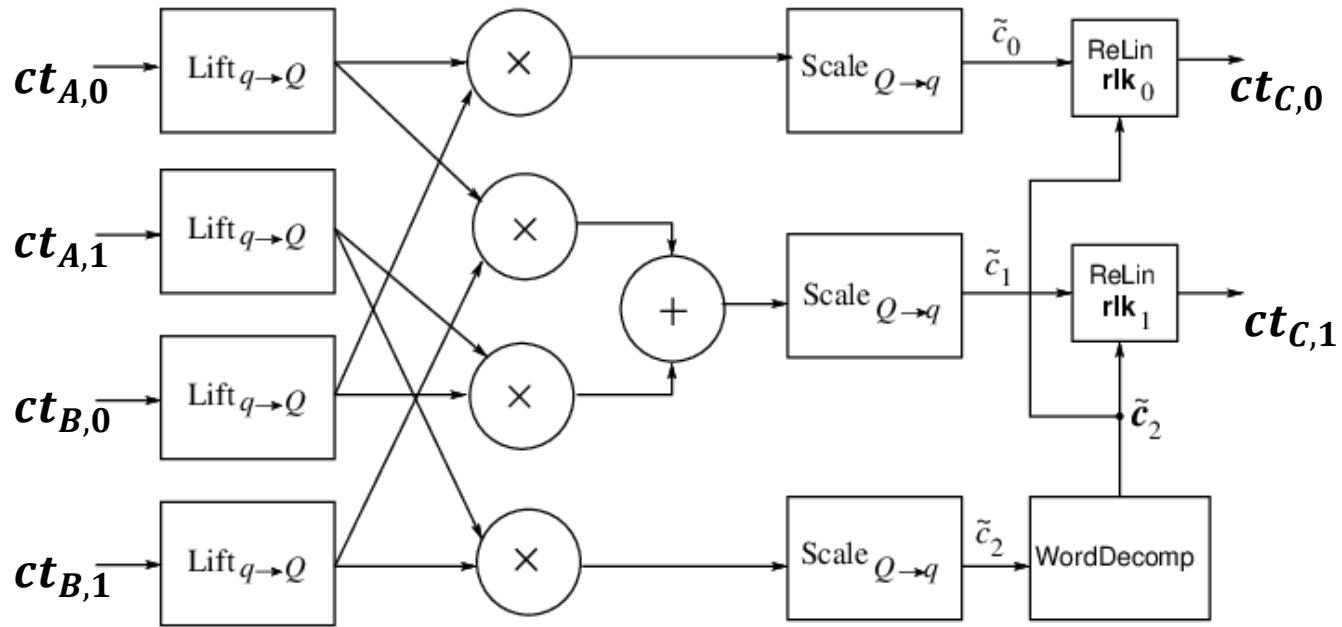
2. Process on cloud



3. Decrypt locally
private key s



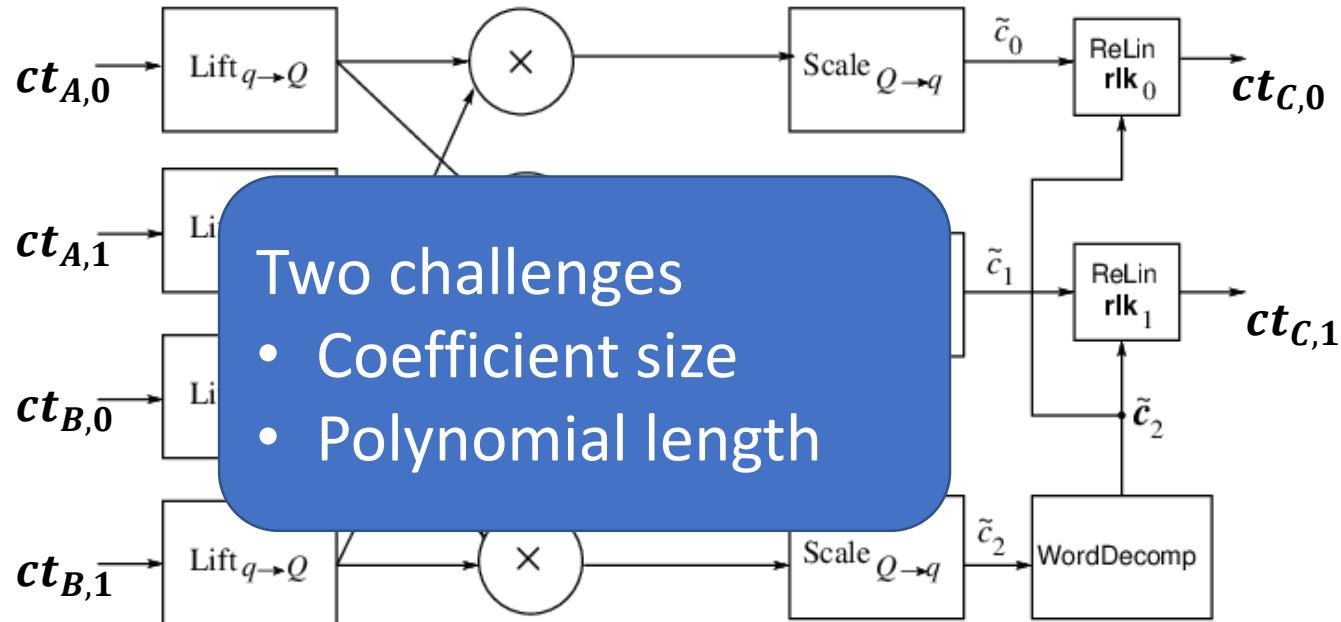
Homomorphic Multiplication



Uses multiple computation blocks:

- Lift
- Polynomial multiplication
- Scale

Homomorphic Multiplication. How complex?



Low complexity applications:

- Polynomials have 4,000 coeffs.
- Coeffs are ~180 bit wide

Medium complexity applications:

- Polynomials have 32,000 coeffs.
- Coeffs are ~1200 bit wide

Lattice-based key exchange schemes

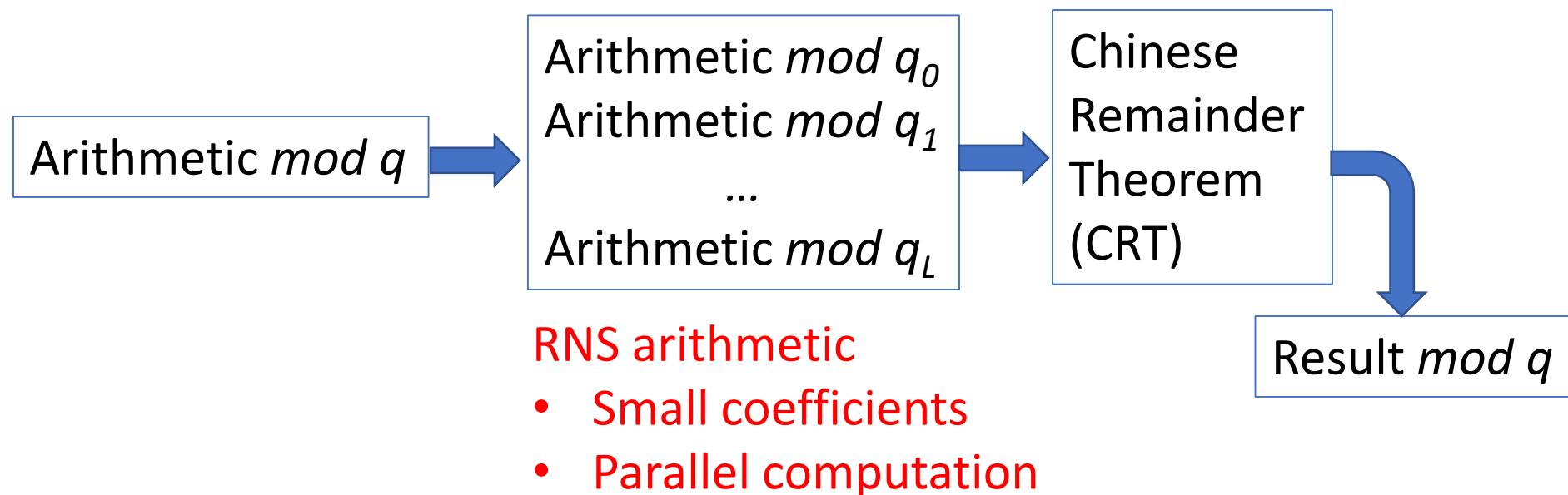
- Polynomials with 256 or 512 coeffs
- Coeff size ~10 bits

Hardware accelerators for Homomorphic Computation

→ Arithmetic of large coefficients

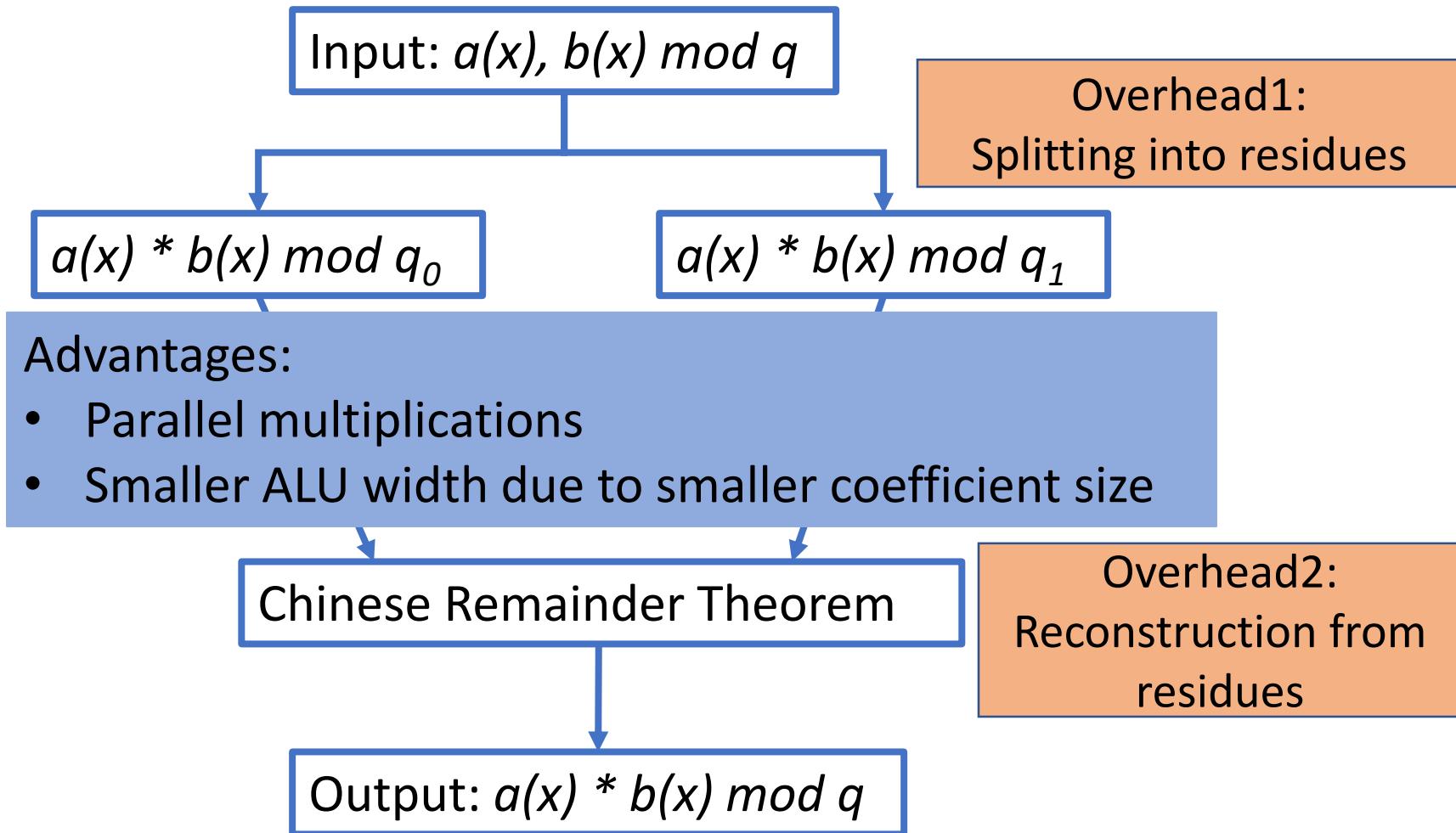
Application of Residue Number System

- We need to compute arithmetic modulo q
- Let $q = \prod q_i$ where q_i are coprime
- Then we can work with Residue Number System (RNS)

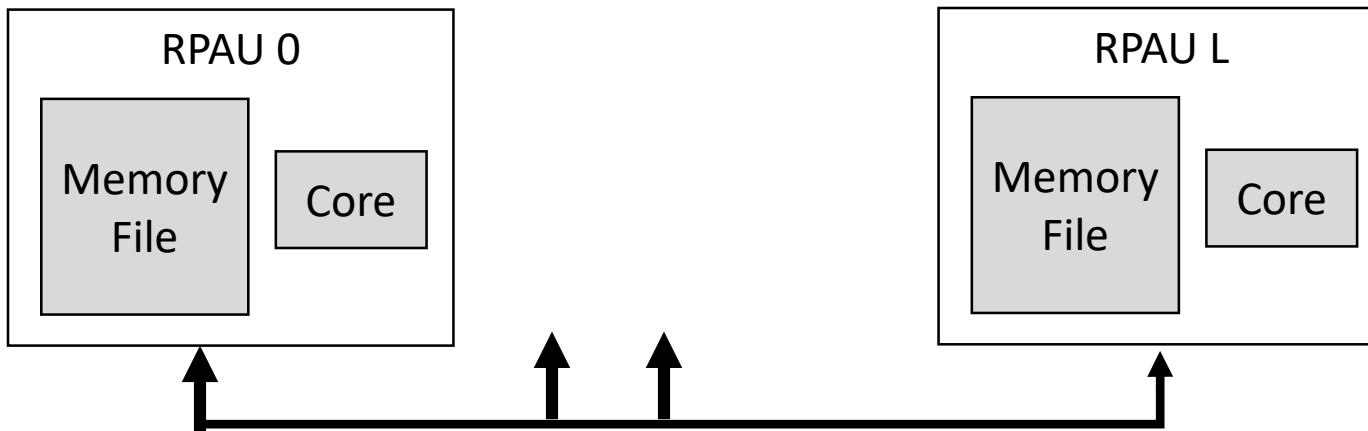


Example: polynomial multiplication

Let $q = q_0 \cdot q_1$ where q_0 and q_1 are of equal bit-length



On Hardware Parallel processing using multiple Residue Polynomial Arithmetic Unit (RPAU)

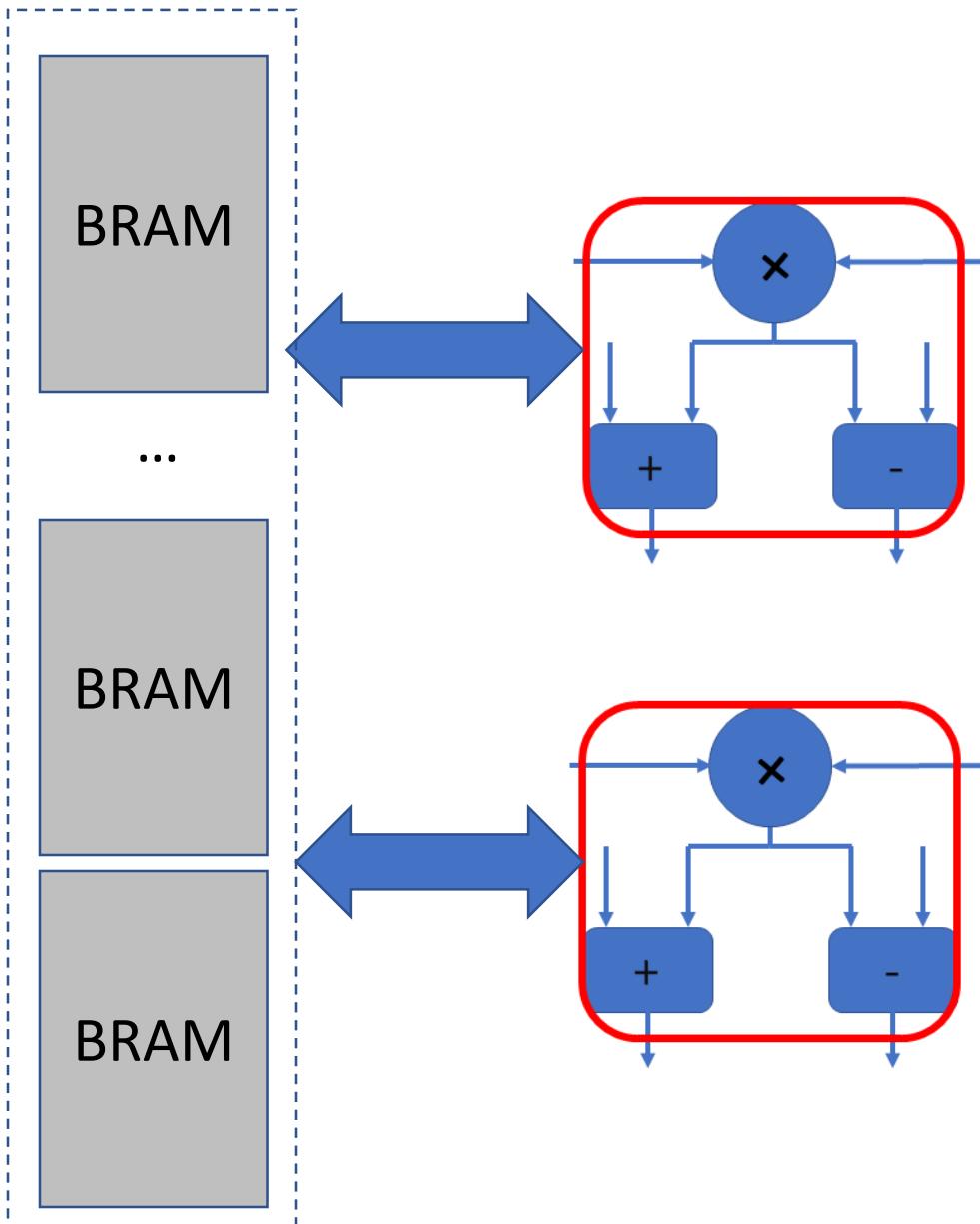


Number of RPAUs is a design parameter

Hardware accelerators for Homomorphic Computation

- Arithmetic of large coefficients
- Arithmetic of large polynomials

Polynomial multiplication: multiple butterfly cores



Single core NTT too slow!

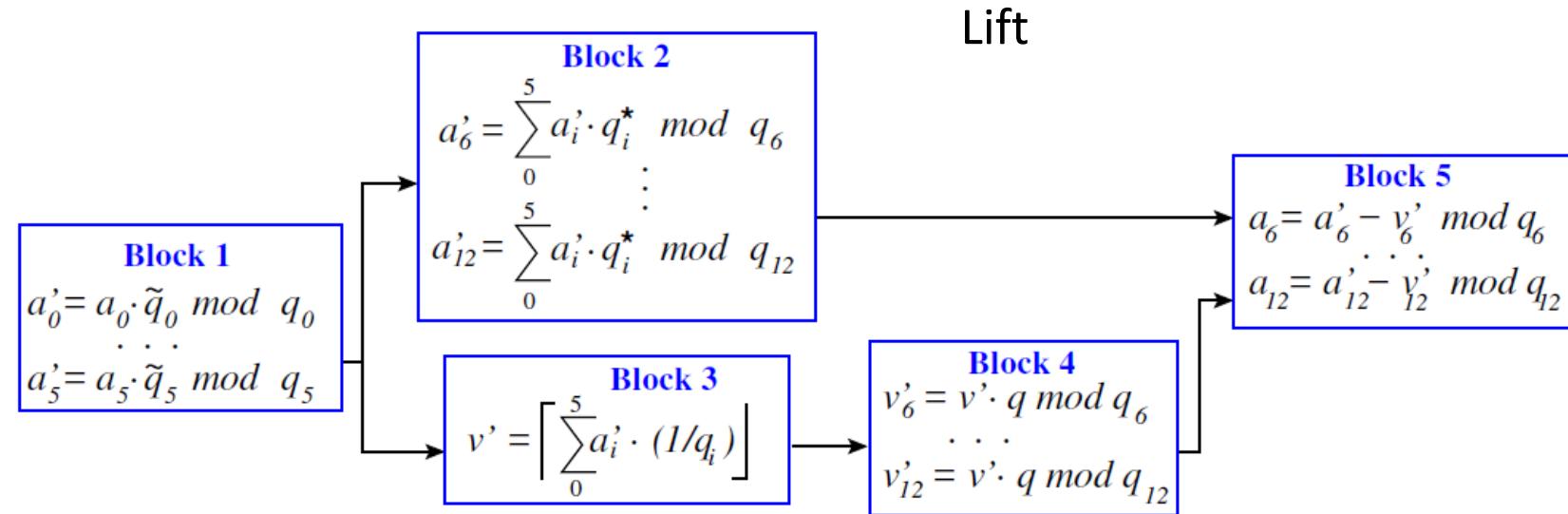
Design Challenges:

- Long routing
- Memory access conflicts

Memory access parallelism



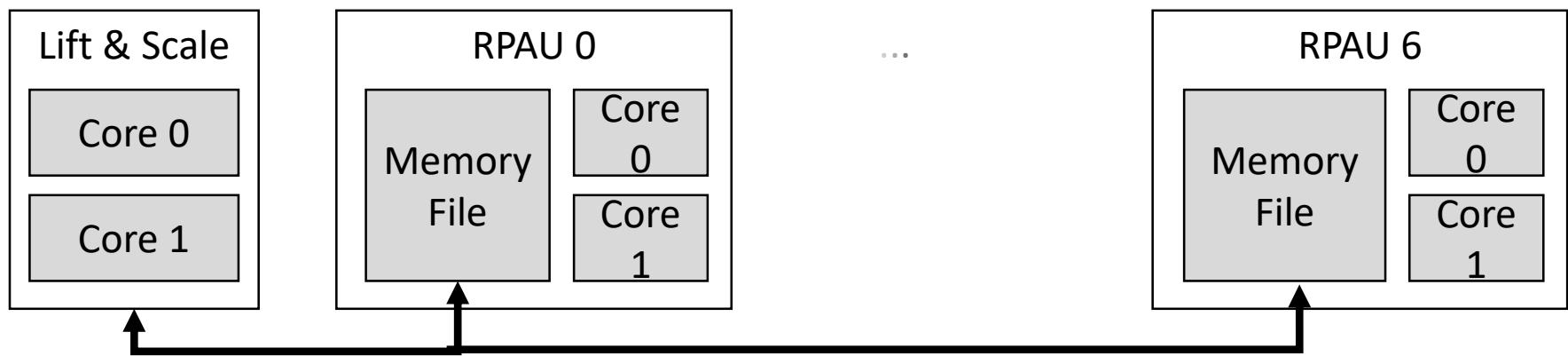
Block Level Pipelining



- Separate building blocks for block-level pipeline
- Realize a resource shared architecture
 - Reduces the area requirement
 - Increase the computation time

Execution Units

- Two parallel cores for Lift and Scale
- Seven Residue Polynomial Arithmetic Unit (RPAU)

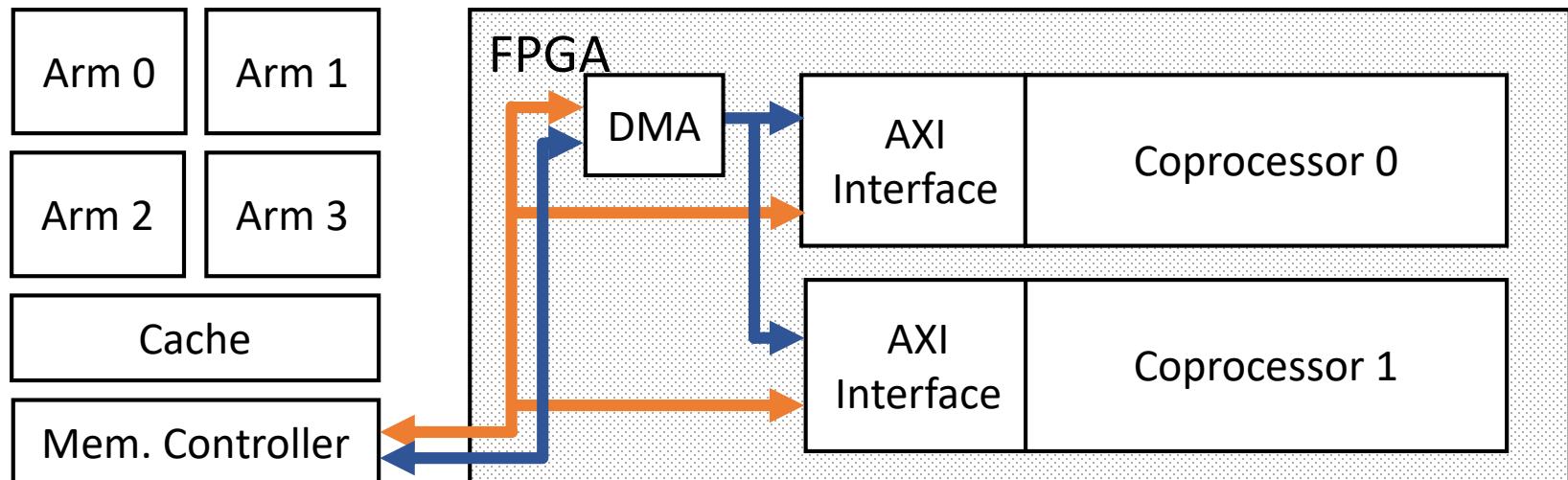


Parameter:

- Ciphertext polynomial degree 4096
- Ciphertext coefficient size 180

Arm + FPGA Implementation

Zynq UltraScale+ MPSoC ZCU102



Source code public on Github

Performance of High-Level Operations

Operation	Speed (cycles)	Speed (msec)
Add in HW	31,339	0.026
Multiply in HW	5,349,567	4.458
Send two ciphertext to HW	434,013	0.362
Receive result ciphertext from HW	216,697	0.180

Measurements are in cycles of CPU clocked at 1200 MHz

Coprocessor is clocked at 200 MHz

Publication: HPCA 2019

400 homomorphic multiplications per sec (2 cores)

Faster than Tesla K80 GPU

Resource Utilization

	LUTs	REGs	BRAMs	DSPs
	# of used instances % utilization			
Two Coprocessors & Interface	133,692 49	60,312 11	815 89	416 16
A Single Coprocessor & Interface	63,522 23	25,622 5	388 43	208 8

Conclusions so far

- Ring-LWE is efficient in hardware and software
- But, there are security concerns due to special structure

$$\begin{pmatrix} a_0 & -a_3 & -a_2 & -a_1 \\ a_1 & a_0 & -a_3 & -a_2 \\ a_2 & a_1 & a_0 & -a_3 \\ a_3 & a_2 & a_1 & a_0 \end{pmatrix} * \begin{pmatrix} s_0 \\ s_1 \\ s_2 \\ s_3 \end{pmatrix} + \begin{pmatrix} e_0 \\ e_1 \\ e_2 \\ e_3 \end{pmatrix} \approx \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \end{pmatrix} \pmod{q}$$

Special structure in matrix

Interpolating LWE and ring-LWE: Module LWE

$$\begin{array}{c}
 \left(\begin{array}{cccc} a_0 & -a_3 & -a_2 & -a_1 \\ a_1 & a_0 & -a_3 & -a_2 \\ a_2 & a_1 & a_0 & -a_3 \\ a_3 & a_2 & a_1 & a_0 \end{array} \right) \quad \left(\begin{array}{cccc} a_8 & -a_{11} & -a_{10} & -a_9 \\ a_9 & a_8 & -a_{11} & -a_{10} \\ a_{10} & a_9 & a_8 & -a_{11} \\ a_{11} & a_{10} & a_7 & a_8 \end{array} \right) * \left(\begin{array}{c} s_0 \\ s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \\ s_6 \\ s_7 \end{array} \right) + \left(\begin{array}{c} e_0 \\ e_1 \\ e_2 \\ e_3 \\ e_4 \\ e_5 \\ e_6 \\ e_7 \end{array} \right) \approx \left(\begin{array}{c} b_0 \\ b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \\ b_6 \\ b_7 \end{array} \right) \\
 \left(\begin{array}{cccc} a_4 & -a_7 & -a_6 & -a_5 \\ a_5 & a_4 & -a_7 & -a_6 \\ a_6 & a_5 & a_4 & -a_7 \\ a_7 & a_6 & a_5 & a_4 \end{array} \right) \quad \left(\begin{array}{cccc} a_{12} & -a_{15} & -a_{14} & -a_{13} \\ a_{13} & a_{12} & -a_{15} & -a_{14} \\ a_{14} & a_{13} & a_{12} & -a_{15} \\ a_{15} & a_{14} & a_{13} & a_{12} \end{array} \right)
 \end{array}$$



$$\left(\begin{array}{cc} a_{0,0}(x) & a_{0,1}(x) \\ a_{1,0}(x) & a_{1,1}(x) \end{array} \right) * \left(\begin{array}{c} s_0(x) \\ s_1(x) \end{array} \right) + \left(\begin{array}{c} e_0(x) \\ e_1(x) \end{array} \right) \approx \left(\begin{array}{c} b_0(x) \\ b_1(x) \end{array} \right) \pmod{q} \pmod{x^4 + 1}$$



The screenshot shows the NIST Computer Security Division website. At the top, the NIST logo and "National Institute of Standards and Technology Information Technology Laboratory" are displayed. To the right is a search bar with "SEARCH:" and a "Search" button, along with links for "CONTACT" and "SITE MAP". Below the header, the "Computer Security Division" and "Computer Security Resource Center" are prominently featured. A navigation menu below the header includes "CSRC Home", "About", "Projects / Research", "Publications", and "News & Events". The main content area is titled "Post-Quantum Cryptography Project" and includes sections for "Documents", "Workshops / Timeline", "Federal Register Notices", "Email Listserve", "PQC Project Contact", and "Archive Information". Another section titled "Post-Quantum Cryptography Standardization" lists "Call for Proposals Announcement", "Call for Proposals", "Submission Requirements", and "Minimum Acceptability Requirements". Above the main content, a breadcrumb trail reads "CSRC HOME > GROUPS > CT > POST-QUANTUM CRYPTOGRAPHY PROJECT".

POST-QUANTUM CRYPTO PROJECT

NEWS -- December 15, 2016: The National Institute of Standards and Technology (NIST) is now accepting submissions for quantum-resistant public-key cryptographic algorithms. The deadline for submission is **November 30, 2017**. Please see the Post-Quantum Cryptography Standardization menu at left for the complete submission requirements and evaluation criteria.

In recent years, there has been a substantial amount of research on quantum computers – machines that exploit quantum mechanical phenomena to solve mathematical problems that are difficult or intractable for conventional computers. If large-scale quantum computers are ever built, they will be able to break many of the public-key cryptosystems currently in use. This would seriously compromise the confidentiality and integrity of digital communications on the Internet and elsewhere. The goal of *post-quantum cryptography* (also called quantum-resistant cryptography) is to develop cryptographic systems that are secure against both quantum and classical computers, and can

Saber: Module-LWR based key exchange, CPA-secure encryption and CCA-secure KEM

**a lattice-based candidate for NIST standardization
moved to second round!**

Jointly designed by EE and Math team!

SABER: flexibility and efficiency

- Saber uses module-LWR problem
- Polynomials are always of 256 coefficients [**Efficient pol. arithmetic**]
- **Flexibility:** matrix dimensions is parameterizable
 - 2-by-2 for 115-bit post-quantum security



Light SABER

- 3-by-3 for 180-bit post-quantum security



SABER

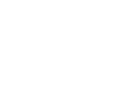
- 4-by-4 for 245-bit post-quantum security



Fire SABER

SABER: Parameter set

$$\left[\frac{p}{q} \begin{pmatrix} a_{0,0}(x) & \dots & a_{0,k-1}(x) \\ \dots & \dots & \dots \\ a_{k-1,0}(x) & \dots & a_{k-1,k-1}(x) \end{pmatrix} * \begin{pmatrix} s_0(x) \\ \dots \\ s_{k-1}(x) \end{pmatrix} \right] \approx \begin{pmatrix} b_0(x) \\ b_{k-1}(x) \end{pmatrix} \pmod{x^{256} + 1}$$

- Polynomials of fixed size 256 coefficients 
- Flexible dimension $k = 2, 3$ or 4 
- How to choose p and q ? 

Learning with rounding (LWR)

A problem with rounding:

$$\left\lfloor \frac{p}{q} \quad \text{Uniform in } [0, q-1] \right\rfloor \quad \text{where } p < q$$

Prime q introduces rounding bias

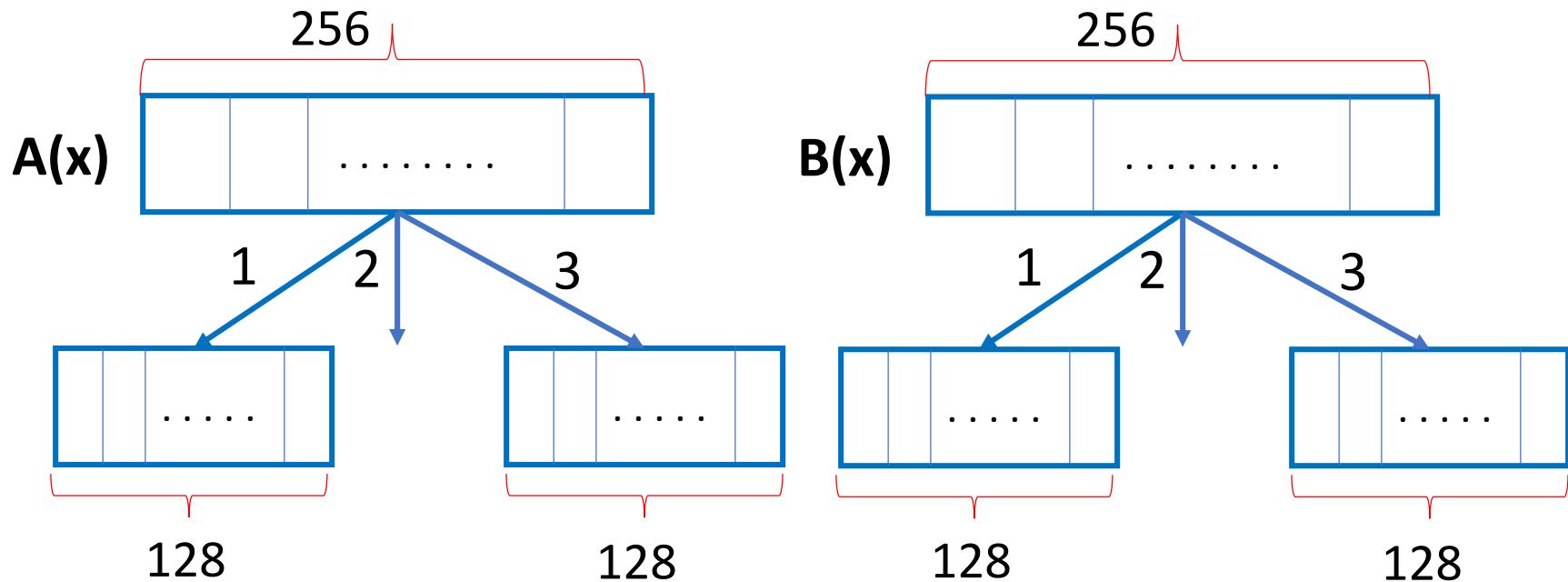


- Cannot use prime q 
- Hence, no NTT-based fast polynomial multiplication
- + No modular reduction + Easy rounding

→ We need to use generic polynomial multiplication algorithm

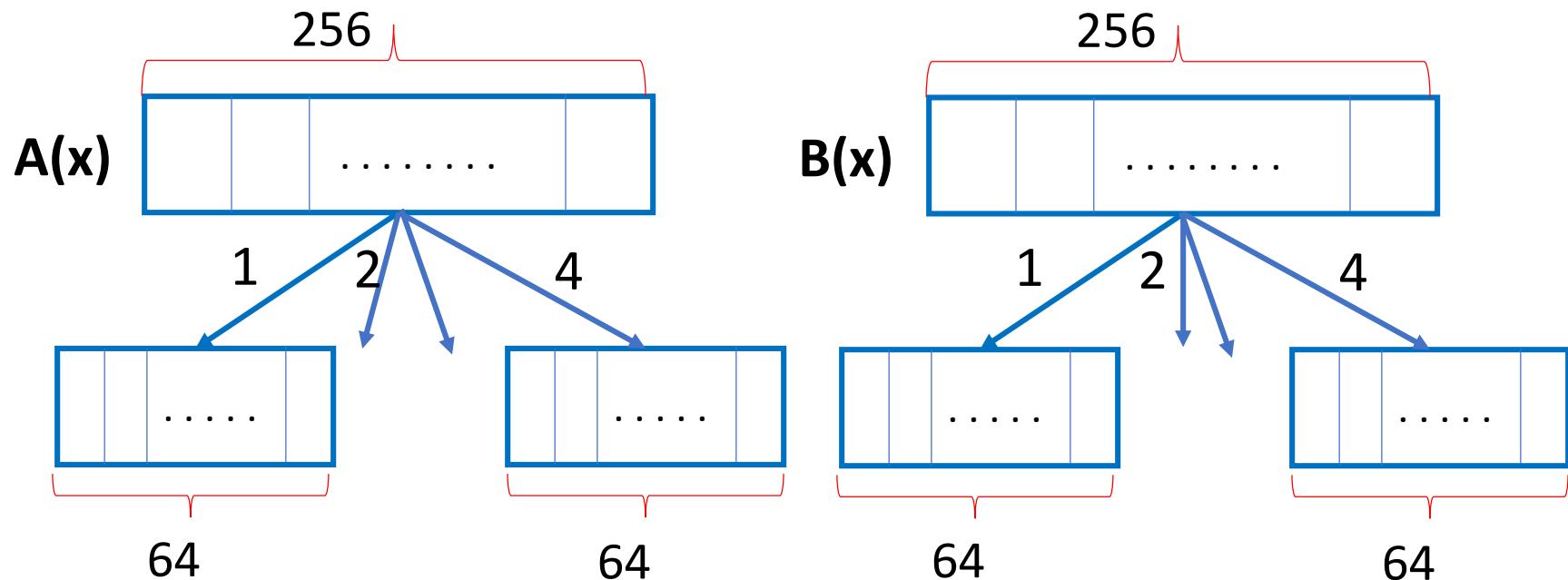
Next best polynomial multiplication algorithms

- Karatsuba multiplication $O(n^{\log_2 3})$



Next best polynomial multiplication algorithms

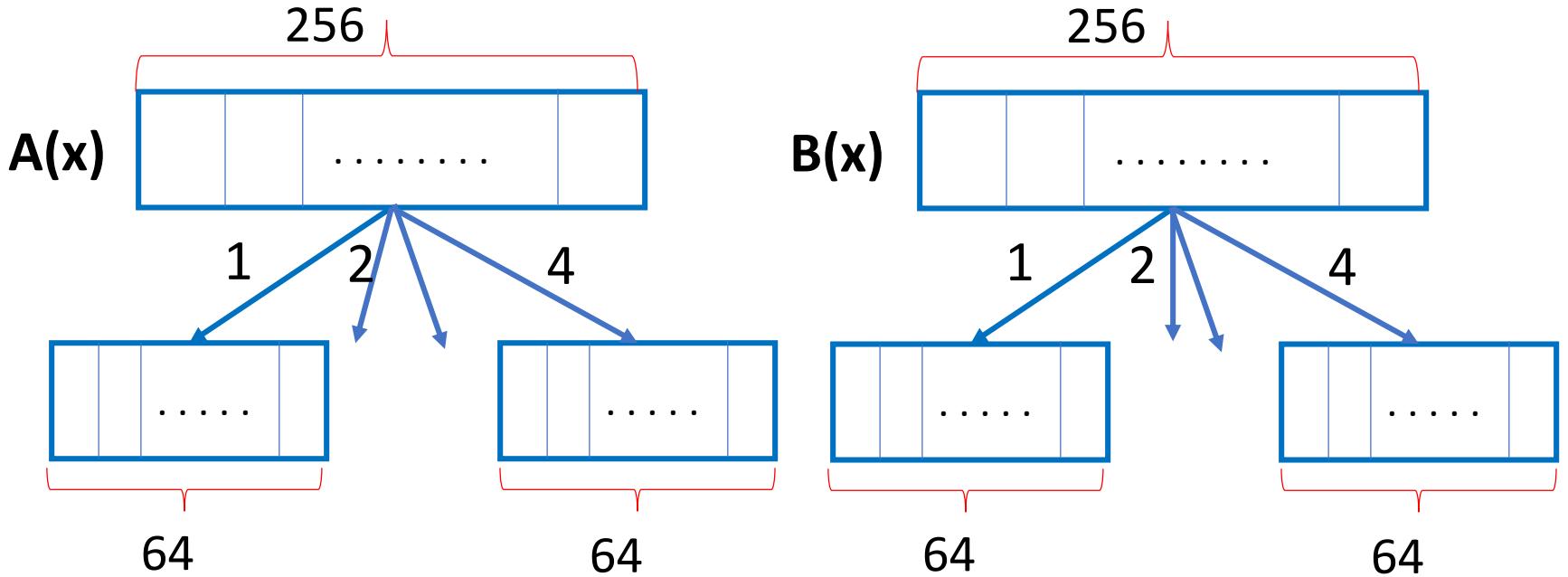
- Toom-Cook multiplication



Toom-Cook 4 Way needs 7 multiplications

Karatsuba would need 9 multiplications

Toom-Cook 4 Way: step-by-step: splitting



Splitting operand into 4 polynomials

Take $y = x^{64}$

$$A(y) = A_3 y^3 + A_2 y^2 + A_1 y + A_0$$

$$B(y) = B_3 y^3 + B_2 y^2 + B_1 y + B_0$$

Toom-Cook 4 Way: step-by-step: evaluation

$$w_1 = A(\infty) * B(\infty) = A_3 * B_3$$

$$w_2 = A(2) * B(2) = (A_0 + 2 \cdot A_1 + 4 \cdot A_2 + 8 \cdot A_3) * (B_0 + 2 \cdot B_1 + 4 \cdot B_2 + 8 \cdot B_3)$$

$$w_3 = A(1) * B(1) = (A_0 + A_1 + A_2 + A_3) * (B_0 + B_1 + B_2 + B_3)$$

$$w_4 = A(-1) * B(-1) = (A_0 - A_1 + A_2 - A_3) * (B_0 - B_1 + B_2 - B_3)$$

$$w_5 = A\left(\frac{1}{2}\right) * B\left(\frac{1}{2}\right) = (8 \cdot A_0 + 4 \cdot A_1 + 2 \cdot A_2 + A_3) * (8 \cdot B_0 + 4 \cdot B_1 + 2 \cdot B_2 + B_3)$$

$$w_6 = A\left(\frac{-1}{2}\right) * B\left(\frac{-1}{2}\right) = (8 \cdot A_0 - 4 \cdot A_1 + 2 \cdot A_2 - A_3) * (8 \cdot B_0 - 4 \cdot B_1 + 2 \cdot B_2 - B_3)$$

$$w_7 = A(0) * B(0) = A_0 * B_0$$

Linear operations

+

Seven multiplications are computed

Toom-Cook 4 Way: step-by-step: interpolation

```
// Interpolation
```

$$w_2 = w_2 + w_5$$

$$w_6 = w_6 - w_5$$

$$w_4 = (w_4 - w_3)/2$$

$$w_5 = w_5 - w_1 - 64 \cdot w_7$$

$$w_3 = w_3 + w_4$$

$$w_5 = 2 \cdot w_5 + w_6$$

$$w_2 = w_2 - 65 \cdot w_3$$

$$w_3 = w_3 - w_7 - w_1$$

$$w_2 = w_2 + 45 \cdot w_3$$

$$w_5 = (w_5 - 8 \cdot w_3)/24$$

$$w_6 = w_6 + w_2$$

$$w_2 = (w_2 + 16 \cdot w_4)/18$$

$$w_3 = w_3 - w_5$$

$$w_4 = -(w_4 + w_2)$$

$$w_6 = (30 \cdot w_2 - w_6)/60$$

$$w_2 = w_2 - w_6$$

```
return  $w_1 \cdot y^6 + w_2 \cdot y^5 + w_3 \cdot y^4 + w_4 \cdot y^3 + w_5 \cdot y^2 + w_6 \cdot y + w_7;$ 
```

Linear operations

This number has a role
to play

Linear operations

Advanced Vector Extensions (AVX)

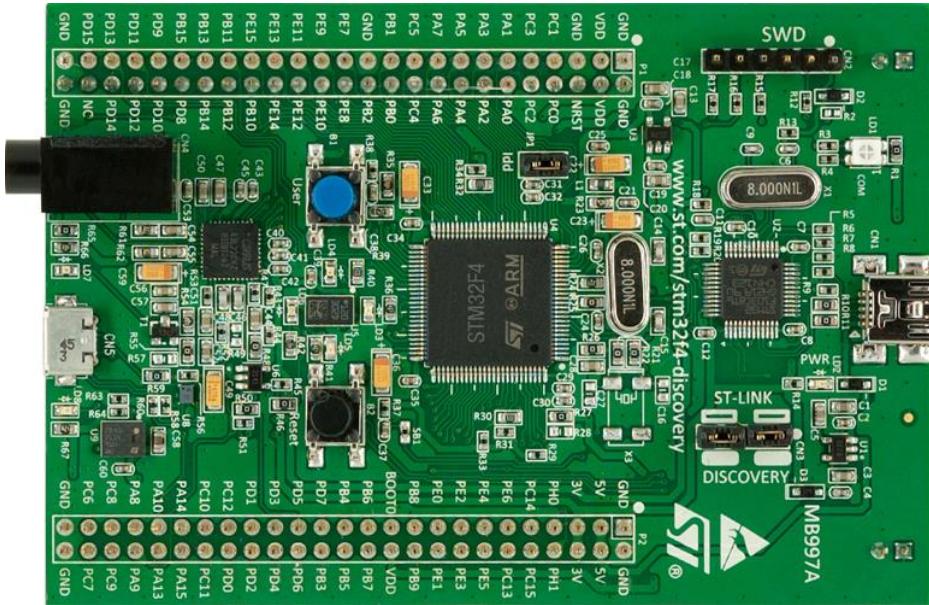


Intrinsics Guide

```
__m256i _mm256_abs_epi16 (__m256i a)
__m256i _mm256_add_epi16 (__m256i a, __m256i b)
__m256i _mm256_adds_epi16 (__m256i a, __m256i b)
__m256i _mm256_blend_epi16 (__m256i a, __m256i b, const int imm8)
__m128i _mm_broadcastw_epi16 (__m128i a)
__m256i _mm256_broadcastw_epi16 (__m128i a)
__m256i _mm256_cmpeq_epi16 (__m256i a, __m256i b)
__m256i _mm256_cmpgt_epi16 (__m256i a, __m256i b)
__m256i _mm256_cvtepi16_epi32 (__m128i a)
__m256i _mm256_cvtepi16_epi64 (__m128i a)
__m256i _mm256_cvtepi8_epi16 (__m128i a)
__m256i _mm256_cvtepu8_epi16 (__m128i a)
int _mm256_extract_epi16 (__m256i a, const int index)
__m256i _mm256_hadd_epi16 (__m256i a, __m256i b)
__m256i _mm256_hadds_epi16 (__m256i a, __m256i b)
__m256i _mm256_hsub_epi16 (__m256i a, __m256i b)
__m256i _mm256_hsubs_epi16 (__m256i a, __m256i b)
```

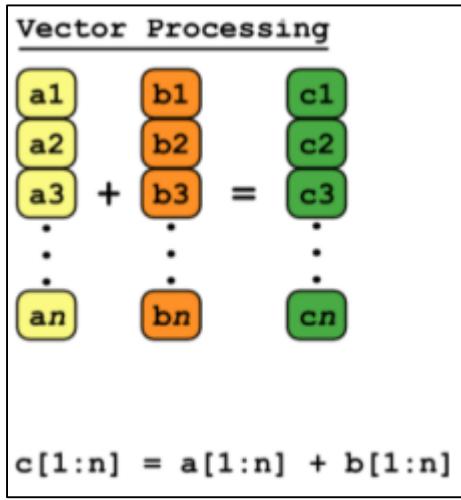
Vectorized instructions for 16-bit operands

DSP instructions



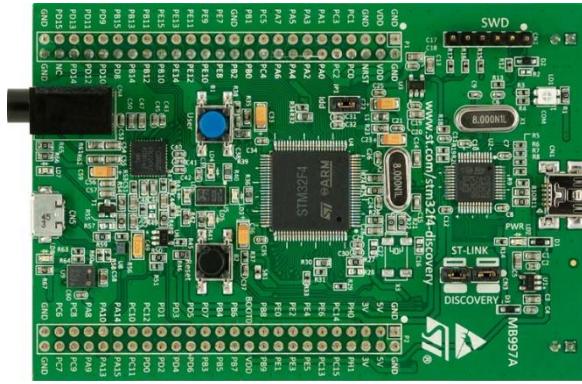
ARM Cortex-M4

- Popular 32-bit microcontroller
- Has DSP instructions for half-word operations



AVX

+



Microcontroller with DSP

Keep coefficients smaller/equal to 16 bits to use

- `_epi16()` AVX intrinsics in high-end platforms
- DSP instructions in low-end microcontrollers

Options for q: $2^{16}, 2^{15}, 2^{14}, 2^{13} \dots$ etc



Toom-Cook 4 Way: step-by-step: interpolation

```
// Interpolation
```

$$w_2 = w_2 + w_5$$

$$w_6 = w_6 - w_5$$

$$w_4 = (w_4 - w_3)/2$$

$$w_5 = w_5 - w_1 - 64 \cdot w_7$$

$$w_3 = w_3 + w_4$$

$$w_5 = 2 \cdot w_5 + w_6$$

$$w_2 = w_2 - 65 \cdot w_3$$

$$w_3 = w_3 - w_7 - w_1$$

$$w_2 = w_2 + 45 \cdot w_3$$

$$w_5 = (w_5 - 8 \cdot w_3)/24$$



This number has a role
to play

$$w_6 = w_6 + w_2$$

$$w_2 = (w_2 + 16 \cdot w_4)/18$$

$$w_3 = w_3 - w_5$$

$$w_4 = -(w_4 + w_2)$$

$$w_6 = (30 \cdot w_2 - w_6)/60$$

$$w_2 = w_2 - w_6$$

return $w_1 \cdot y^6 + w_2 \cdot y^5 + w_3 \cdot y^4 + w_4 \cdot y^3 + w_5 \cdot y^2 + w_6 \cdot y + w_7;$

Division by 24 in Toom-Cook Interpolation

$$w_5 = (w_5 - 8 \cdot w_3)/24$$

- $24 = 8 \cdot 3$
- We are working in R_q where $q = 2^i$
- 3 has inverse in mod q
E.g. $3^{-1} \text{ mod } 2^{15} \rightarrow 10923$
- So, division by 3 is same as multiplying by $3^{-1} \text{ mod } q$

Division by 24 in Toom-Cook Interpolation

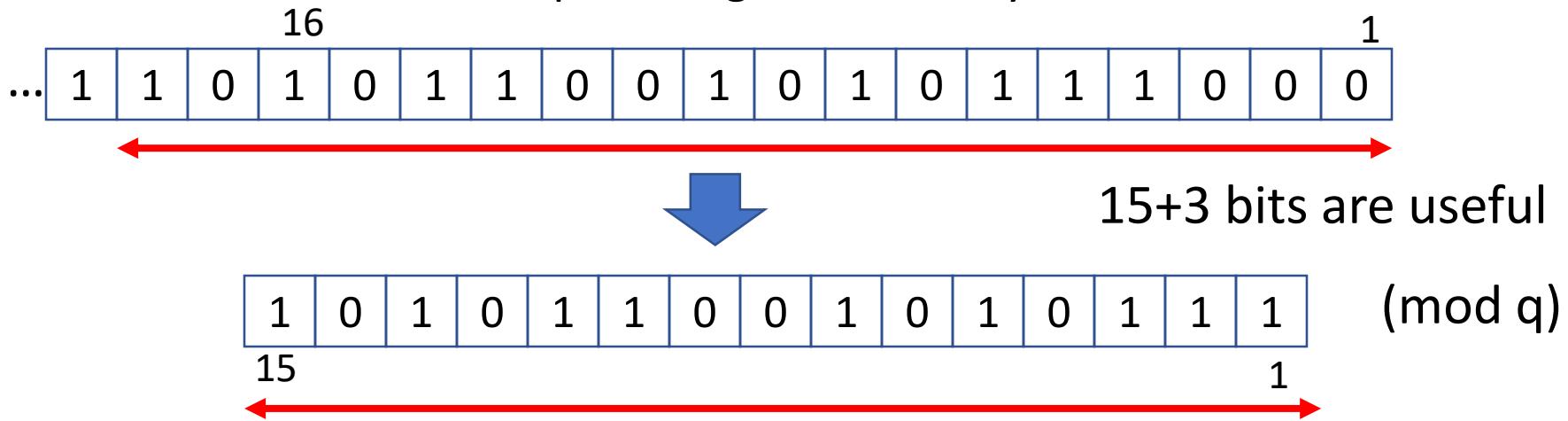
$$w_5 = (w_5 - 8 \cdot w_3)/24$$

- $24 = 8 \cdot 3$
- We are working in R_q where $q = 2^i$
- But, 8 does not have inverse in mod q

Only option: do actual division

Working with $q = 2^{15}$

Example: integer division by $8=2^3$

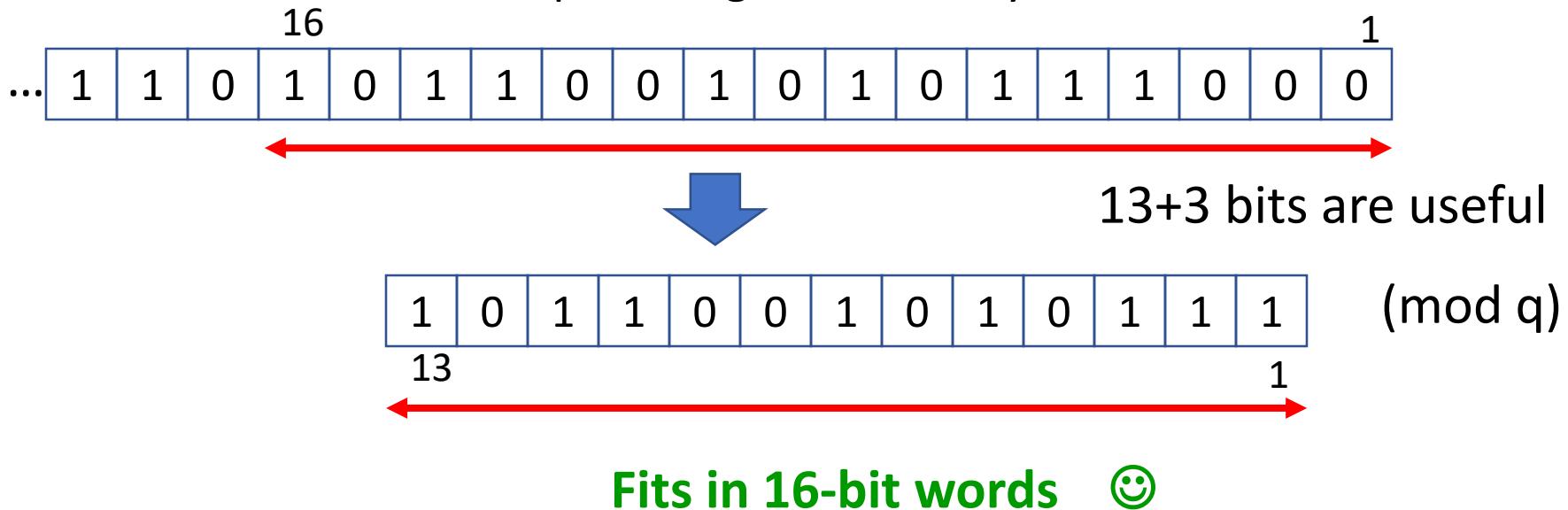


In 16-bit Computer:

- Difficult to implement: requires careful arithmetic of two words
- Slower arithmetic

Working with $q = 2^{13}$

Example: integer division by $8=2^3$



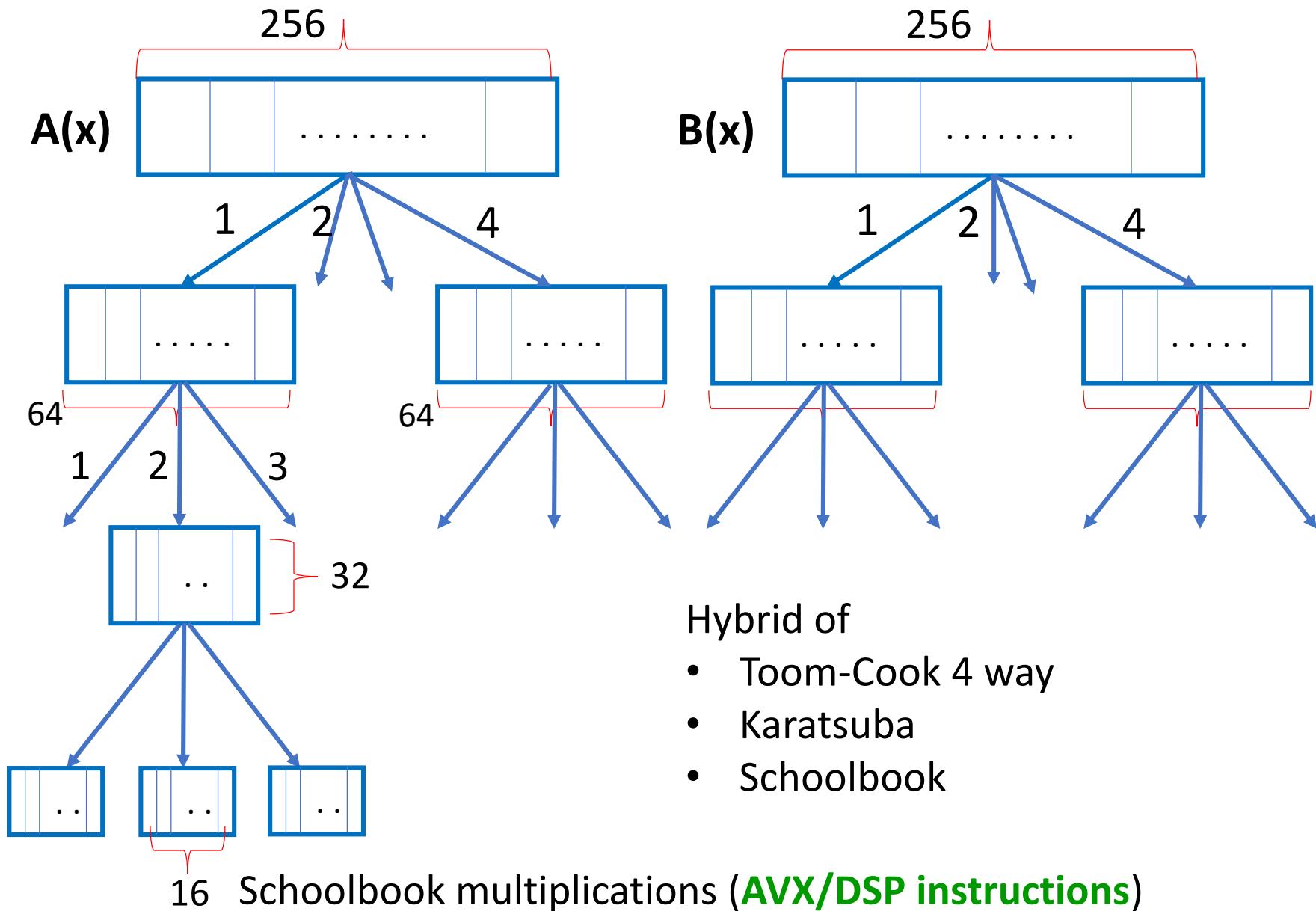
In 16-bit Computer:

- **Easy** to implement
- **Less complicated** arithmetic

Saber Parameters

- Polynomial length $n = 256$
- $q = 2^{13}$
- $p = 2^{10}$

Polynomial multiplication in Saber



School-book multiplication using AVX

Consider one polynomial multiplication

$$\begin{array}{r} \dots b_3 x^3 + b_2 x^2 + b_1 x + b_0 \\ \hline \dots a_3 x^3 + a_2 x^2 + a_1 x + a_0 \\ \hline a_0 b_2 & a_0 b_1 & a_0 b_0 \\ \cdot \quad \cdot \quad \cdot & a_1 b_1 & a_1 b_0 \\ & a_2 b_0 \end{array}$$

School-book multiplication using AVX

Consider 16 polynomial multiplications

$$a^{15}(x) * b^{15}(x)$$

$$\dots b_3^{15} x^3 + b_2^{15} x^2 + b_1^{15} x + b_0^{15}$$

$$\dots a_3^{15} x^3 + a_2^{15} x^2 + a_1^{15} x + a_0^{15}$$

$$a_0^{15} b_2^{15} \quad a_0^{15} b_1^{15} \quad a_0^{15} b_0^{15}$$

$$\cdot \quad \cdot \quad \cdot \quad a_1^{15} b_1^{15} \quad a_1^{15} b_0^{15}$$

$$a_2^{15} b_0^{15}$$

$$a^0(x) * b^0(x)$$

$$\dots b_3^0 x^3 + b_2^0 x^2 + b_1^0 x + b_0^0$$

$$\dots a_3^0 x^3 + a_2^0 x^2 + a_1^0 x + a_0^0$$

$$a_0^0 b_2^0 \quad a_0^0 b_1^0 \quad a_0^0 b_0^0$$

$$\cdot \quad \cdot \quad \cdot \quad a_1^0 b_1^0 \quad a_1^0 b_0^0$$

$$a_2^0 b_0^0$$

and 16x vectorized processor

School-book multiplication using AVX

Consider 16 polynomial multiplications

$$\dots b_3^{15} x^3 + b_2^{15} x^2 + b_1^{15} x + b_0^{15}$$

$$\dots a_3^{15} x^3 + a_2^{15} x^2 + a_1^{15} x + a_0^{15}$$

$$a_0^{15} b_2^{15} \quad a_0^{15} b_1^{15} \quad a_0^{15} b_0^{15}$$

$$\cdot \quad \cdot \quad \cdot \quad a_1^{15} b_1^{15} \quad a_1^{15} b_0^{15}$$

$$a_2^{15} b_0^{15}$$

$$\dots b_3^0 x^3 + b_2^0 x^2 + b_1^0 x + b_0^0$$

$$\dots a_3^0 x^3 + a_2^0 x^2 + a_1^0 x + a_0^0$$

$$a_0^0 b_2^0 \quad a_0^0 b_1^0 \quad a_0^0 b_0^0$$

$$\cdot \quad \cdot \quad \cdot \quad a_1^0 b_1^0 \quad a_1^0 b_0^0$$

$$a_2^0 b_0^0$$

and 16x vectorized processor

$$\mathbf{A0} = \boxed{a_0^{15}} \dots \boxed{a_0^2} \boxed{a_0^1} \boxed{a_0^0}$$

$$\mathbf{B0} = \boxed{b_0^{15}} \dots \boxed{b_0^2} \boxed{b_0^1} \boxed{b_0^0}$$

School-book multiplication using AVX

Consider 16 polynomial multiplications

$$\begin{array}{l} \dots b_3^{15}x^3 + b_2^{15}x^2 + b_1^{15}x + b_0^{15} \\ \dots a_3^{15}x^3 + a_2^{15}x^2 + a_1^{15}x + a_0^{15} \\ \hline a_0^{15}b_2^{15} \quad a_0^{15}b_1^{15} \quad a_0^{15}b_0^{15} \\ \cdot \quad \cdot \quad \cdot \quad a_1^{15}b_1^{15} \quad a_1^{15}b_0^{15} \\ \quad \quad \quad a_2^{15}b_0^{15} \end{array}$$

$$\begin{array}{l} \dots b_3^0x^3 + b_2^0x^2 + b_1^0x + b_0^0 \\ \dots a_3^0x^3 + a_2^0x^2 + a_1^0x + a_0^0 \\ \hline a_0^0b_2^0 \quad a_0^0b_1^0 \quad a_0^0b_0^0 \\ \cdot \quad \cdot \quad \cdot \quad a_1^0b_1^0 \quad a_1^0b_0^0 \\ \quad \quad \quad a_2^0b_0^0 \end{array}$$

and 16x vectorized processor

$$A0 = \boxed{a_0^{15}} \dots \boxed{a_0^2} \boxed{a_0^1} \boxed{a_0^0}$$

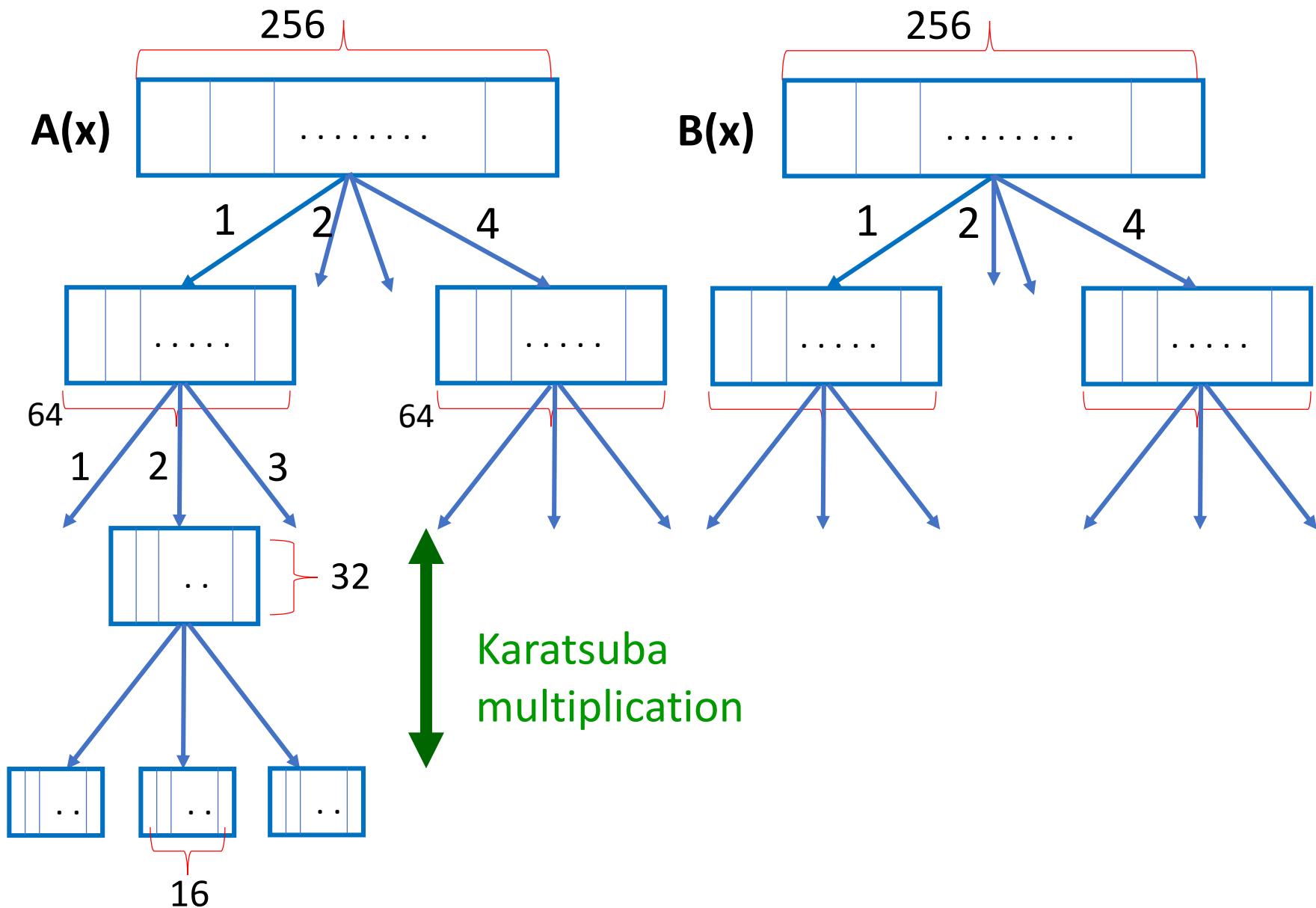
AVX_MUL(A0, B0)

$$A0B0 = \boxed{a_0^{15}b_0^{15}} \dots \boxed{a_0^1b_0^1} \boxed{a_0^0b_0^0}$$

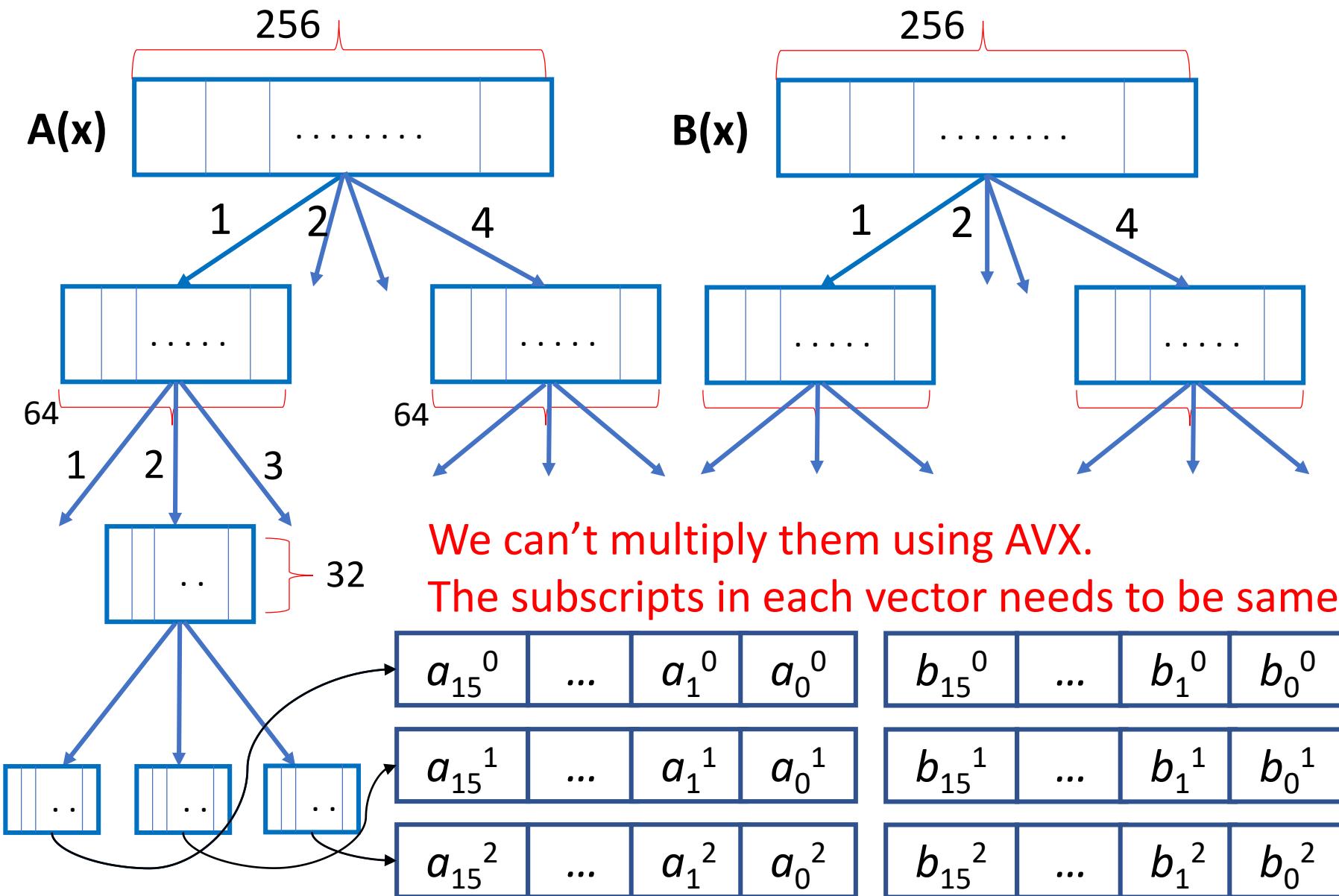
$$B0 = \boxed{b_0^{15}} \dots \boxed{b_0^2} \boxed{b_0^1} \boxed{b_0^0}$$

All in parallel
(assumes all coeffs. are available)

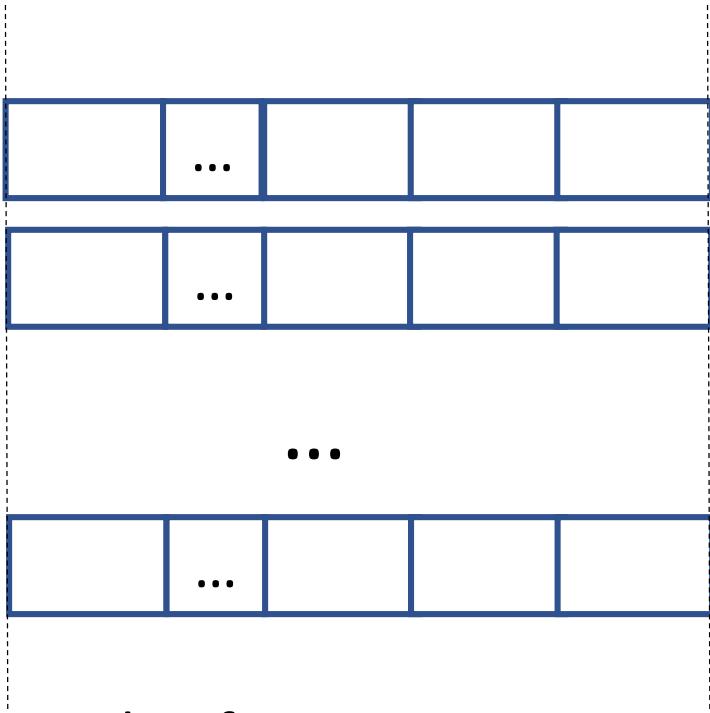
Polynomial multiplication



Polynomial multiplication



Buckets for 16 coeff. polynomials

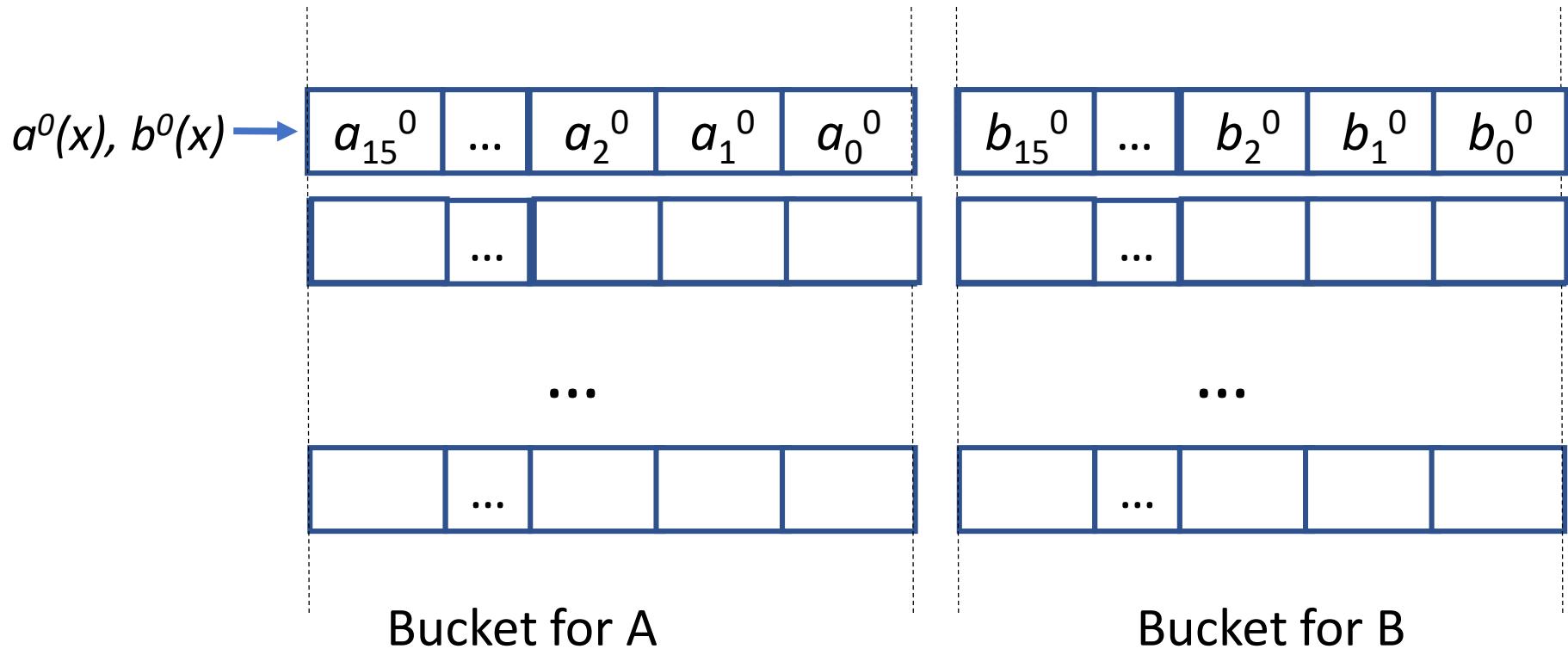


Bucket for A.

- Row contains one 16 coeff. Pol
- There are 16 rows

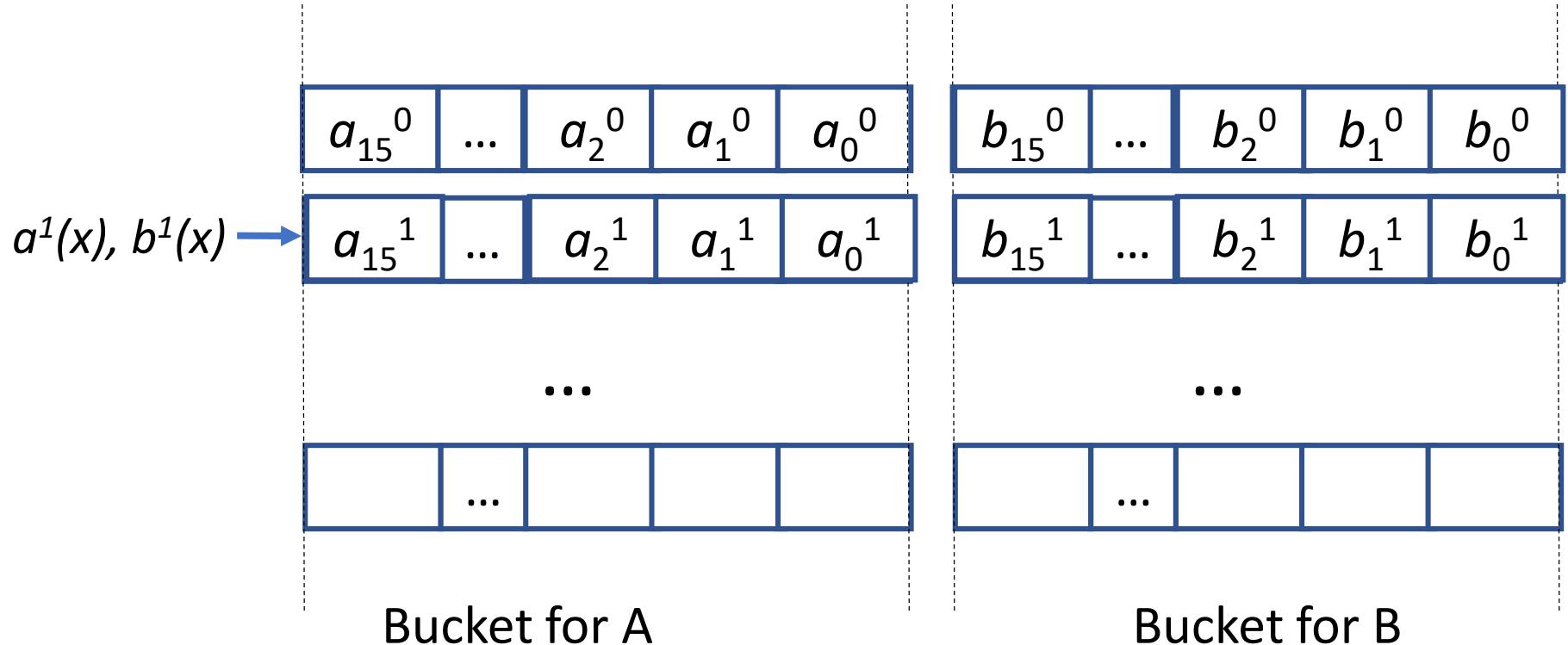
Total 256 coeff. In the bucket

Buckets for 16 coeff. polynomials



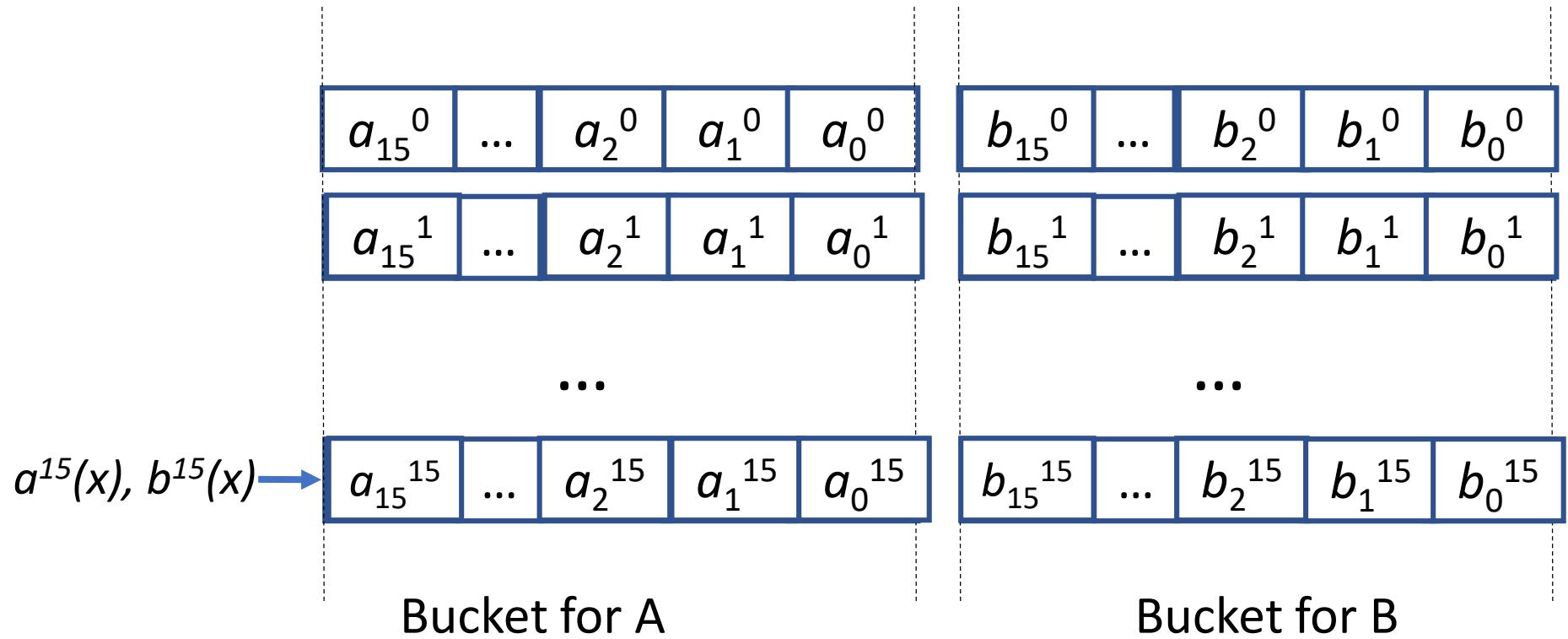
- Multiplications are not performed immediately
- The buckets are filled first

Buckets for 16 coeff. polynomials



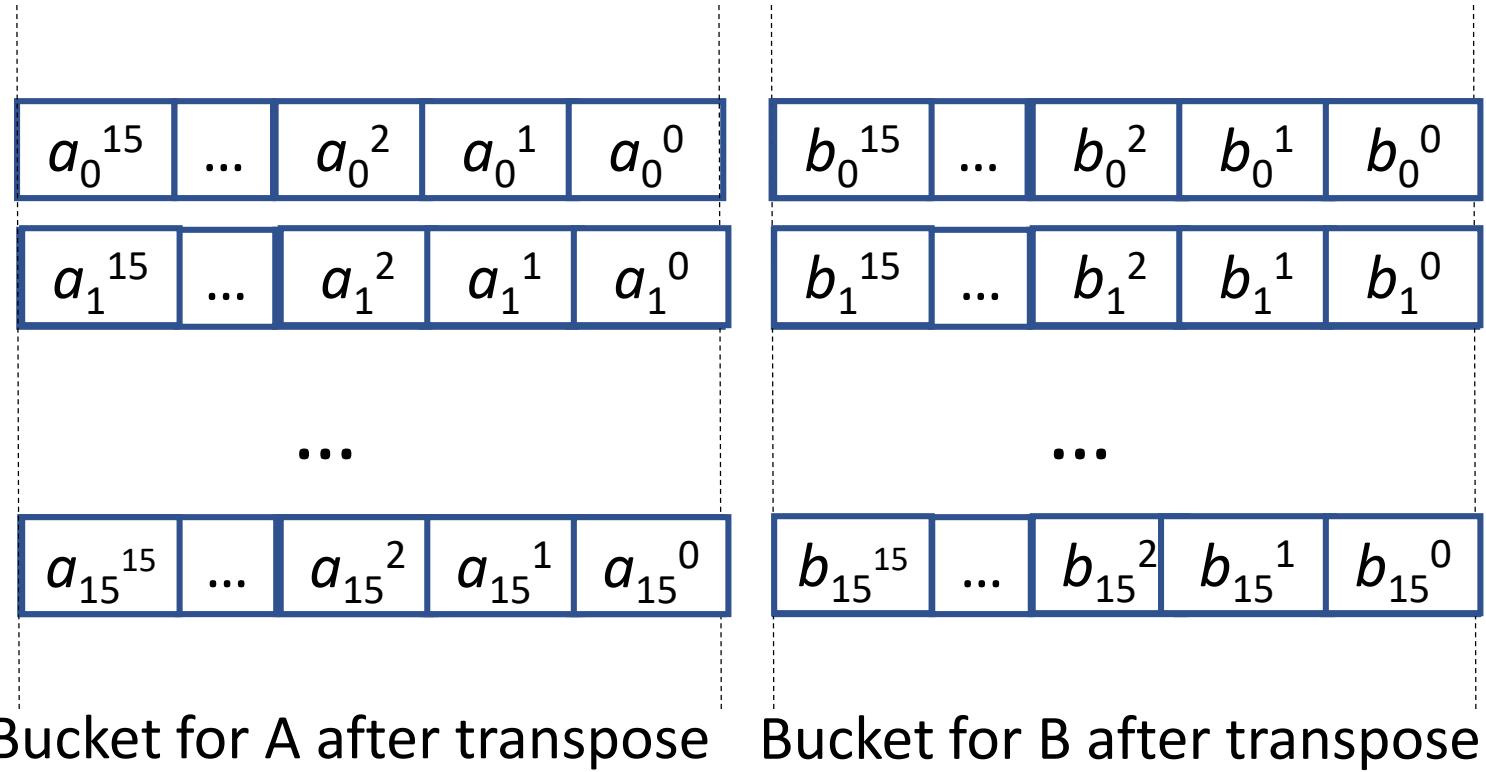
- Multiplications are not performed immediately
- The buckets are filled first

Buckets for 16 coeff. polynomials



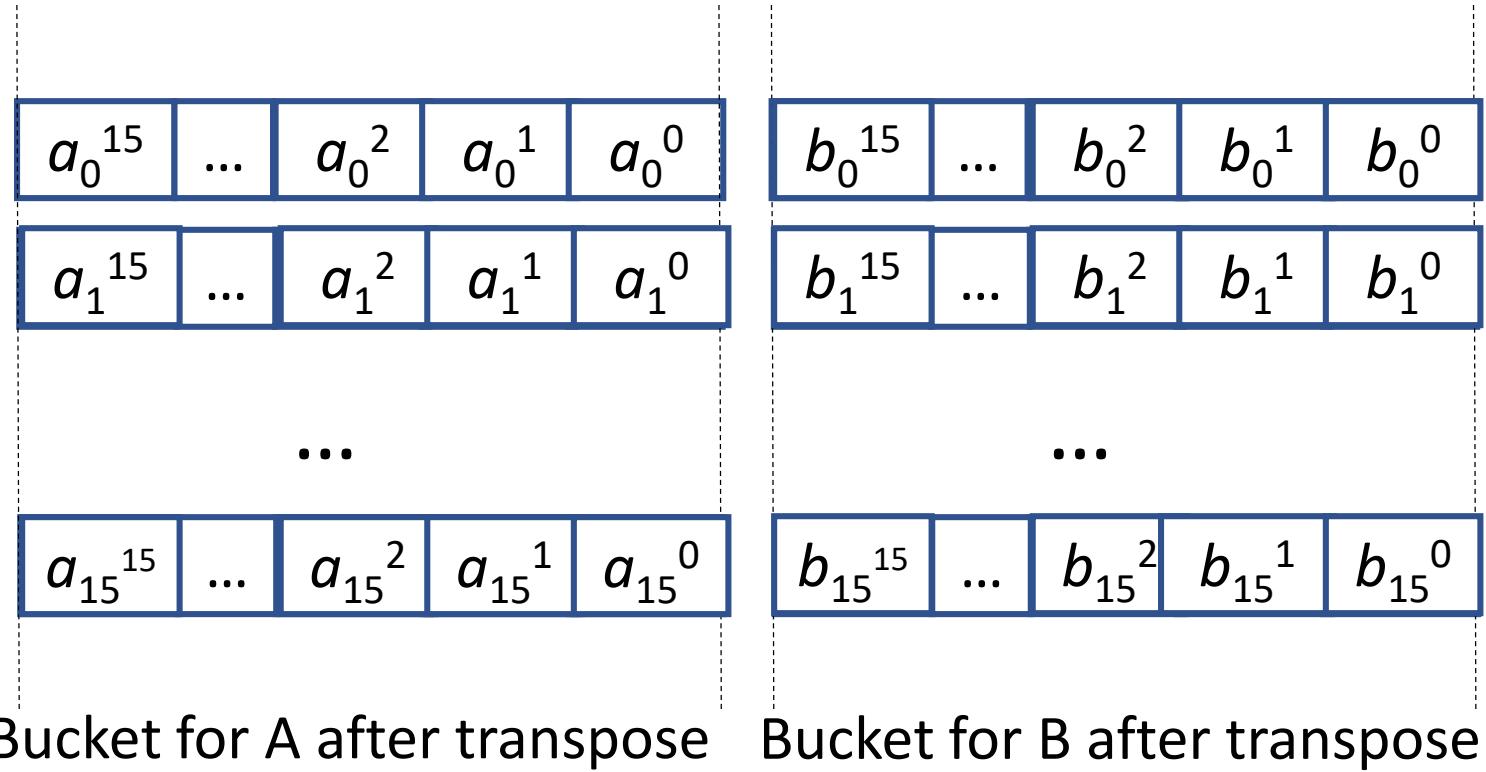
- Multiplications are not performed immediately
- The buckets are filled first

Buckets for 16 coeff. polynomials



- Transpose each bucket
- Transpose using AVX is an interesting optimization problem

Buckets for 16 coeff. polynomials

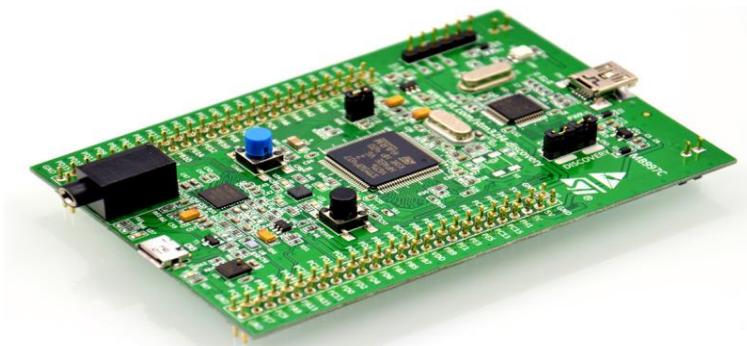


- Now multiply vectors using AVX instruction
- Obtain C which is a vector length 32, containing $32 * 16$ coeffs.
- Finally, transpose C to get 16 result polynomials

School-book multiplication using DSP instructions

Cortex-M4: STM32F4-discovery by STMicroelectronics

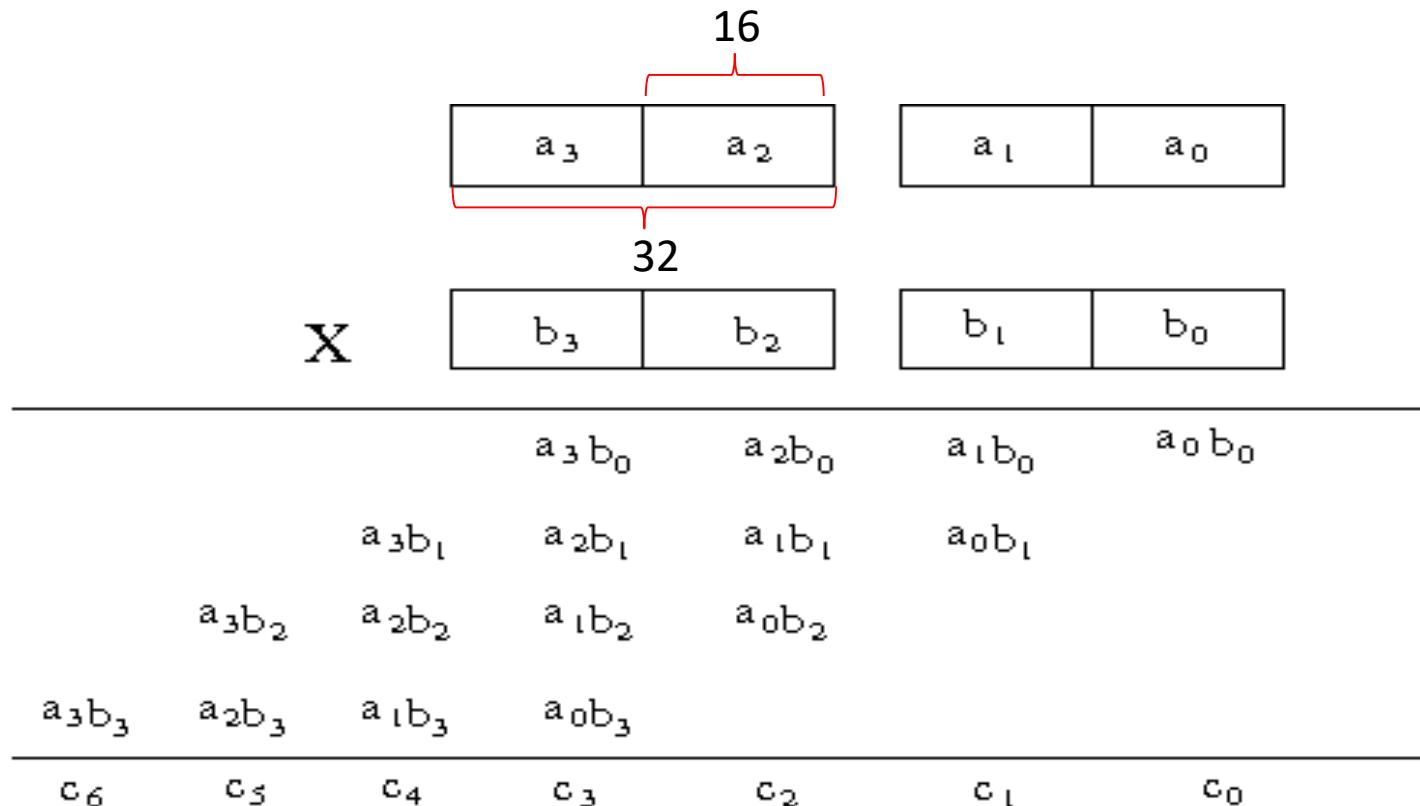
- DSP instructions
- 14 registers fully available
- 192 KB of RAM



School-book multiplication using SMLA instruction

Multiplications between half words

$$\text{SMLA}(B/T)(B/T)(r^a, r^b, r^c, r^d) := r^a \leftarrow r^b_{0/1} * r^c_{0/1} + r^d$$

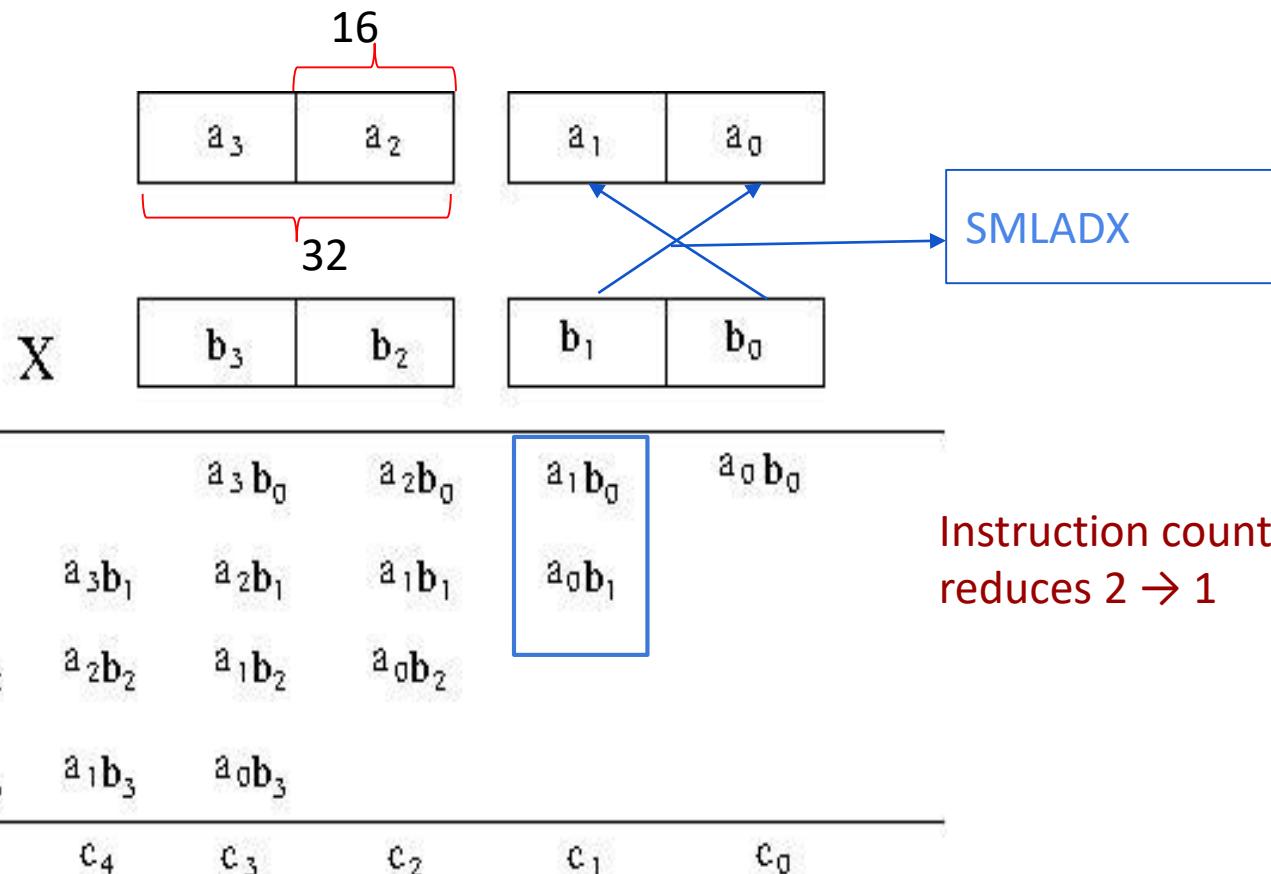


16 instructions !

School-book multiplication using SMLA instruction

DSP instruction for cross multiplication of half words

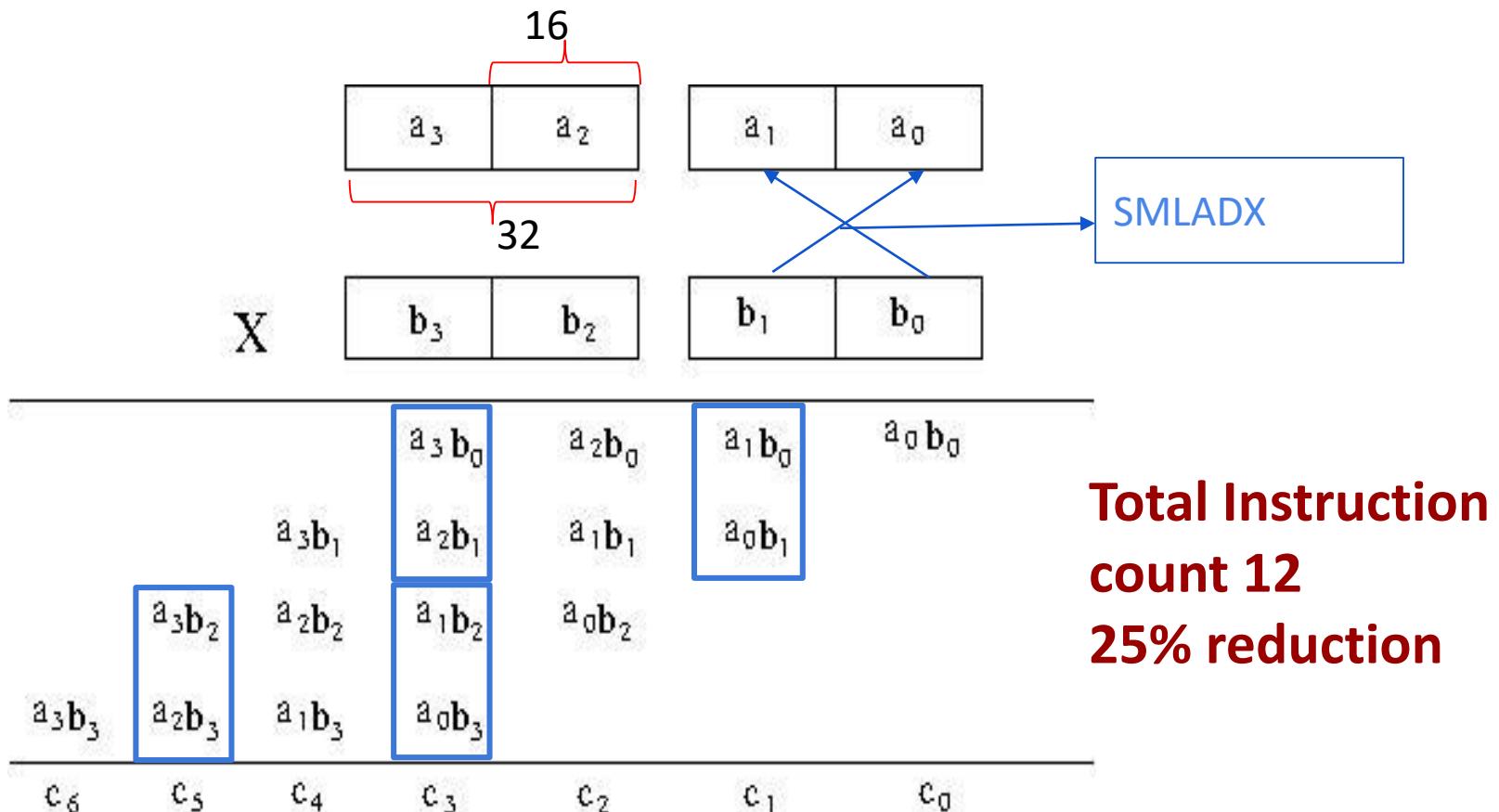
$$\text{SMLADX}(r^a, r^b, r^c, r^d) := r^a \leftarrow r^b_0 * r^c_1 + r^b_1 * r^c_0 + r^d$$



School-book multiplication using SMLA instruction

DSP instruction for cross multiplication of half words

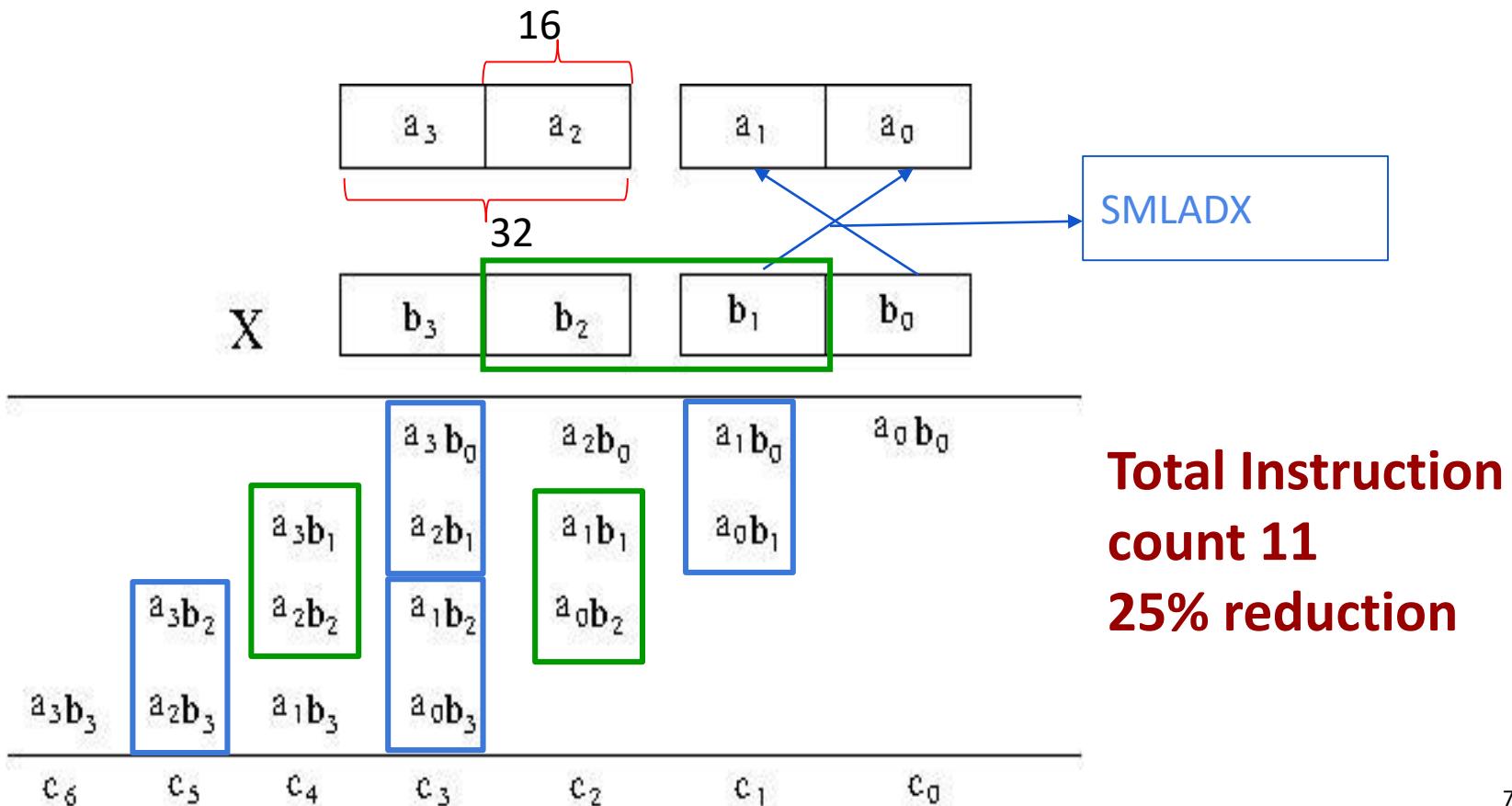
$$\text{SMLADX}(r^a, r^b, r^c, r^d) := r^a \leftarrow r^b_0 * r^c_1 + r^b_1 * r^c_0 + r^d$$



School-book multiplication using SMLA instruction

Pack non-adjacent coefficients in spare register using PKHBT

Apply SMLADX again



For 16x16 Schoolbook multiplication → 37.5% reduction overall

SABER: High-level optimization

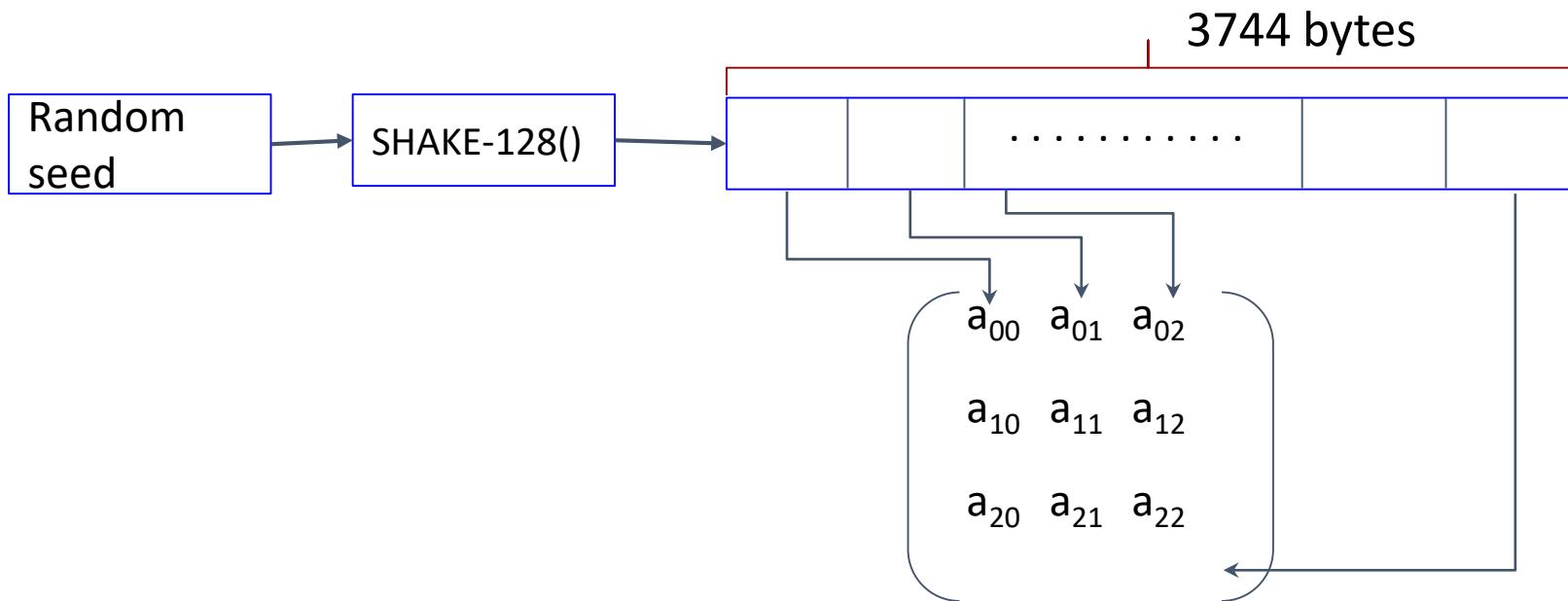
$$\left[\frac{p}{q} \begin{pmatrix} a_{0,0}(x) & \dots & a_{0,k-1}(x) \\ \dots & \dots & \dots \\ a_{k-1,0}(x) & \dots & a_{k-1,k-1}(x) \end{pmatrix} * \begin{pmatrix} s_0(x) \\ \dots \\ s_{k-1}(x) \end{pmatrix} \right] \approx \begin{pmatrix} b_0(x) \\ b_{k-1}(x) \end{pmatrix} \pmod{x^{256} + 1}$$

Matrices and vectors are generated by
expanding random seed using XOF

Matrix and vector generation

- The public matrix A, secrete vectors s and s' require large number of random bytes

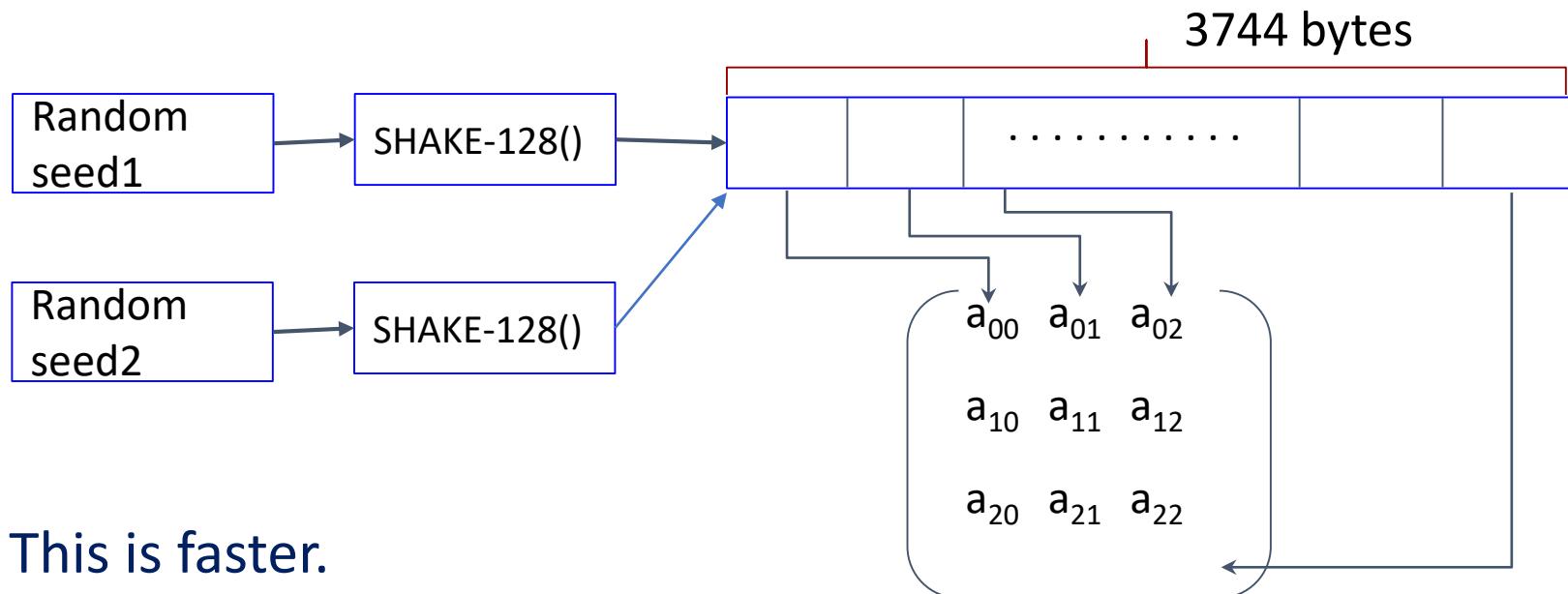
Option 1: Use single SHAKE and generate RAND sequentially



Matrix and vector generation

- The public matrix A, secrete vectors s and s' require large number of random bytes

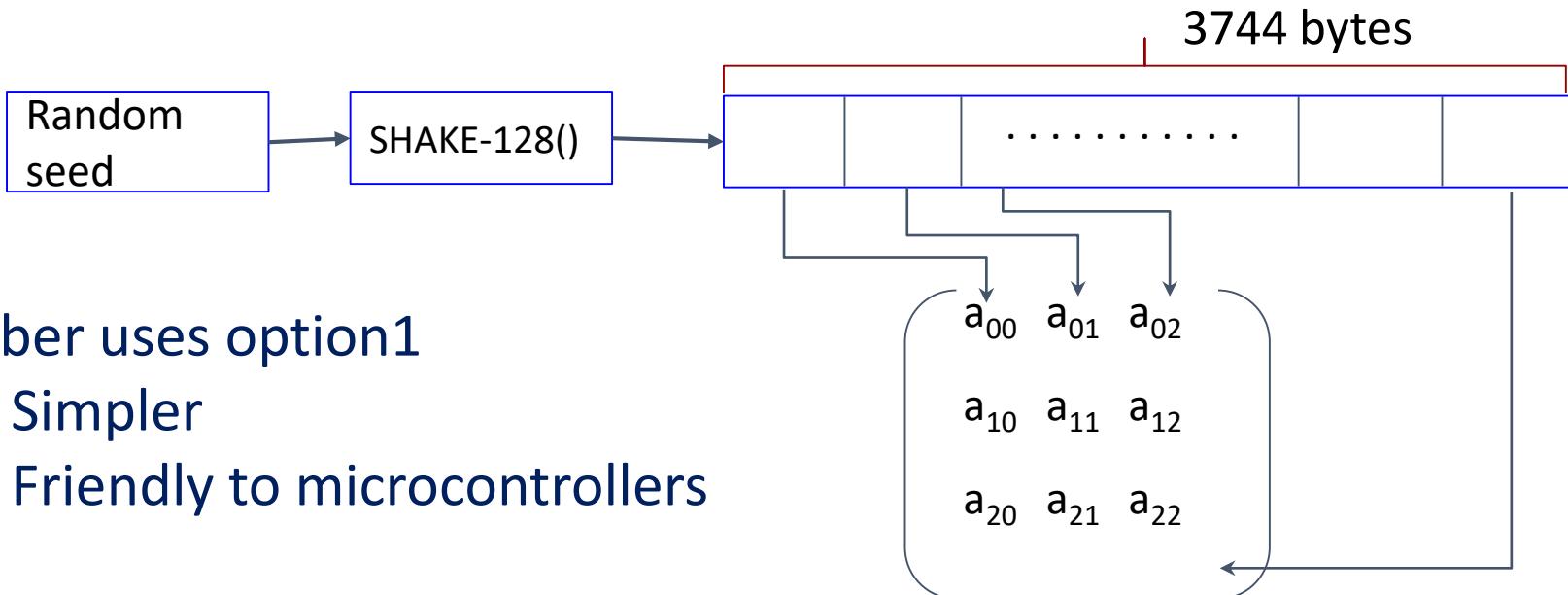
Option 2: Use multiple SHAKE and generate RAND parally



Matrix and vector generation

- The public matrix A, secrete vectors s and s' require large number of random bytes

Option 1: Use single SHAKE and generate RAND sequentially

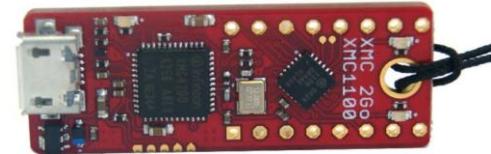


Saber uses option1

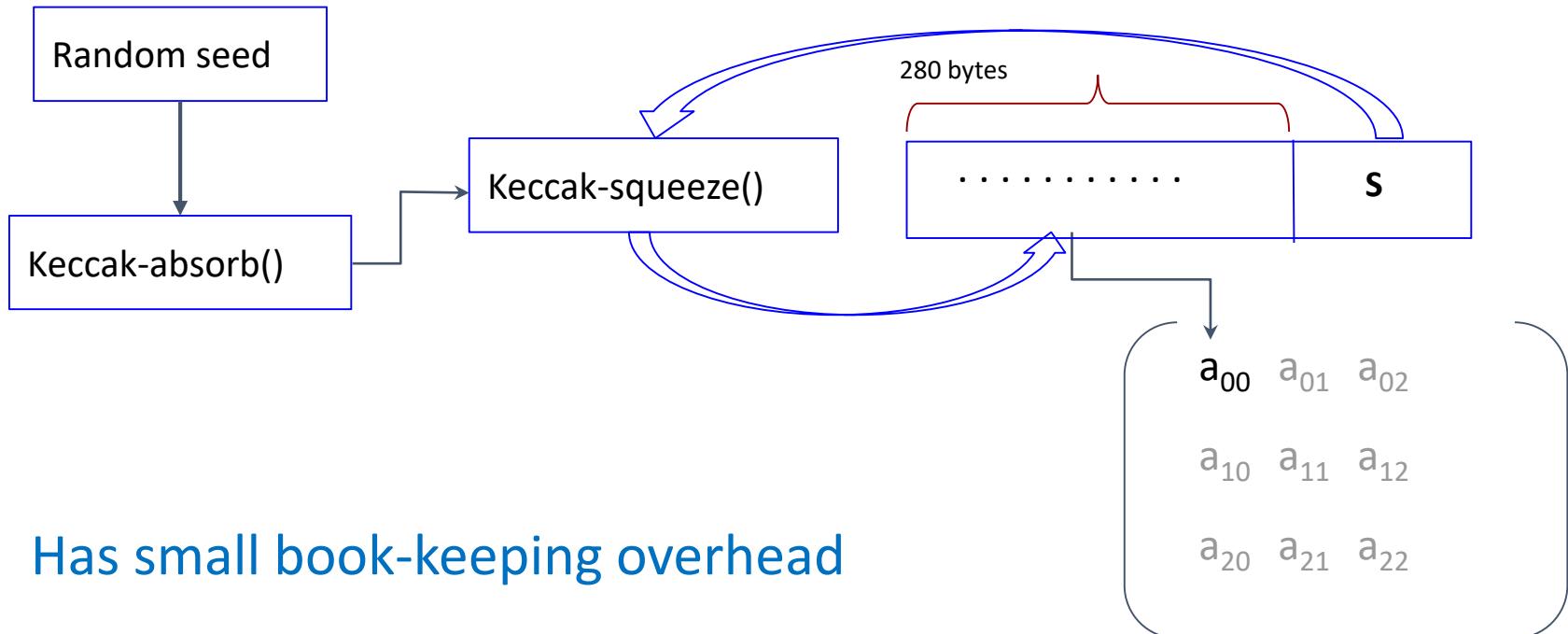
- Simpler
- Friendly to microcontrollers
- Hardware: single Keccak core consumes roughly 50% of area in LPR encryption scheme

Matrix and vector generation: memory efficient

- ‘Just in time’ approach



Cortex M0 with 8KB RAM



Matrix and vector Multiplication

$A \leftarrow \text{XOF}(\text{seed}_A)$

Keygen:

$s \leftarrow \text{XOF}(\text{seed}_s)$

Computes A^*s

Encryption:

$s' \leftarrow \text{XOF}(\text{seed}_{s'})$

Computes $A^T * s'$

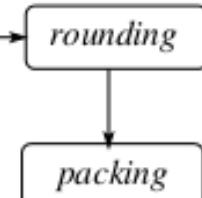
Order of A is different in KeyGen and Enc

Matrix and vector Multiplication

$$\begin{pmatrix} (1) & (2) & (3) \\ & & \end{pmatrix} \cdot \begin{pmatrix} (0) \\ (0) \\ (0) \\ s \end{pmatrix} = \begin{pmatrix} (0) \\ (0) \end{pmatrix}$$

A

Row-major generation of matrix



$$\begin{pmatrix} (1) \\ (2) \\ (3) \\ & \end{pmatrix} \cdot \begin{pmatrix} (0) \\ (0) \\ (0) \\ s \end{pmatrix} = \begin{pmatrix} (0) \\ (0) \\ (0) \\ s' \end{pmatrix}$$

A^T

Column-major generation of matrix

Intermediate polynomials

A is generated ‘one polynomial at a time’

- Row-major has smaller memory requirement

Enc is costlier than KeyGen.

- Use row-major for Enc and col-major for KeyGen
- Round2 spex of Saber does this

Saber: results

- Secret key size 1344 bytes
- Public key size 992 bytes
- Ciphertext size 1088 bytes

Performance on Intel Haswell

Implementation	Key Generation	Encapsulation	Decapsulation
AVX2	104 K	122 K	120 K

Performance on ARM Cortex M (CHES 2018)

Platform	Key Generation	Encapsulation	Decapsulation	Stack
Cortex M4 (DSP)	1.1 M	1.5 M	1.6 M	8 KB
Cortex M0	4.7 M	6.3 M	7.5 M	6 KB

Work in progress: hardware implementation of Saber

Thank you