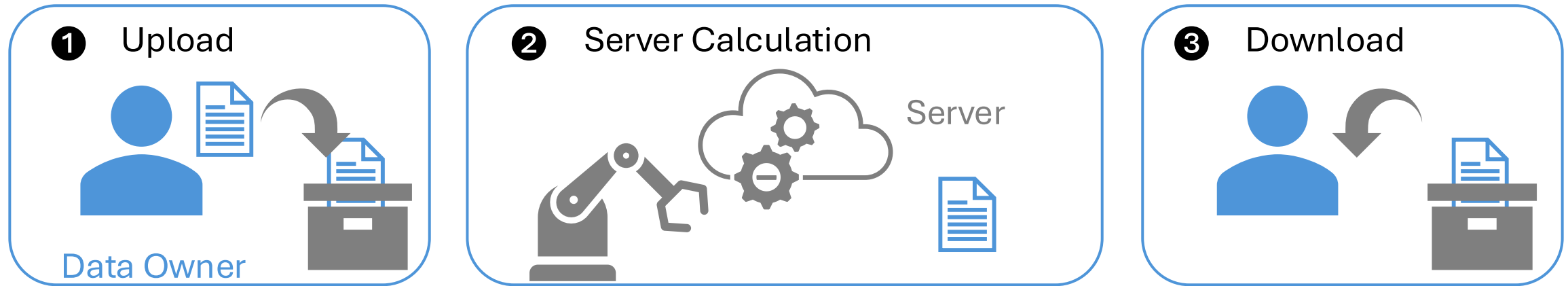


# ReliaFHE: Resilient Design for Fully Homomorphic Encryption Accelerators

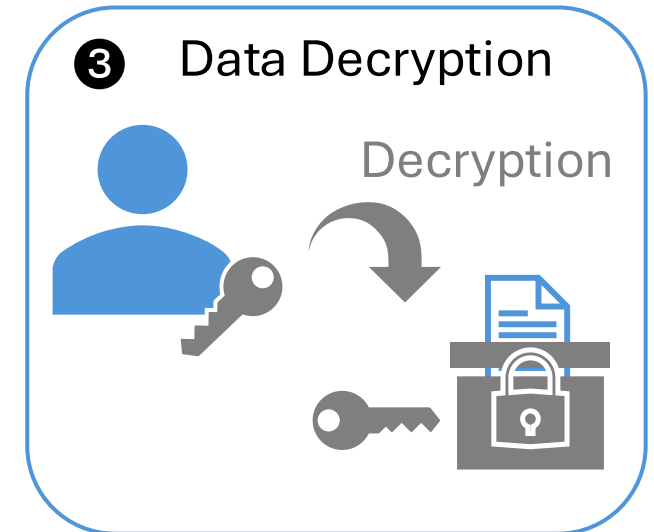
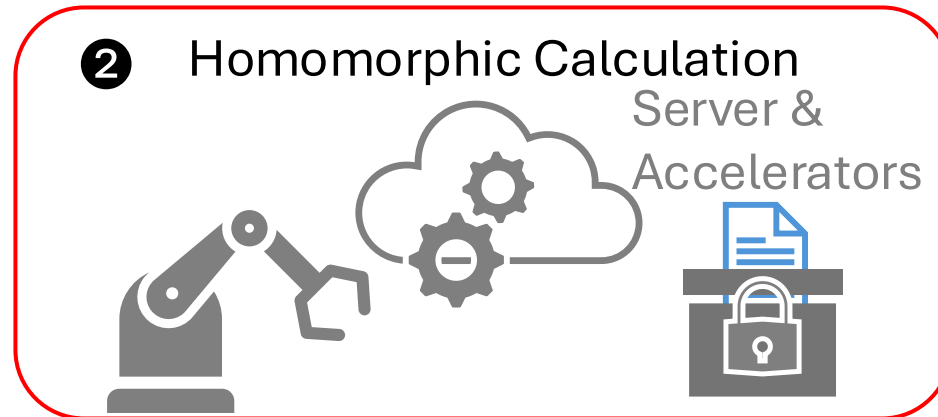
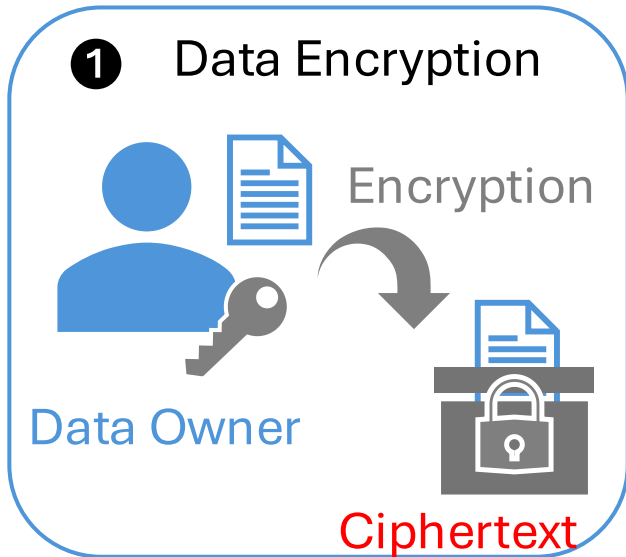
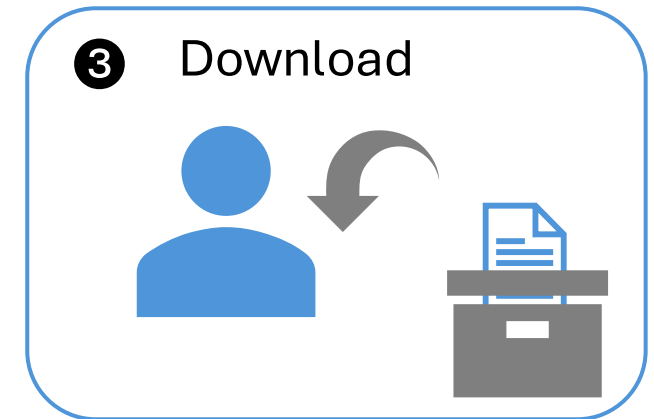
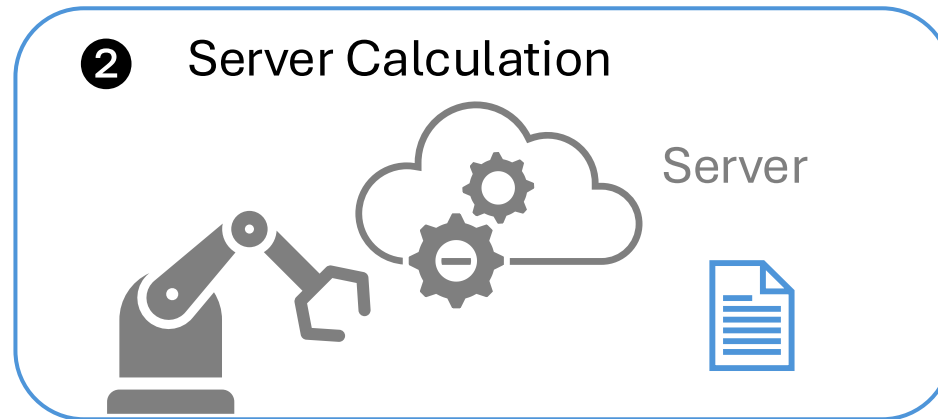
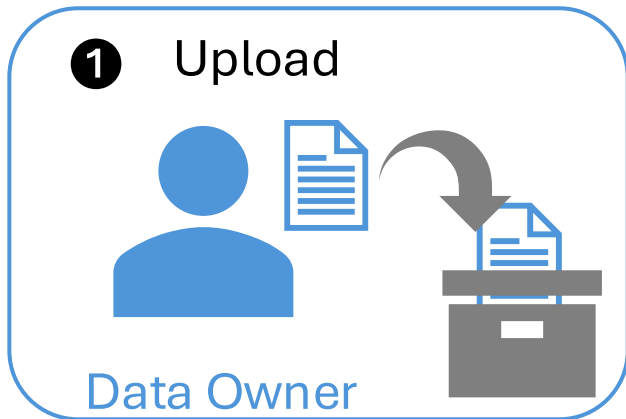
Fan Li, Mayank Kumar, Ruizhi Zhu, Mengxin Zheng, Qian Lou, Xin Xin\*

College of Engineering and Computer Science  
University of Central Florida

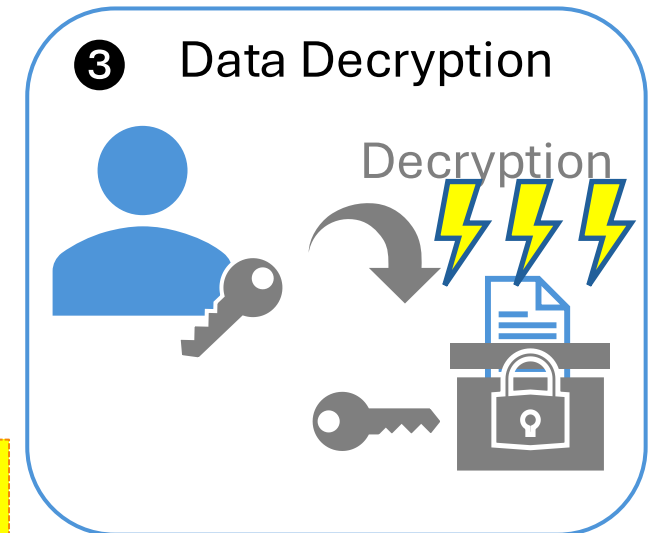
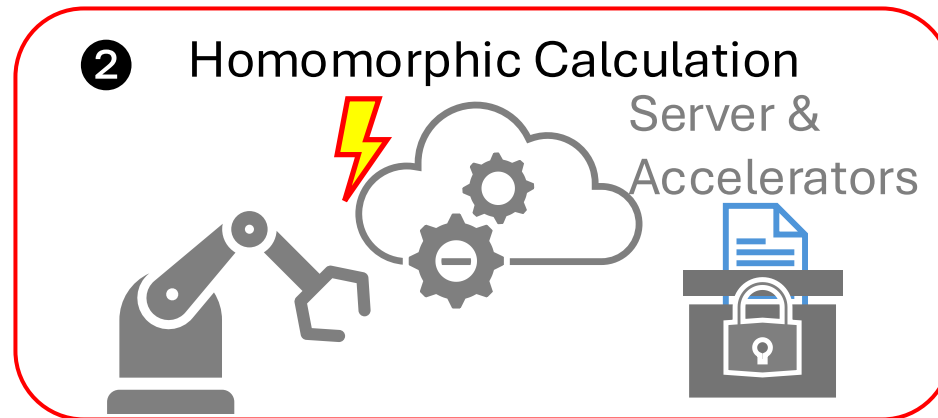
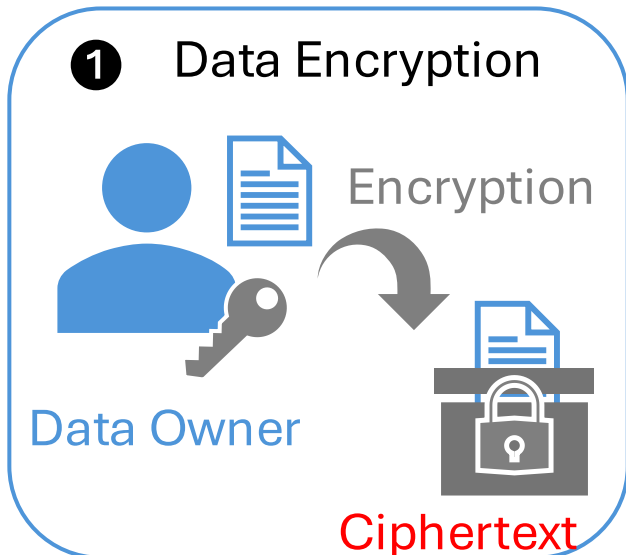
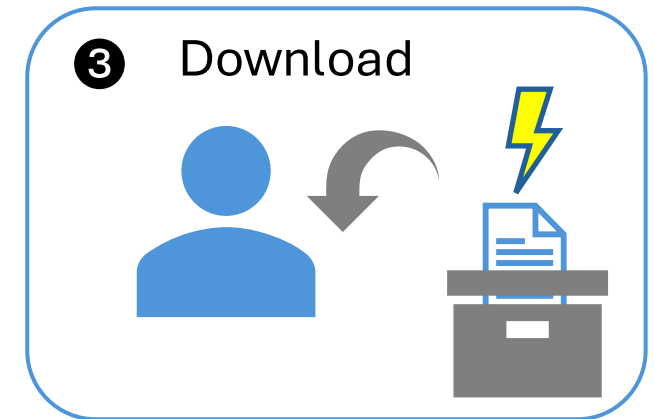
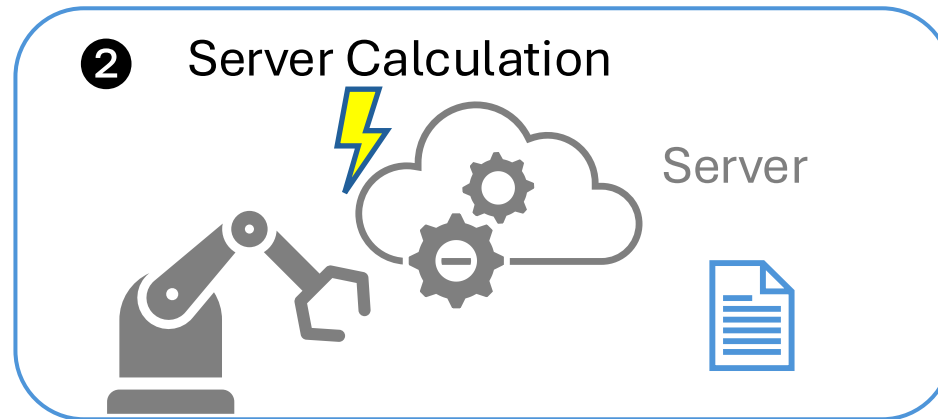
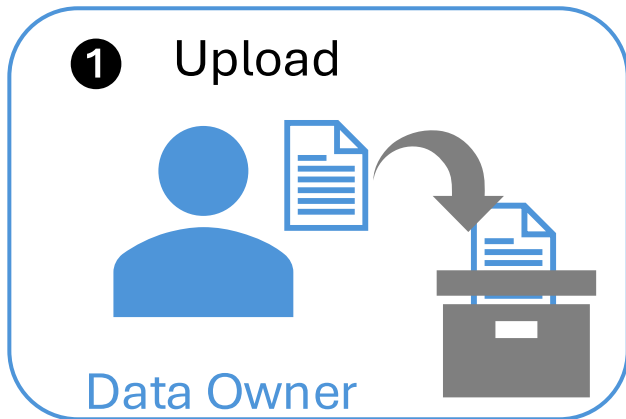
# Fully Homomorphic Encryption (FHE)



# Fully Homomorphic Encryption (FHE)



# Fully Homomorphic Encryption (FHE)

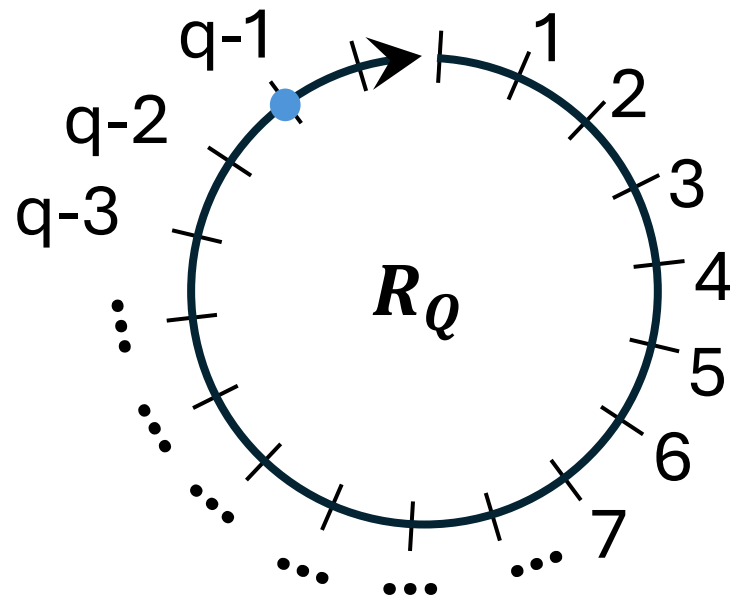


**FHE is more error sensitive**

# Error-Sensitivity Analysis

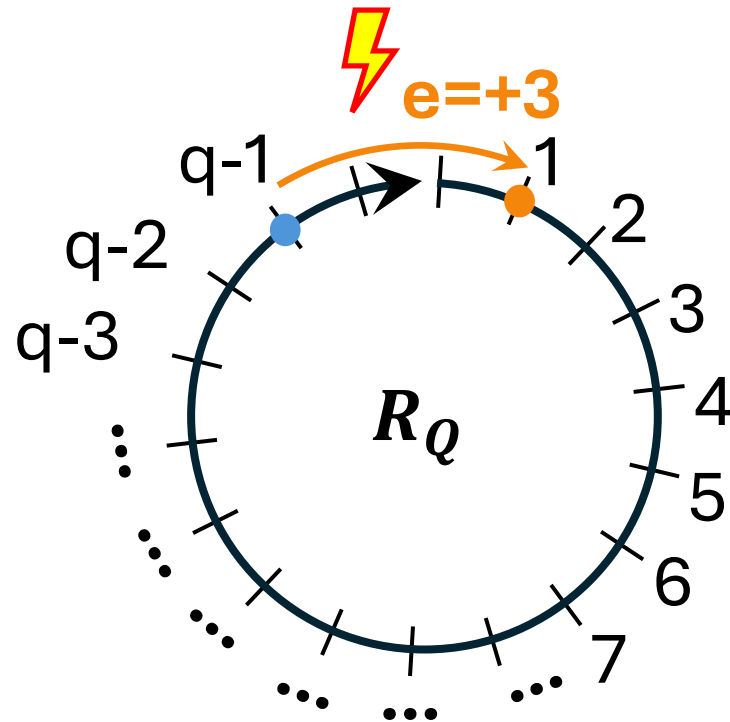
- Modulus Operations
- NTT Representation
- Packing Algorithm

# Error-Sensitivity Analysis



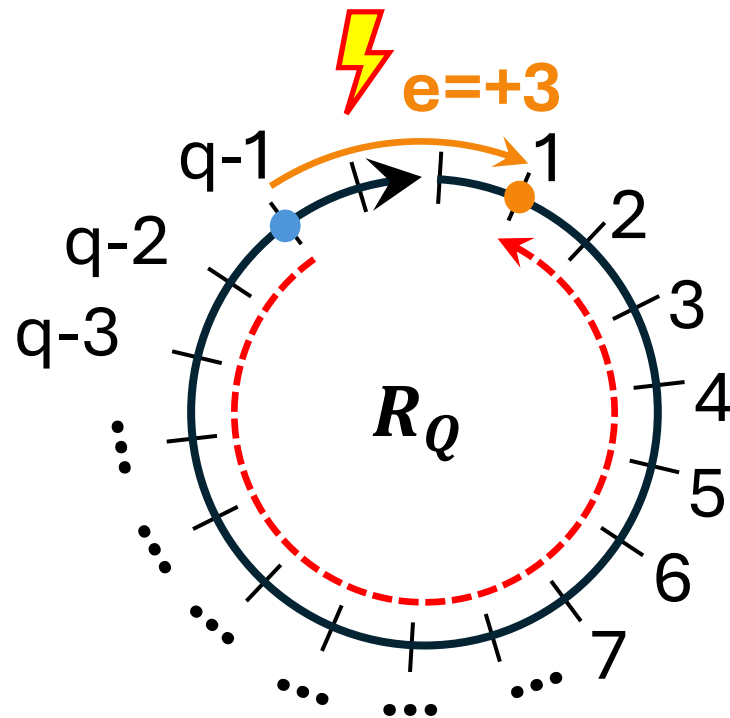
- **Modulus Operations**
- NTT Representation
- Packing Algorithm

# Error-Sensitivity Analysis



- **Modulus Operations**
- NTT Representation
- Packing Algorithm

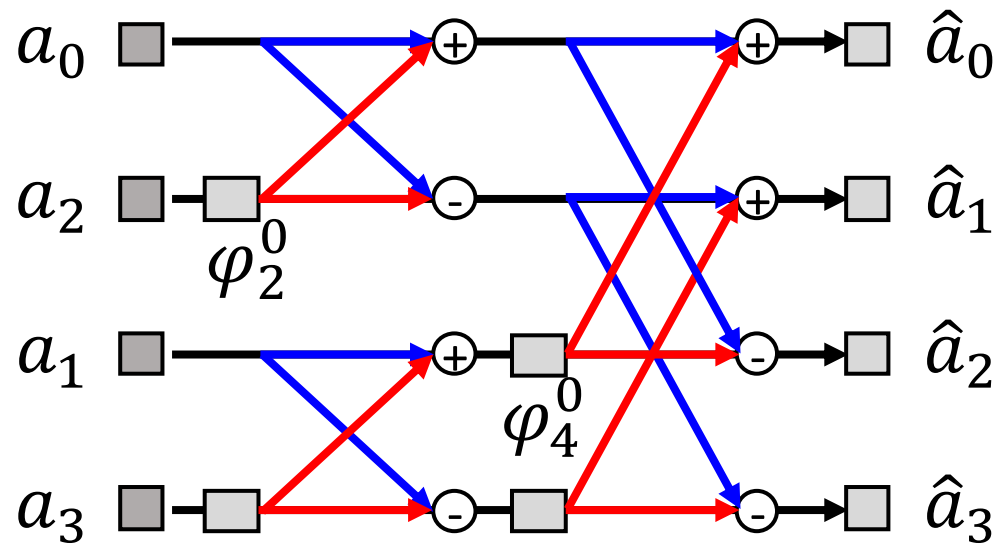
# Error-Sensitivity Analysis



## Error amplification

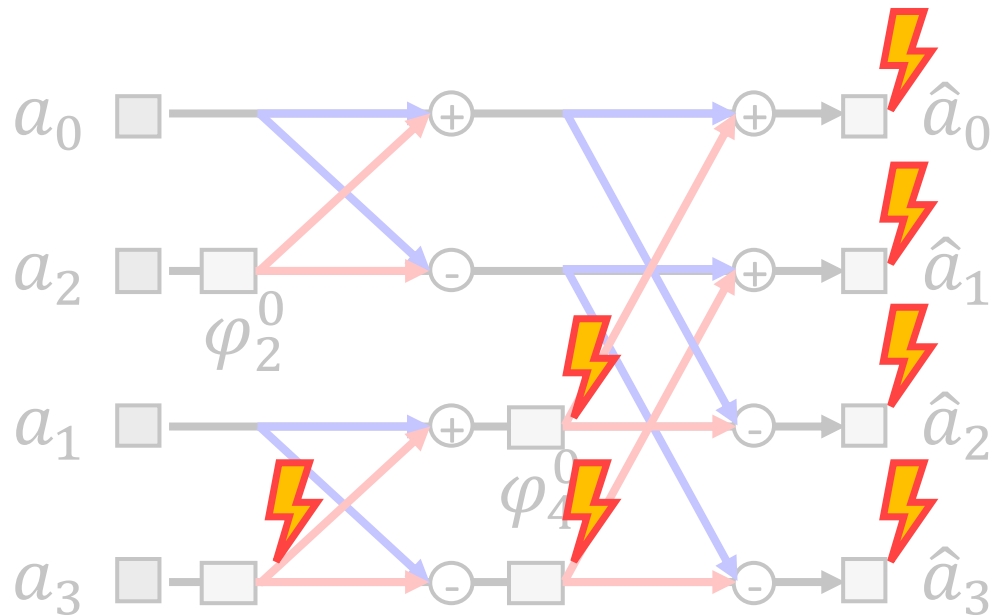
- Modulus Operations
- NTT Representation
- Packing Algorithm

# Error-Sensitivity Analysis



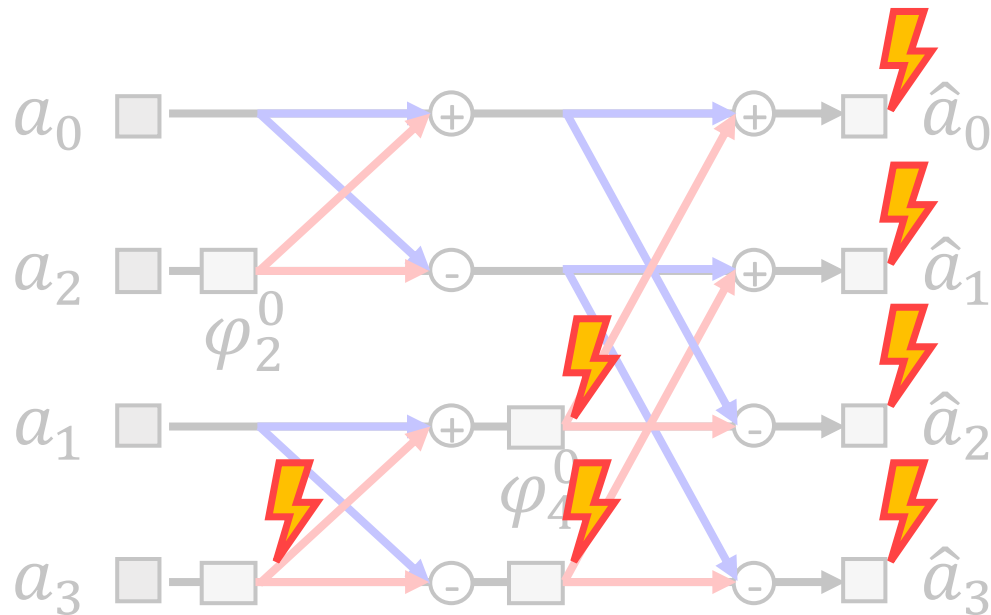
- Modulus Operations
- **NTT Representation**
- Packing Algorithm

# Error-Sensitivity Analysis



- Modulus Operations
- **NTT Representation**
- Packing Algorithm

# Error-Sensitivity Analysis



## Error broadcast

- Modulus Operations
- **NTT Representation**
- Packing Algorithm

# Error-Sensitivity Analysis

## Plaintext

$$\begin{bmatrix} m_{0,0}^{(a)} & m_{0,1}^{(a)} & \dots \\ m_{1,0}^{(a)} & m_{1,1}^{(a)} \text{⚡} & \dots \\ \dots & \dots & \dots \end{bmatrix} \cdot \begin{bmatrix} m_0^b \\ m_1^b \\ \dots \end{bmatrix} = \begin{bmatrix} m_0^{(c)} \\ m_1^{(c)} \text{⚡} \\ \dots \end{bmatrix}$$

- Modulus Operations
- NTT Representation
- **Packing Algorithm**

# Error-Sensitivity Analysis

## Plaintext

$$\begin{bmatrix} m_{0,0}^{(a)} & m_{0,1}^{(a)} & \dots \\ m_{1,0}^{(a)} & m_{1,1}^{(a)} \text{⚡} & \dots \\ \dots & \dots & \dots \end{bmatrix} \cdot \begin{bmatrix} m_0^b \\ m_1^b \\ \dots \end{bmatrix} = \begin{bmatrix} m_0^{(c)} \\ m_1^{(c)} \text{⚡} \\ \dots \end{bmatrix}$$

- Modulus Operations
- NTT Representation

## Ciphertext

$$\left( \begin{bmatrix} ct_0^{(a)} \\ ct_1^{(a)} \\ \dots \end{bmatrix} + \begin{bmatrix} ct_0^{(a)} \\ ct_1^{(a)} \text{⚡} \\ \dots \end{bmatrix} + \dots \right) \cdot \begin{bmatrix} ct_0^{(b)} \\ ct_1^{(b)} \\ \dots \end{bmatrix} \xrightarrow{Dec} \begin{bmatrix} m_0^{(c)} \\ m_1^{(c)} \\ \dots \end{bmatrix}$$

- **Packing Algorithm**

# Error-Sensitivity Analysis

## Plaintext

$$\begin{bmatrix} m_{0,0}^{(a)} & m_{0,1}^{(a)} & \dots \\ m_{1,0}^{(a)} & m_{1,1}^{(a)} \text{⚡} & \dots \\ \dots & \dots & \dots \end{bmatrix} \cdot \begin{bmatrix} m_0^b \\ m_1^b \\ \dots \end{bmatrix} = \begin{bmatrix} m_0^{(c)} \\ m_1^{(c)} \text{⚡} \\ \dots \end{bmatrix}$$

- Modulus Operations
- NTT Representation

## Ciphertext

$$\left( \begin{bmatrix} ct_0^{(a)} \\ ct_1^{(a)} \\ \dots \end{bmatrix} + \begin{bmatrix} ct_0^{(a)} \\ ct_1^{(a)} \text{⚡} \\ \dots \end{bmatrix} + \dots \right) \cdot \begin{bmatrix} ct_0^{(b)} \\ ct_1^{(b)} \\ \dots \end{bmatrix} \xrightarrow{Dec} \begin{bmatrix} m_0^{(c)} \\ m_1^{(c)} \\ \dots \end{bmatrix}$$

- **Packing Algorithm**

# Error-Sensitivity Analysis

## Plaintext

$$\begin{bmatrix} m_{0,0}^{(a)} & m_{0,1}^{(a)} & \dots \\ m_{1,0}^{(a)} & m_{1,1}^{(a)} & \dots \\ \dots & \dots & \dots \end{bmatrix} \cdot \begin{bmatrix} m_0^b \\ m_1^b \\ \dots \end{bmatrix} = \begin{bmatrix} m_0^{(c)} \\ m_1^{(c)} \\ \dots \end{bmatrix}$$

## Error broadcast

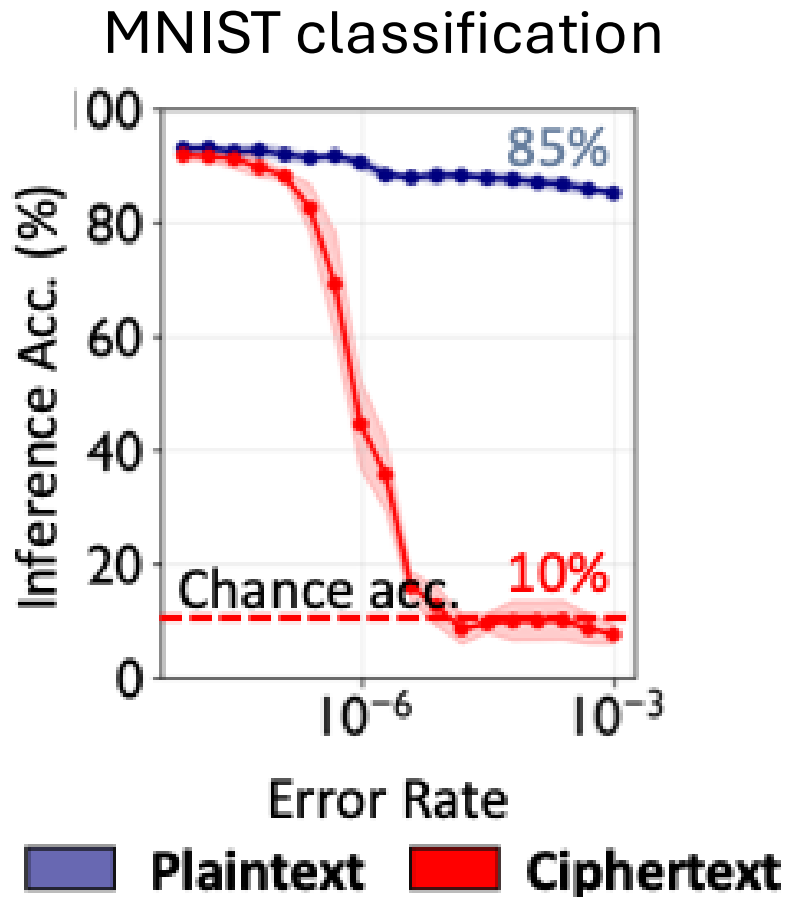
- Modulus Operations
- NTT Representation

## Ciphertext

$$\left( \begin{bmatrix} ct_0^{(a)} \\ ct_1^{(a)} \\ \dots \end{bmatrix} + \begin{bmatrix} ct_0^{(a)} \\ ct_1^{(a)} \\ \dots \end{bmatrix} + \dots \right) \cdot \begin{bmatrix} ct_0^{(b)} \\ ct_1^{(b)} \\ \dots \end{bmatrix} \xrightarrow{Dec} \begin{bmatrix} m_0^{(c)} \\ m_1^{(c)} \\ \dots \end{bmatrix}$$

- Packing Algorithm

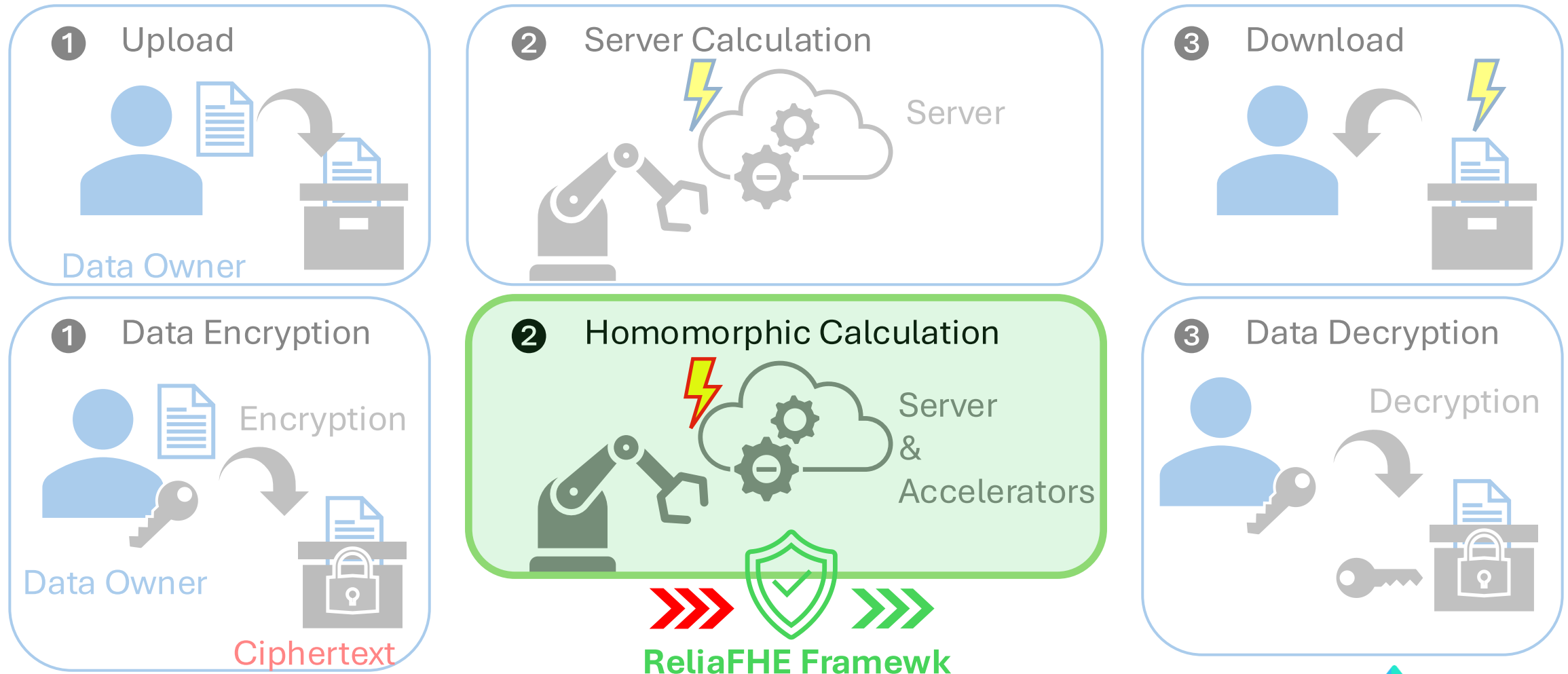
# Error-Sensitivity Analysis



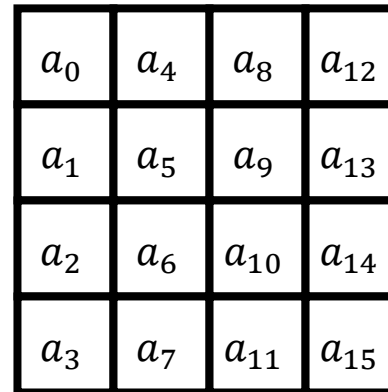
**FHE-based model is more error sensitive**

- Modulus Operations
- NTT Representation
- Packing Algorithm

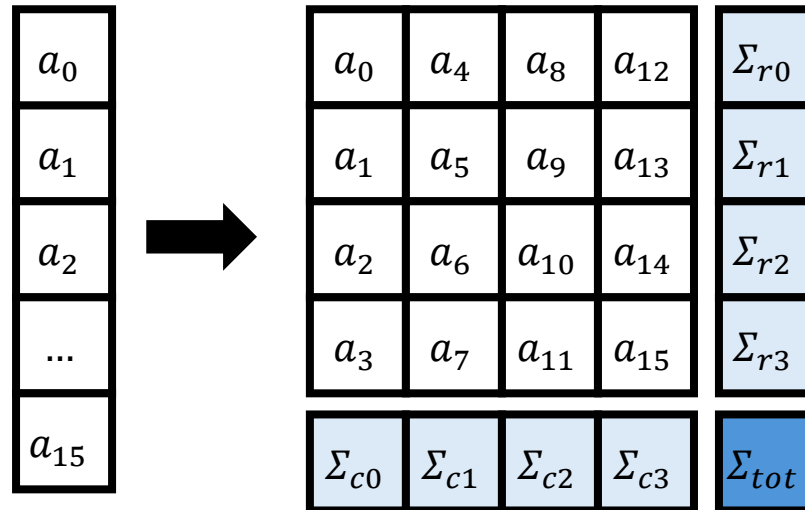
# ReliaFHE Framework



# Storage Protection(e.g. coefficient storage)



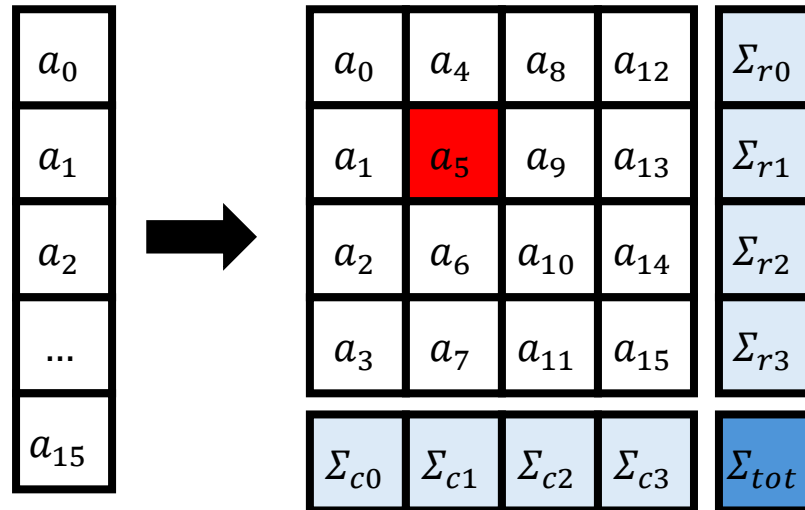
# Storage Protection(e.g. coefficient storage)



$\Sigma_r, \Sigma_c$  : row and col-wise Checksum

$\Sigma_{tot}$ : total Checksum

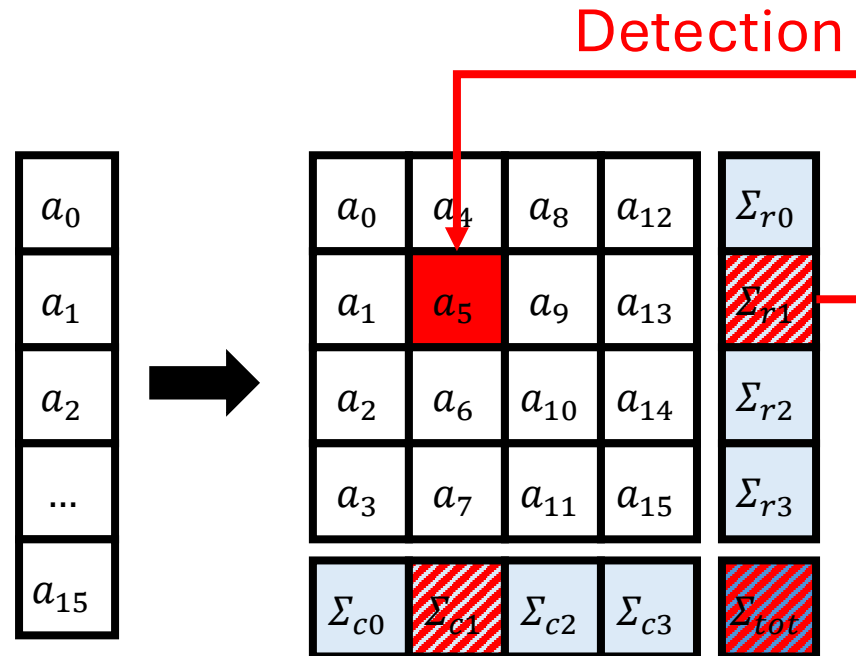
# Storage Protection(e.g. coefficient storage)



$\Sigma_r, \Sigma_c$  : row and col-wise Checksum

$\Sigma_{tot}$ : total Checksum

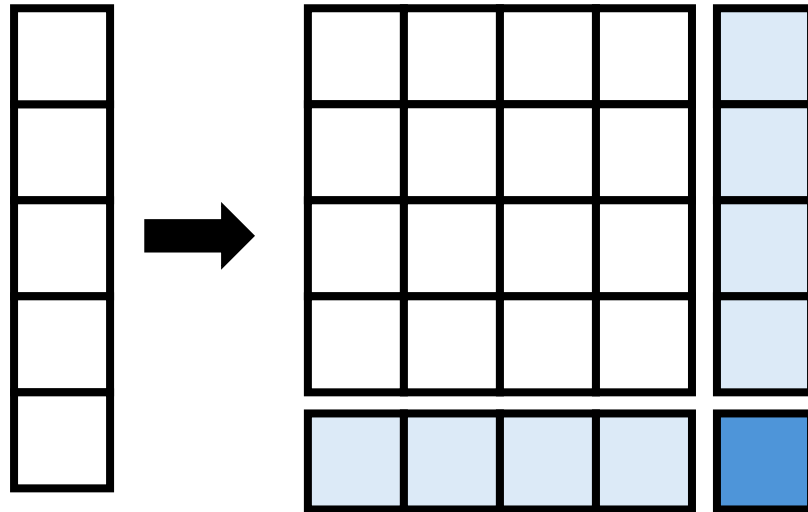
# Storage Protection(e.g. coefficient storage)



$\Sigma_r, \Sigma_c$  : row and col-wise Checksum

$\Sigma_{tot}$ : total Checksum

# Storage Protection(e.g. coefficient storage)



Redundancy:  $O(\sqrt{n})$

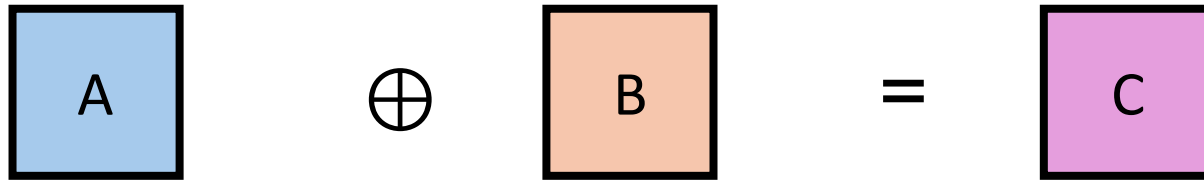
e.g.

$$|v|=65536$$

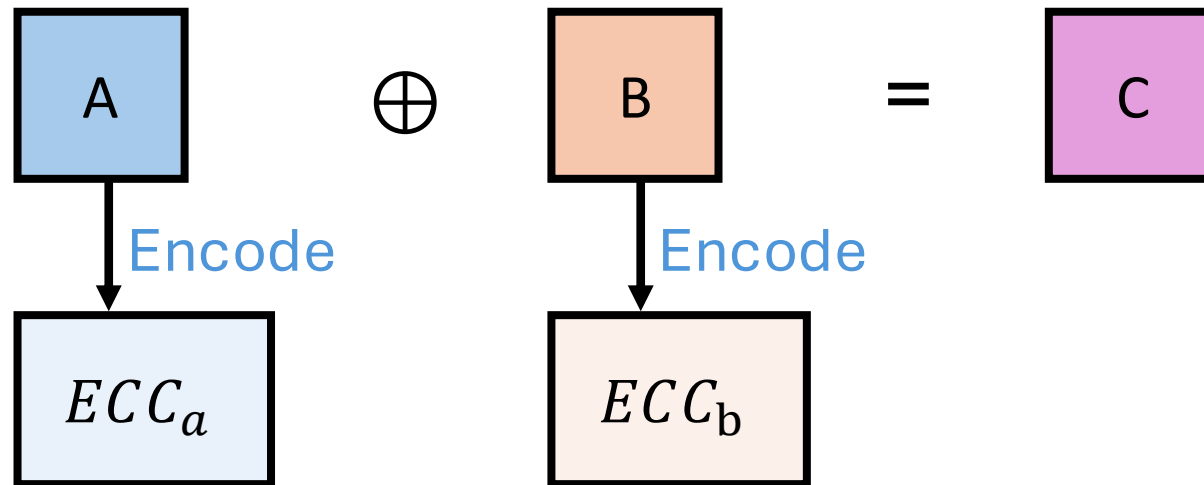
$$|R|= 256$$

0.3% Redundancy

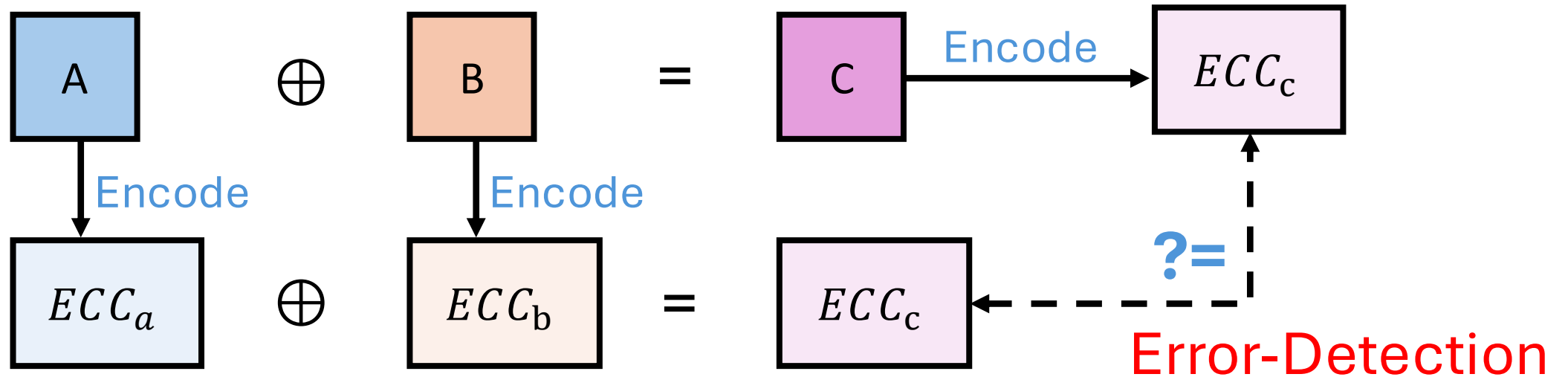
# Computation Protection



# Computation Protection



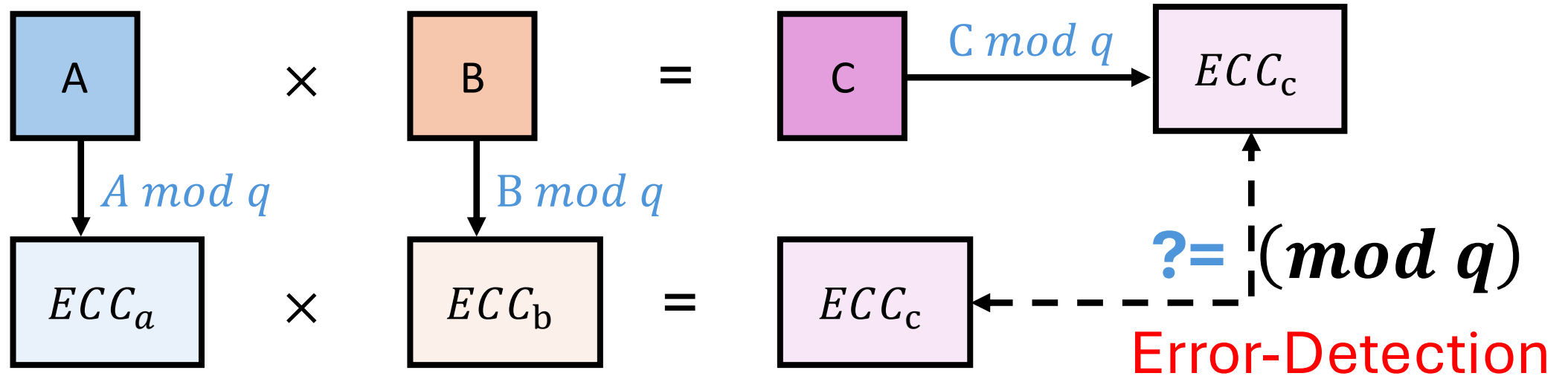
# Computation Protection



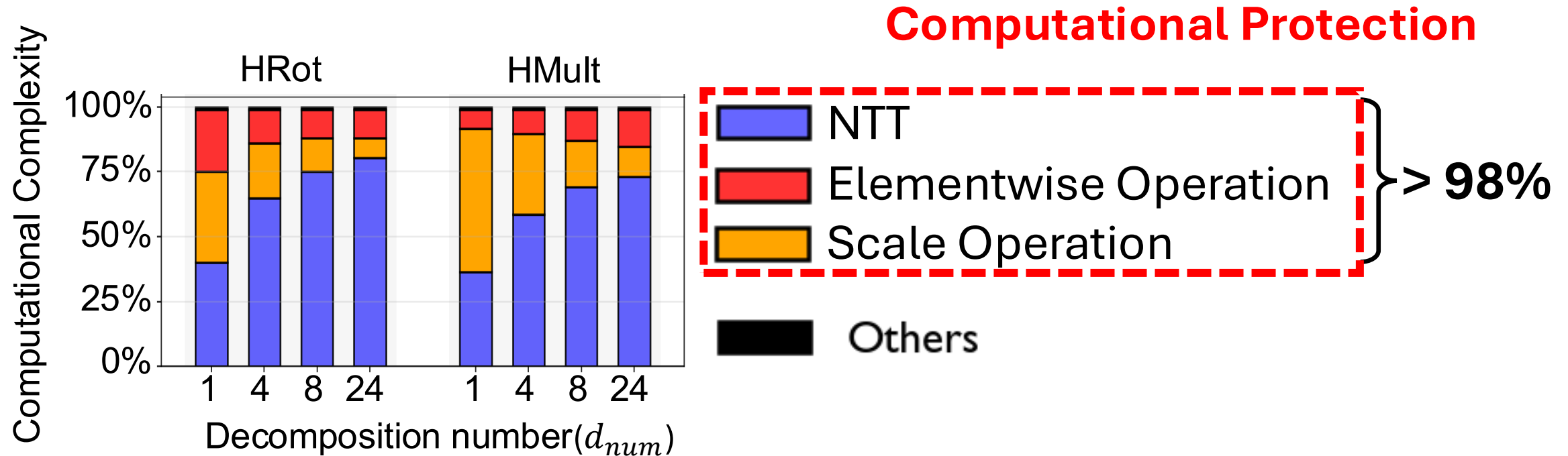
**Construct computational compatible ECC**

# Computation Protection

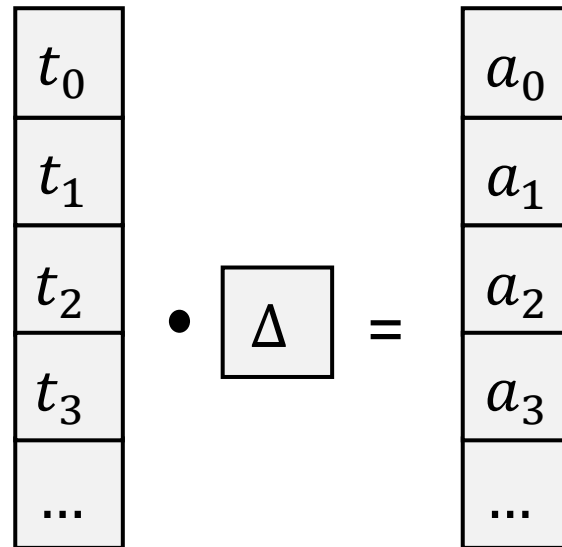
e.g.



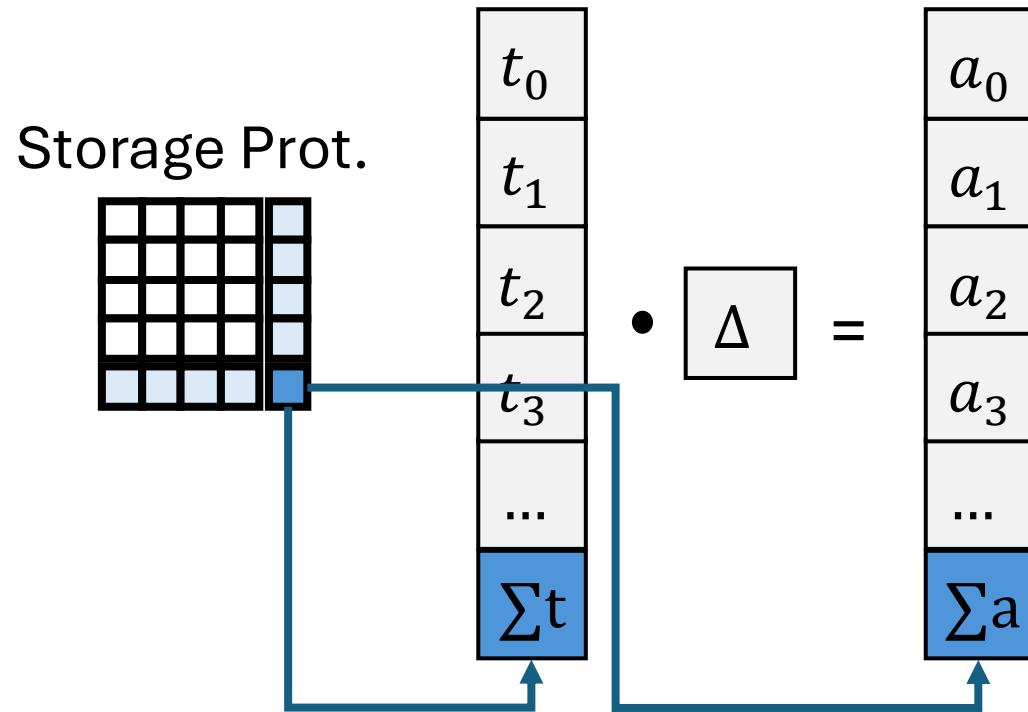
# FHE Computation Kernels



# Scaling Operation (e.g., BConv)

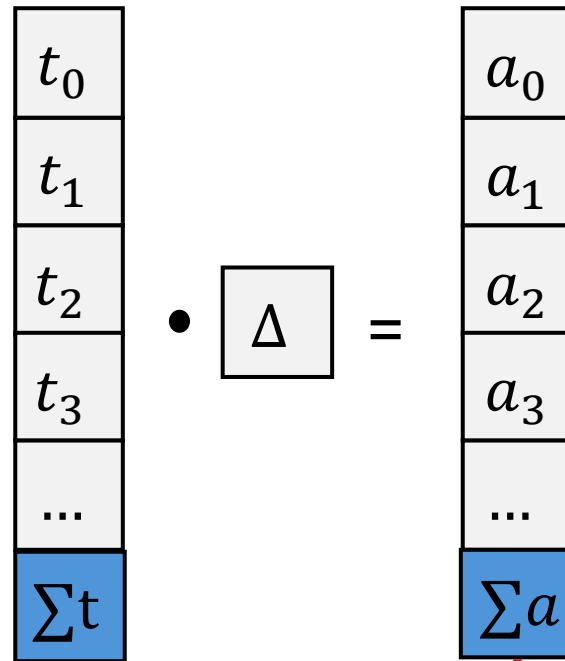


# Scaling Operation (e.g., BConv)



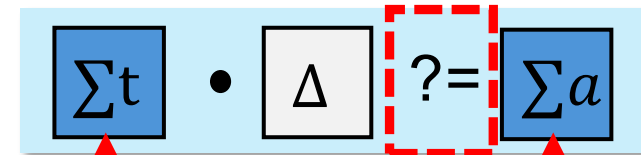
# Scaling Operation (e.g., BConv)

Storage Prot.

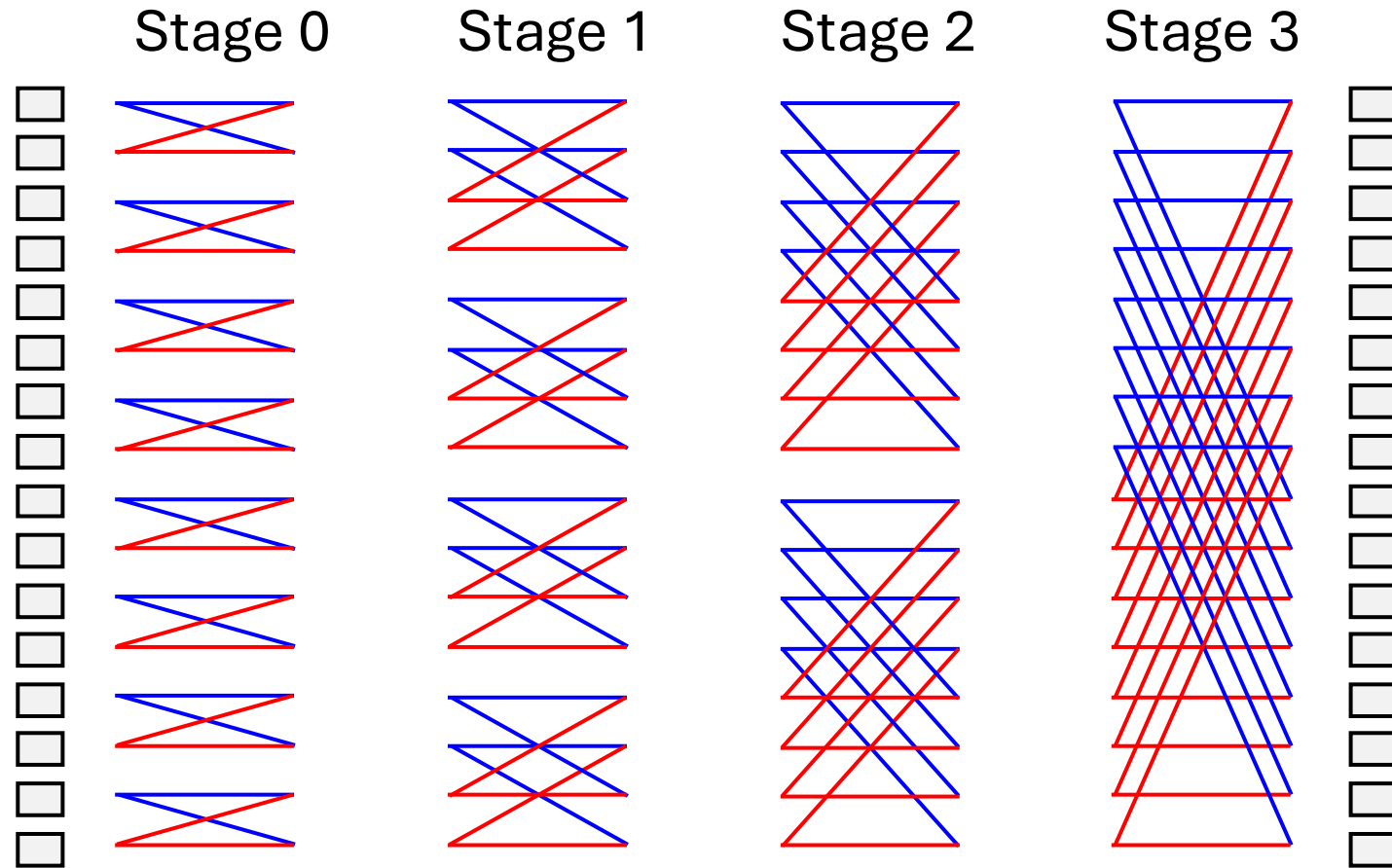


<0.01% Overhead

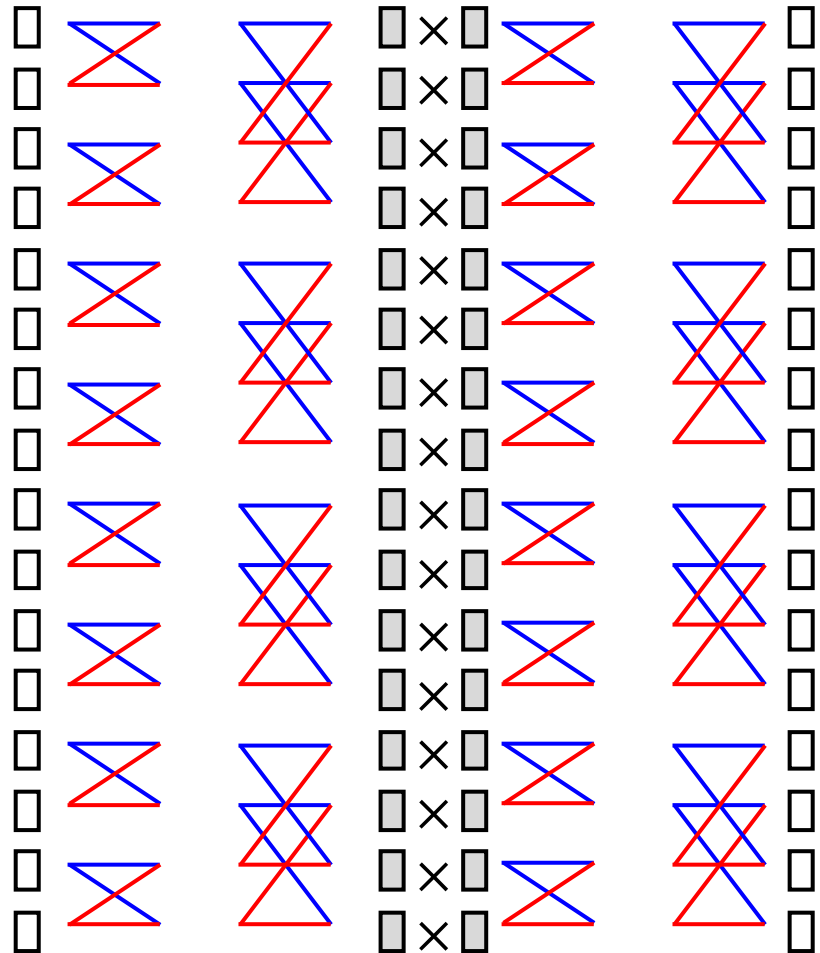
Error-Detection



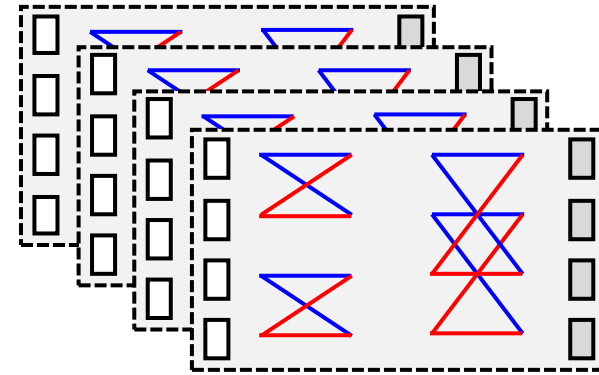
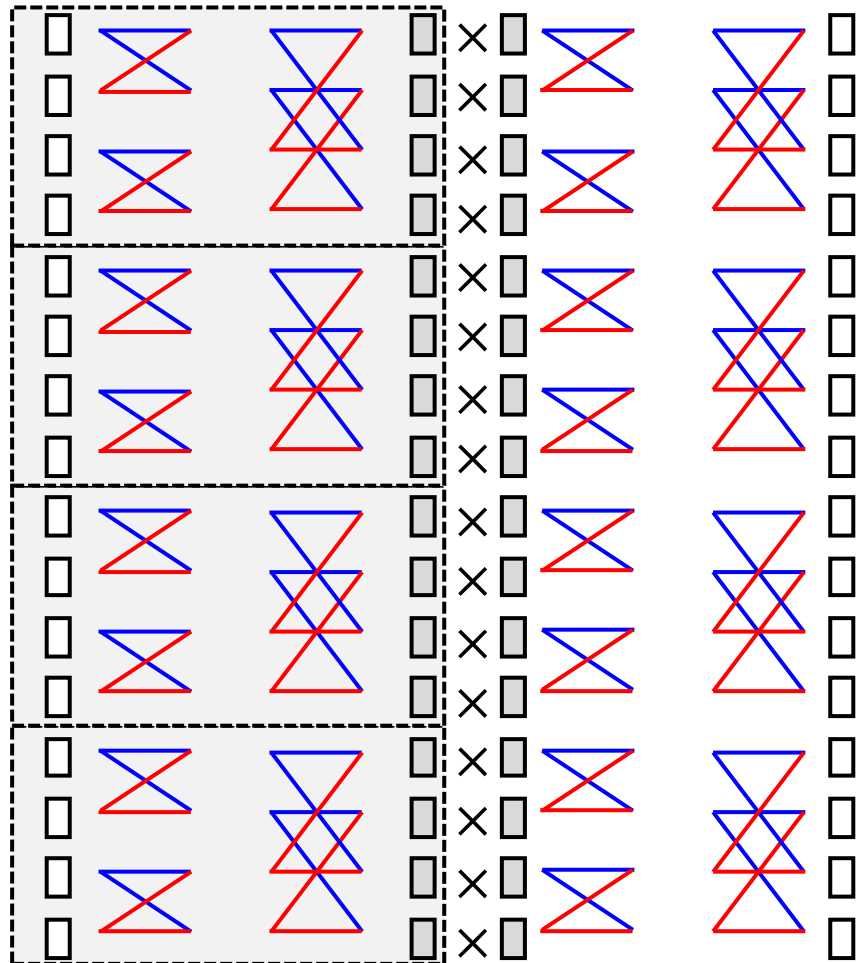
# NTT Protection



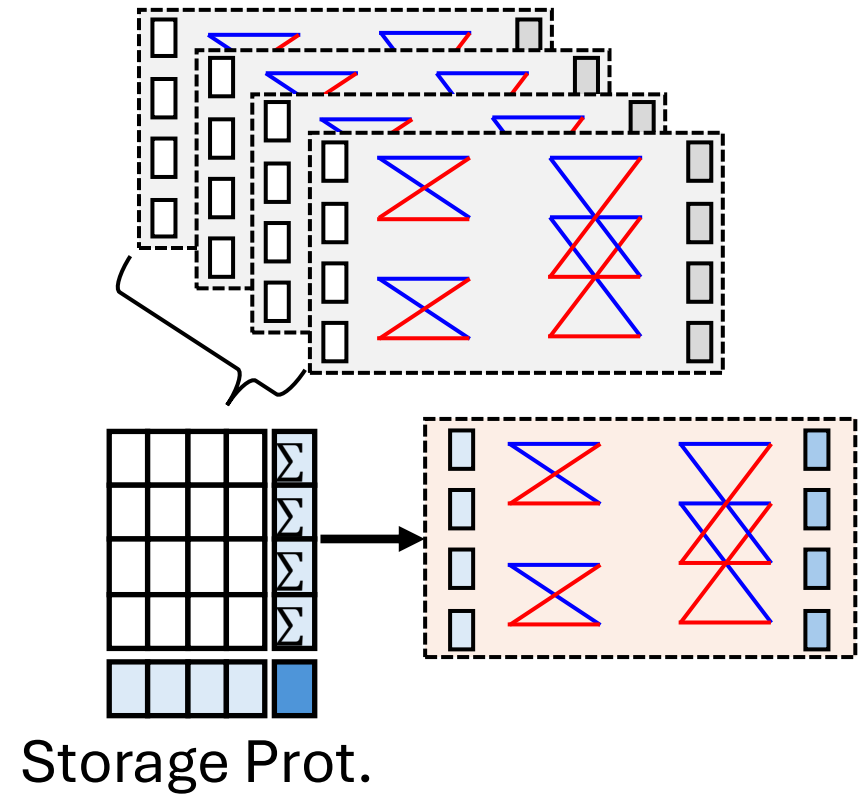
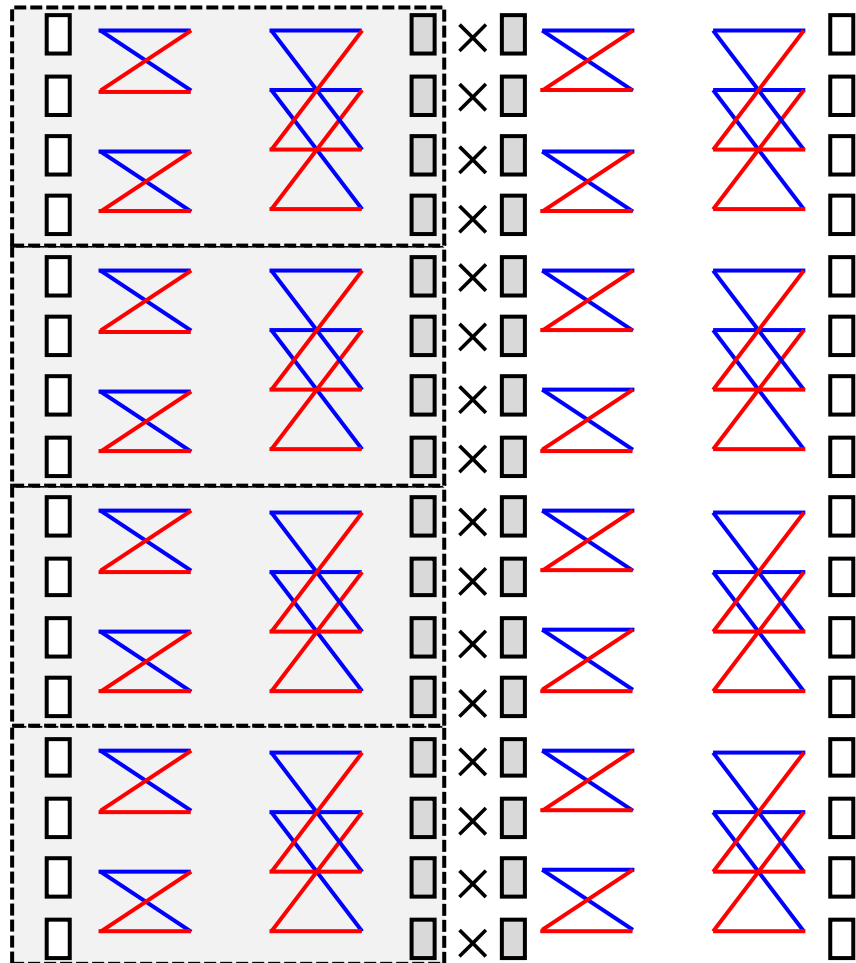
# NTT Protection



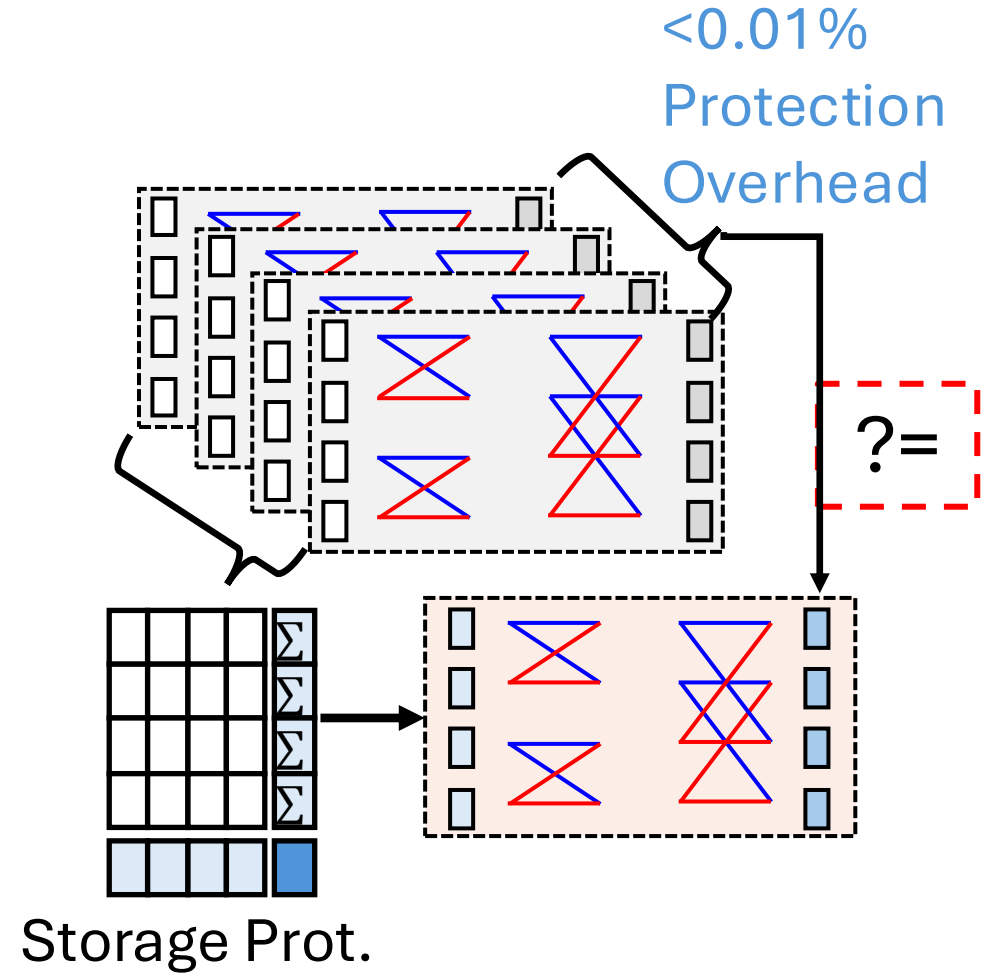
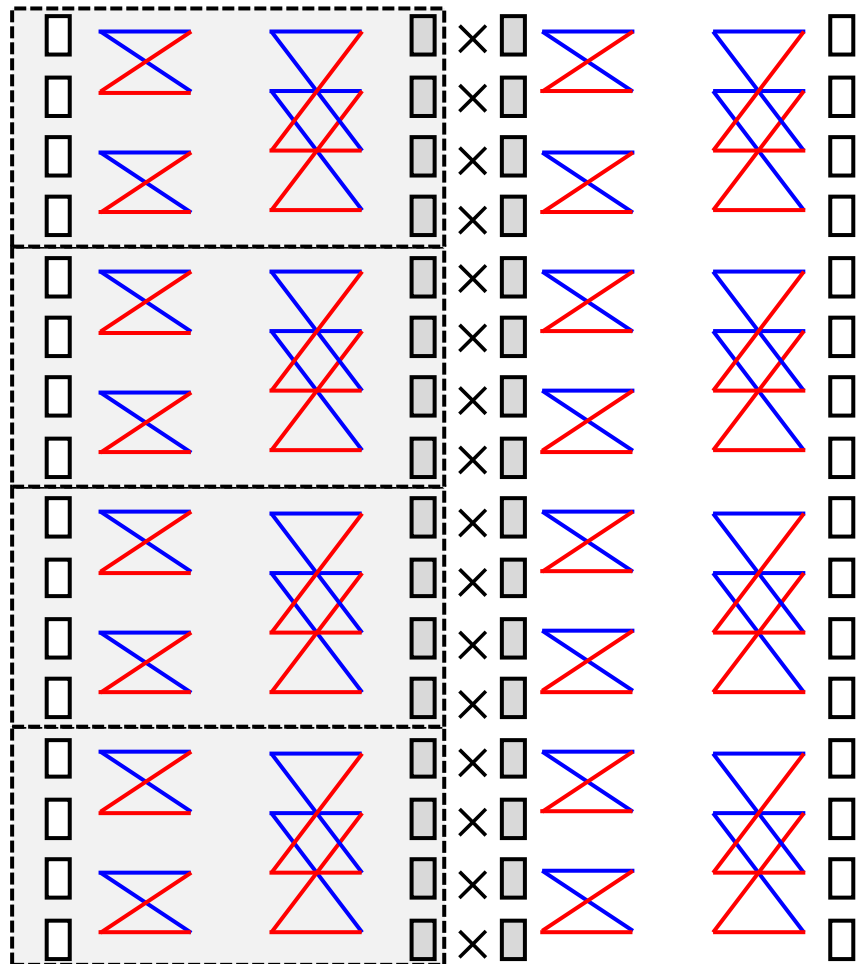
# NTT Protection



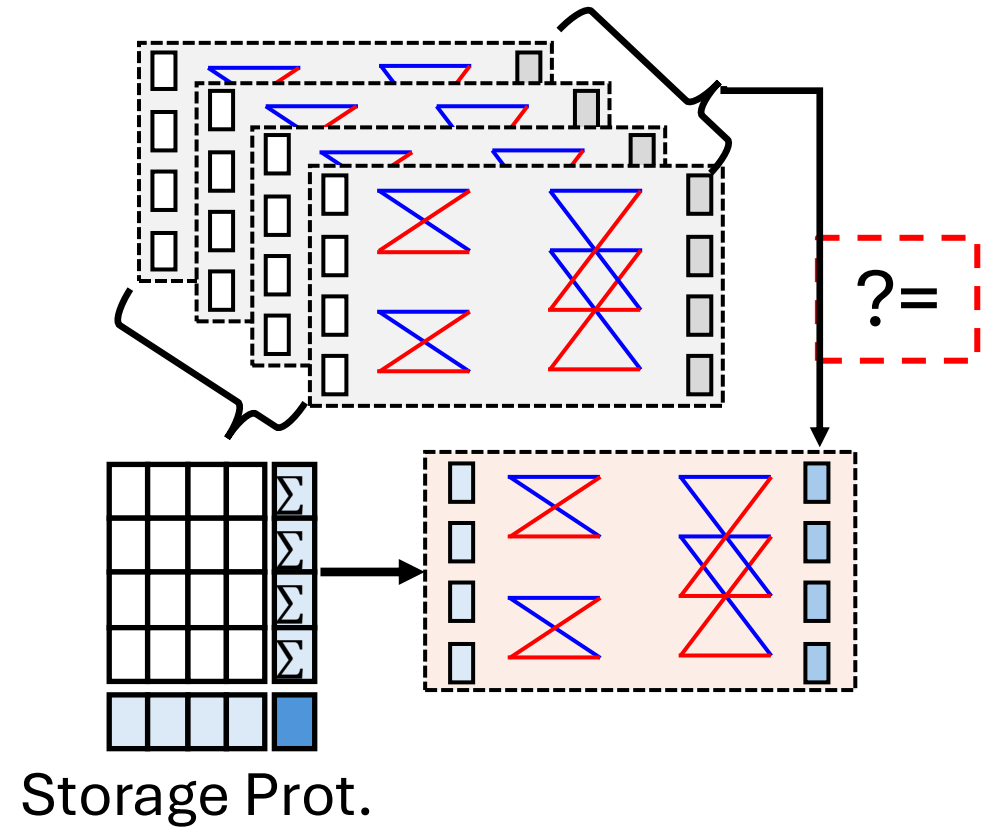
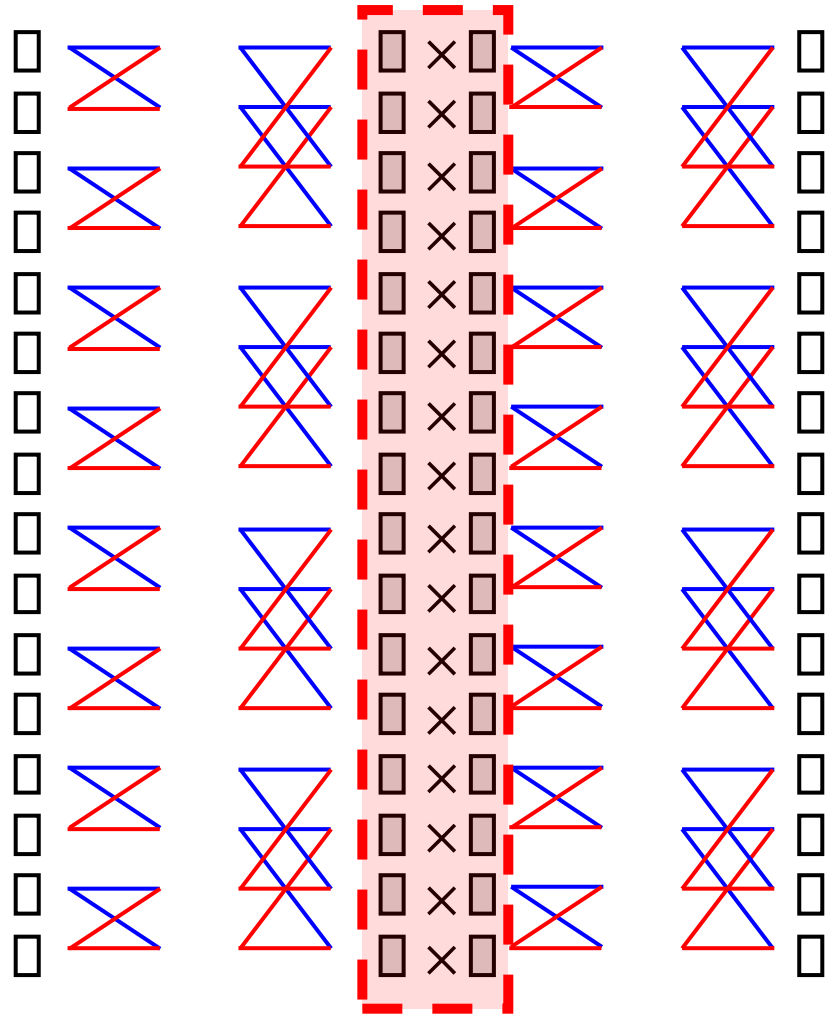
# NTT Protection



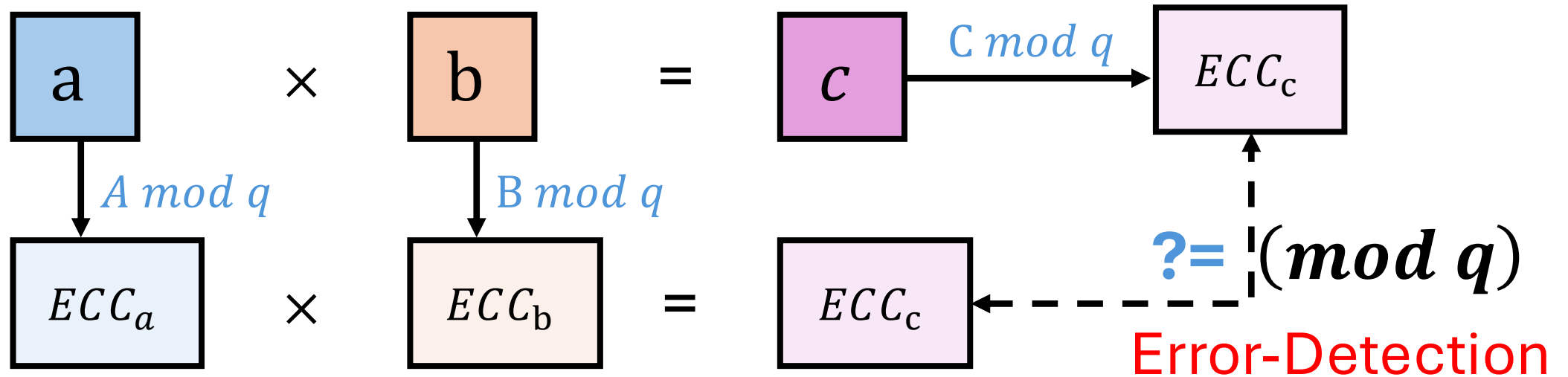
# NTT Protection



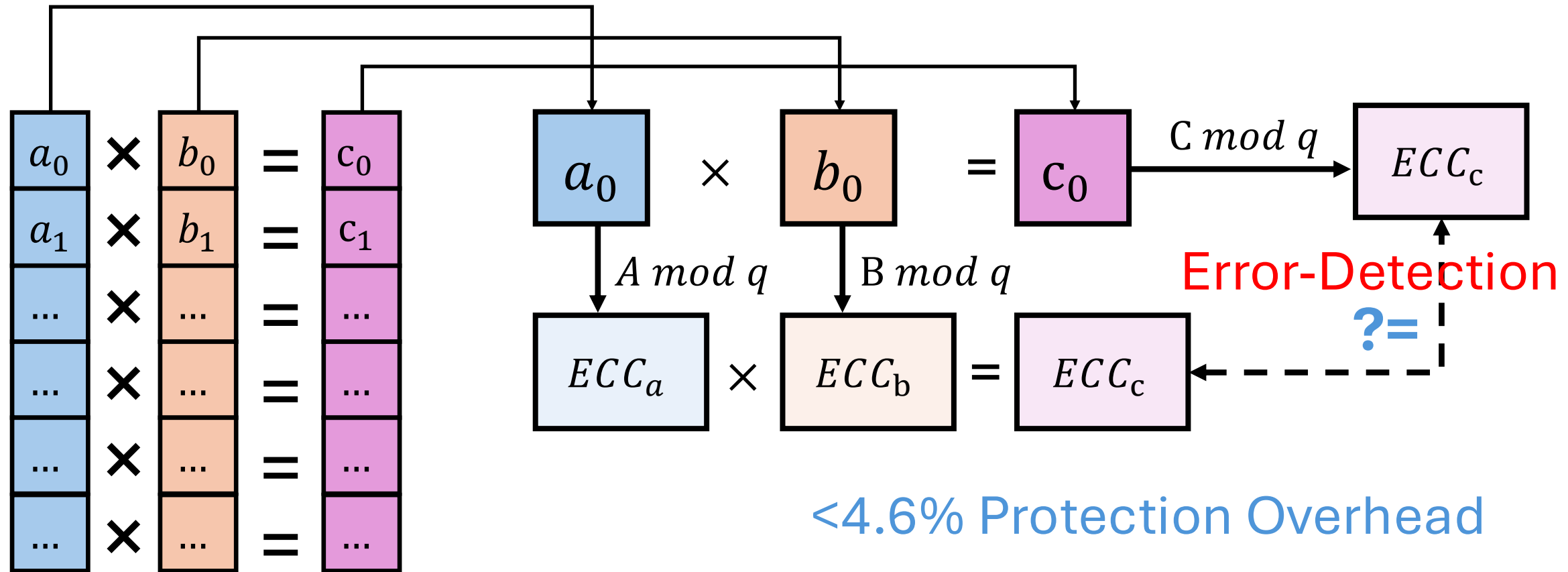
# NTT Protection



# Elementwise Protection (e.g. EvkMult)

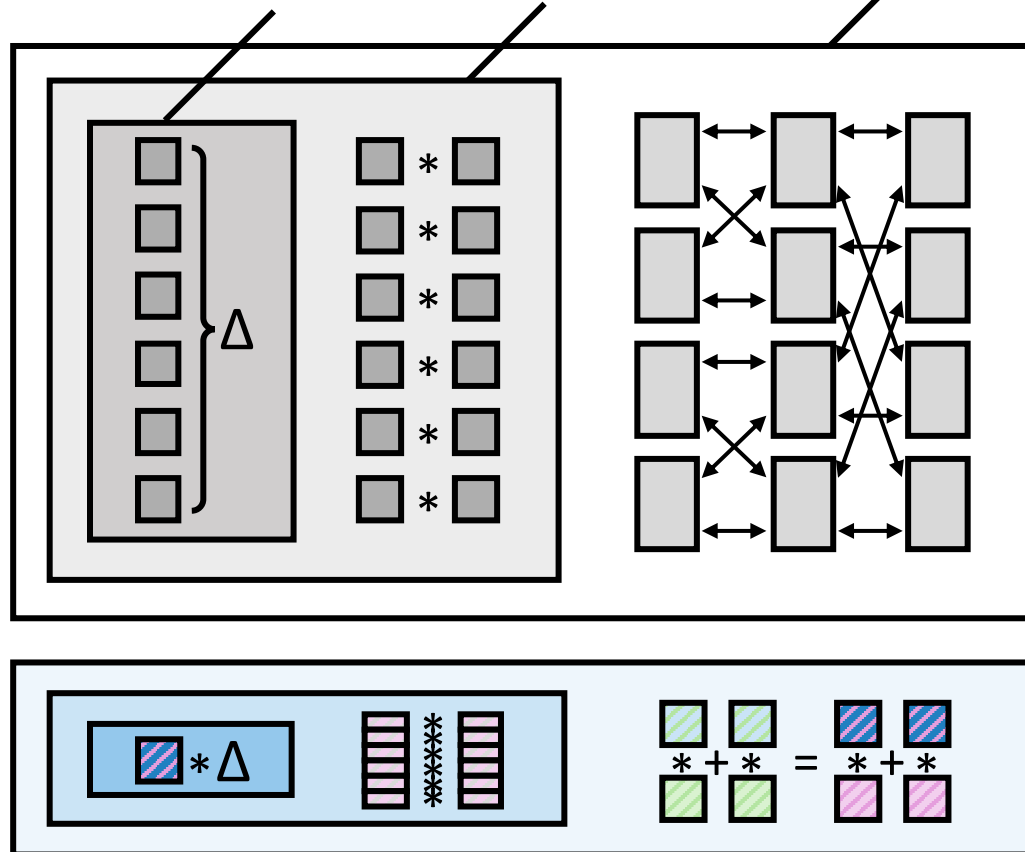


# Elementwise Protection (e.g. EvkMult)



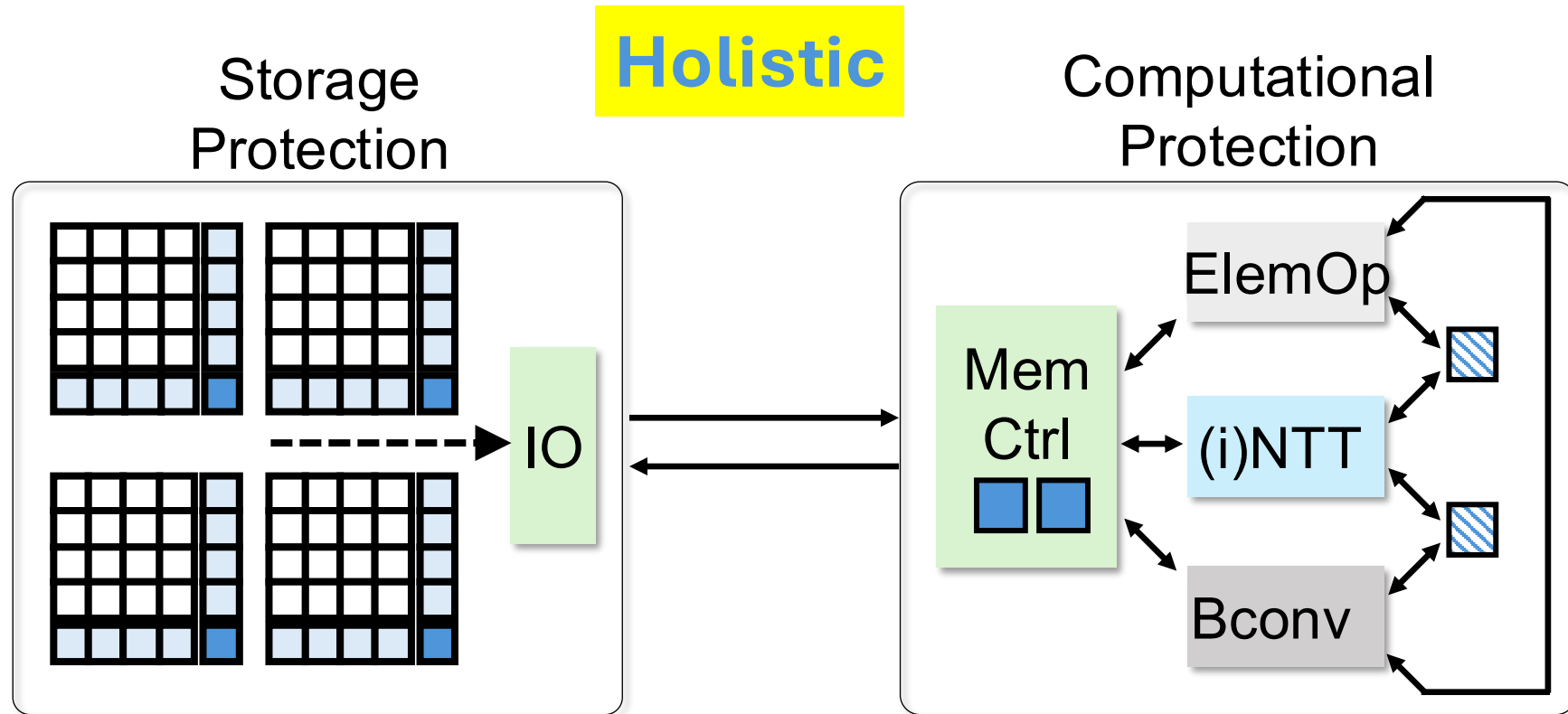
# Conclusion

Scale Op. Elementwise Op. NTT Op.






**Hierarchical**

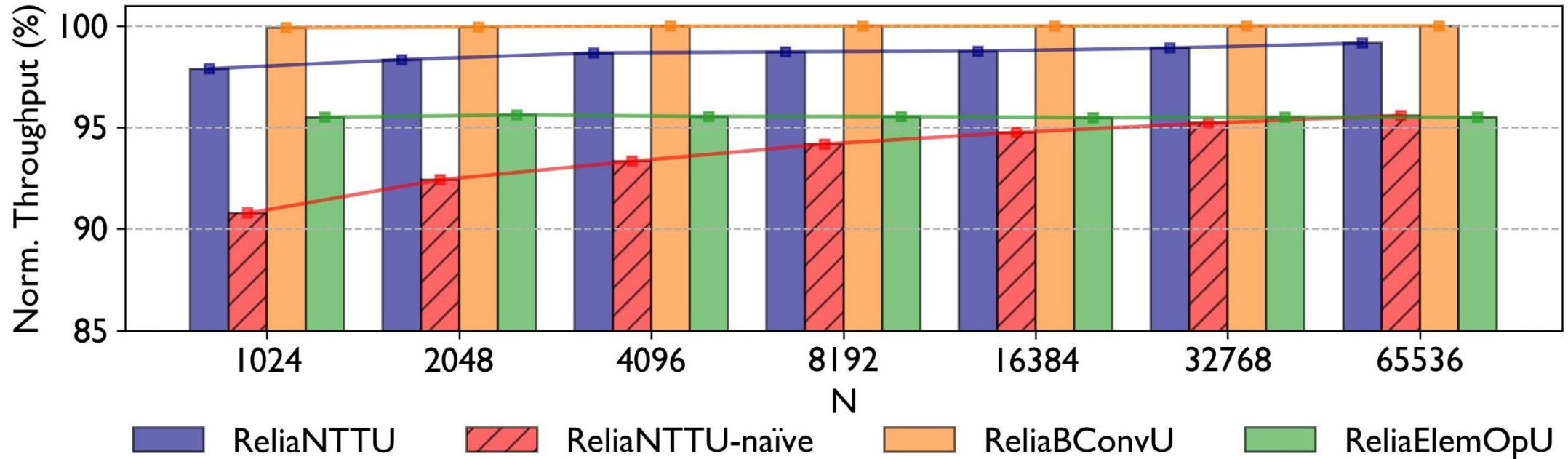
# Conclusion



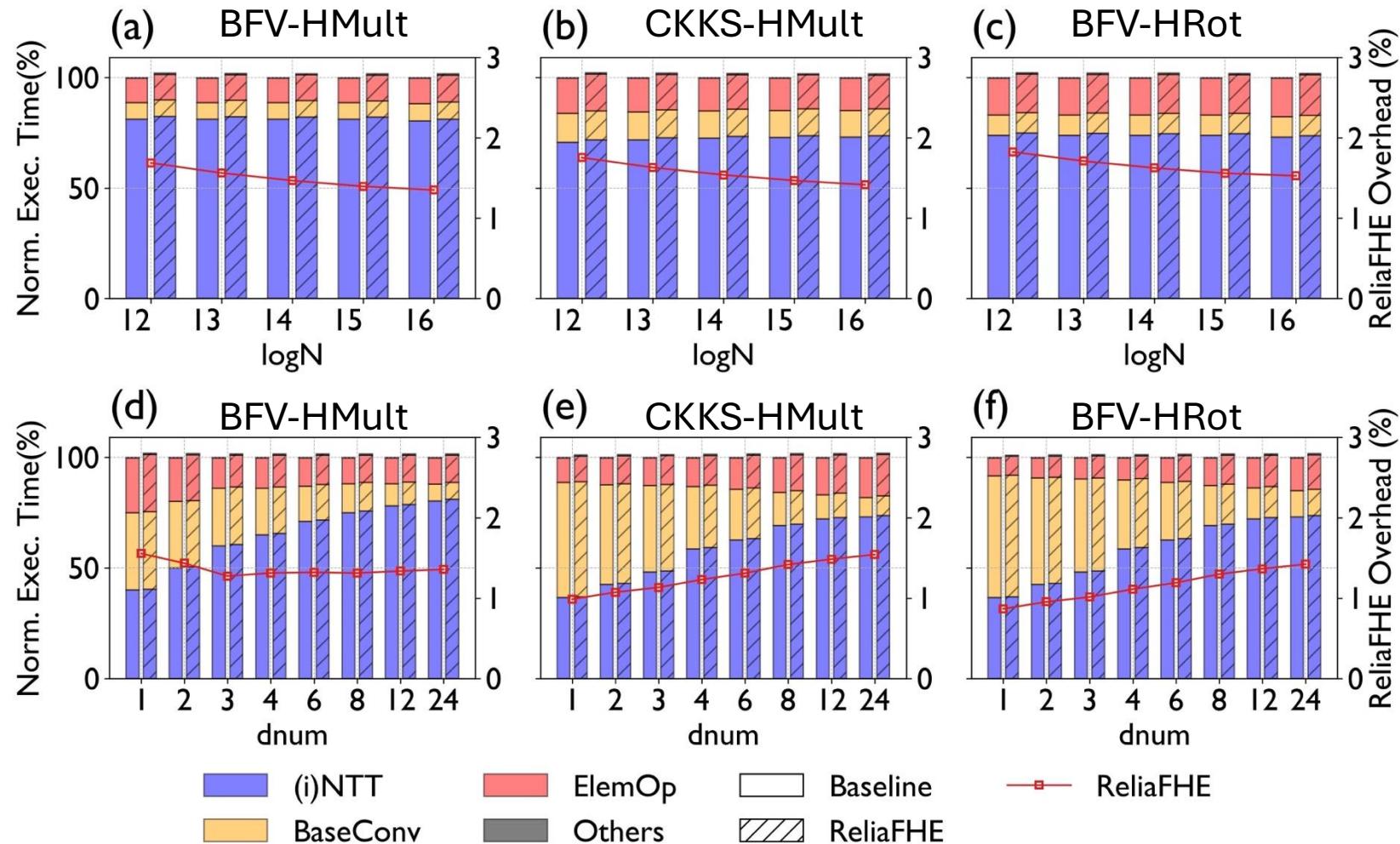
# Conclusion

 Holistic	 Hierarchical	 Lightweight
<b>Complete Protection</b>	<b>Multi-level Defense</b>	<b>Minimal Impact</b>
Integrated security for both <b>Memory</b> and <b>Computation</b> .	Layered granularity: <b>Scale operation</b> → <b>Elementwise Operation</b> → <b>NTT operation</b> .	Exceptional efficiency with only <b>1%~2%</b> overhead.

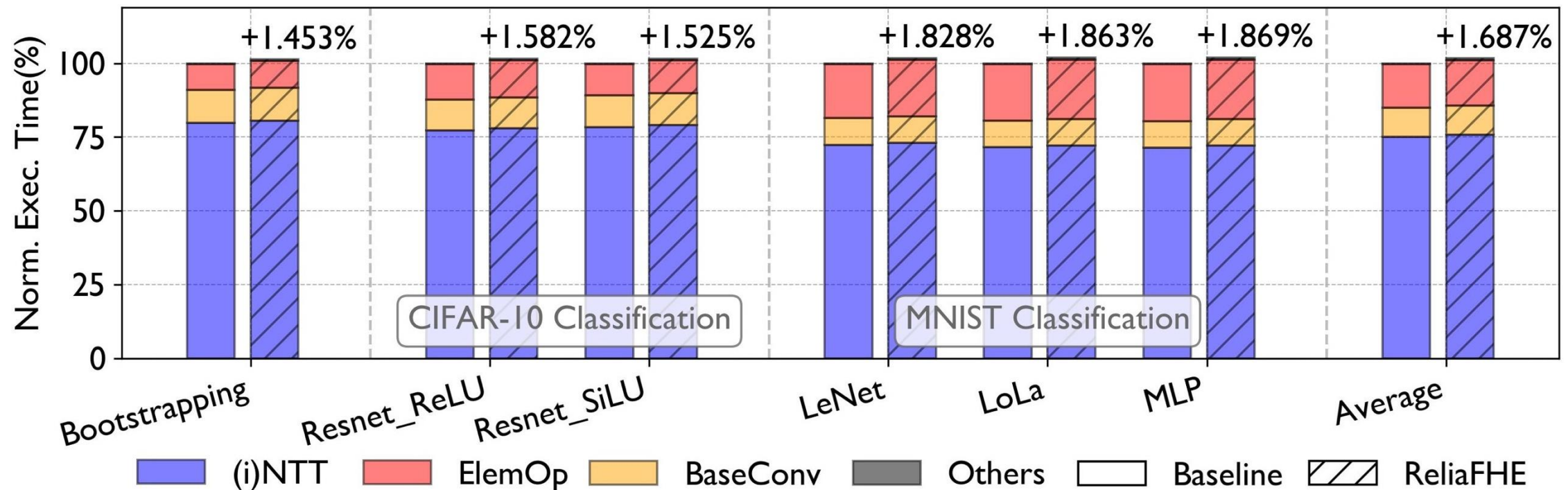
# Performance-Operation Kernels



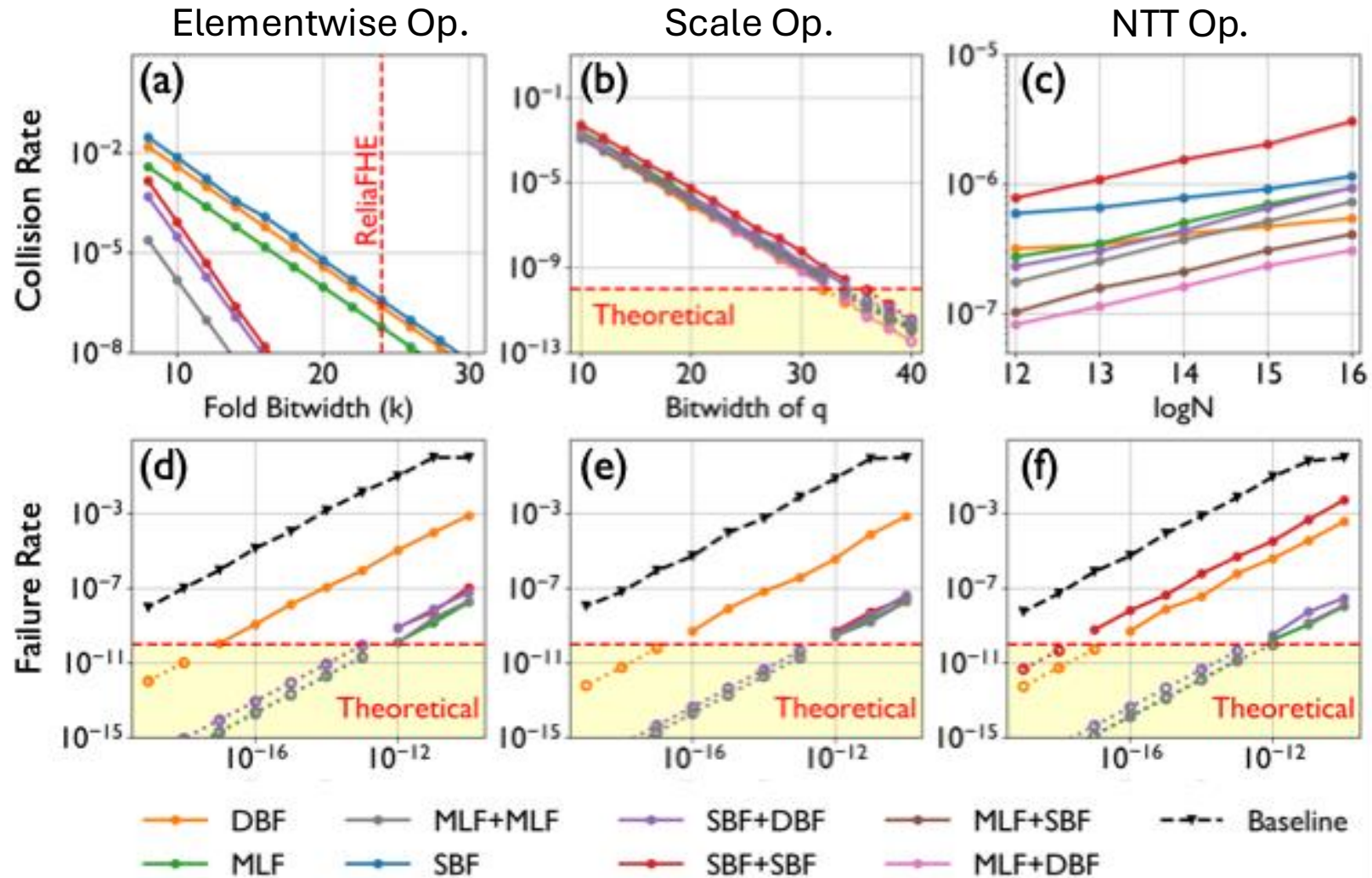
# Performance-Primitive Operations



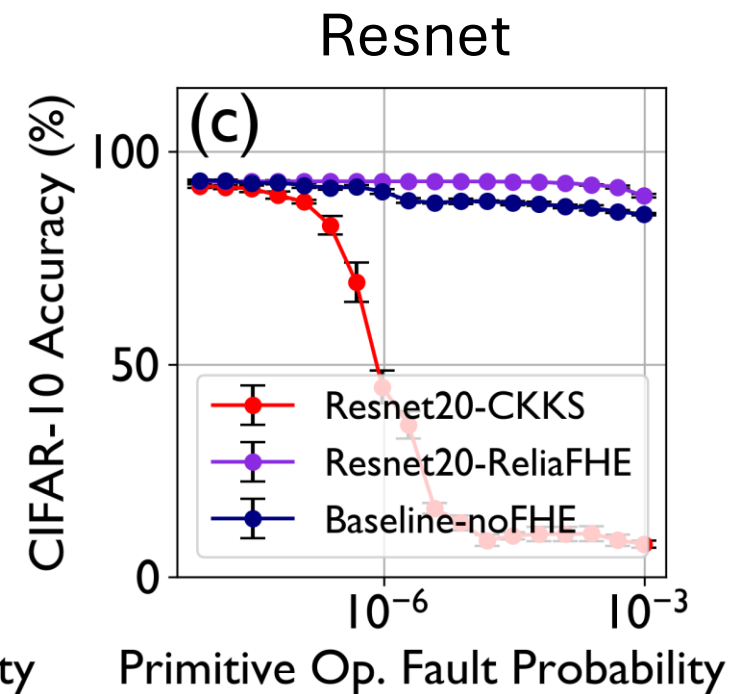
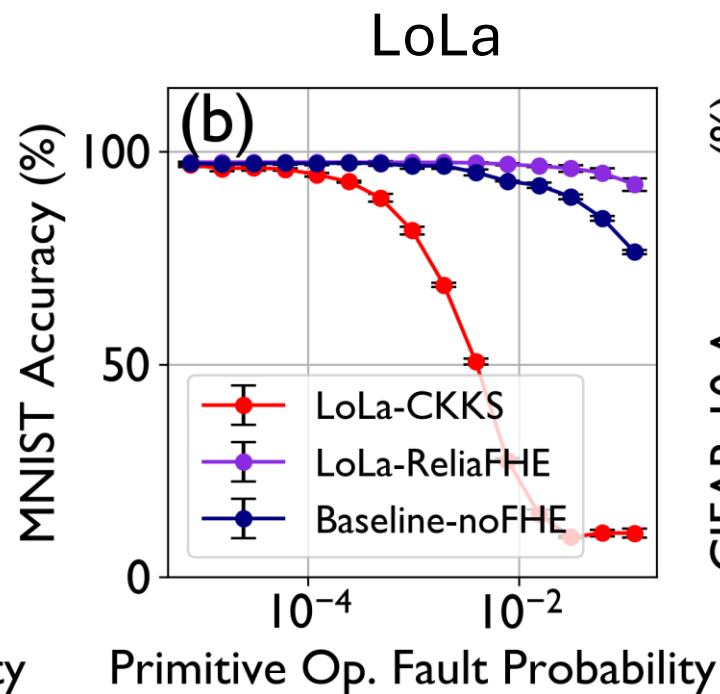
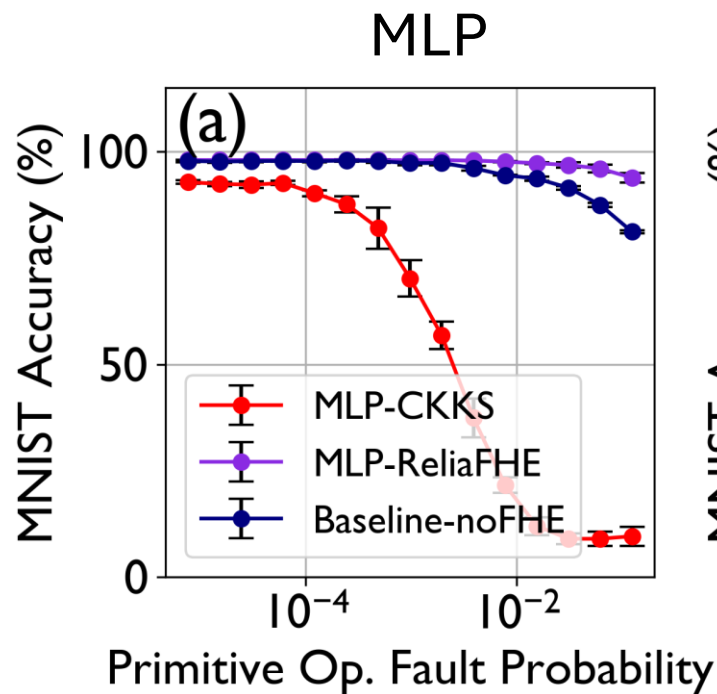
# Performance-Application



# Reliability



# Error-tolerance Ability



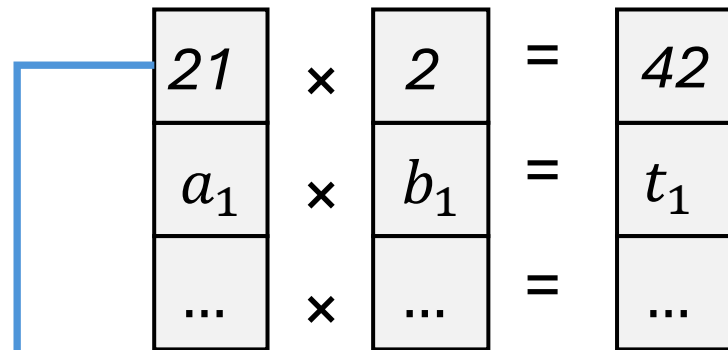
Thank you for listening  
Q&A

Fan Li

# Elementwise Operation (e.g. EvkMult)

$$\begin{array}{c} a_0 \\ a_1 \\ \dots \\ \dots \\ \dots \\ \dots \end{array} \times \begin{array}{c} b_0 \\ b_1 \\ \dots \\ \dots \\ \dots \\ \dots \end{array} = \begin{array}{c} t_0 \\ t_1 \\ \dots \\ \dots \\ \dots \\ \dots \end{array}$$

# Elementwise Operation (e.g. EvkMult)



Error-Detection

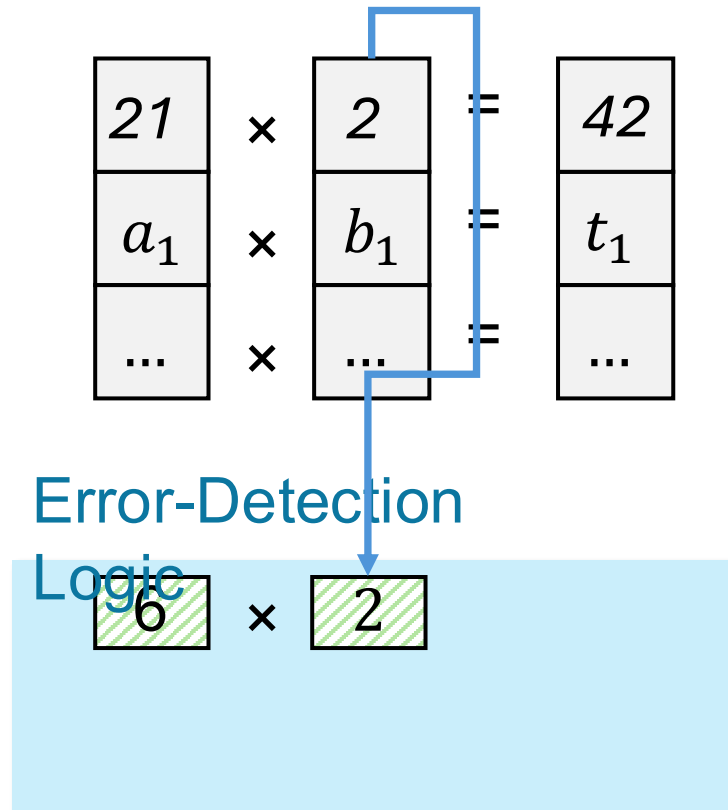
Logic

6

$$21 = b0000\ 0000\ 0001\ 0101$$

$$\begin{aligned} 21_{\Sigma} &= b(0000 + 0000 + 0001 + 0101) \\ &= b0110 \\ &= 6 \end{aligned}$$

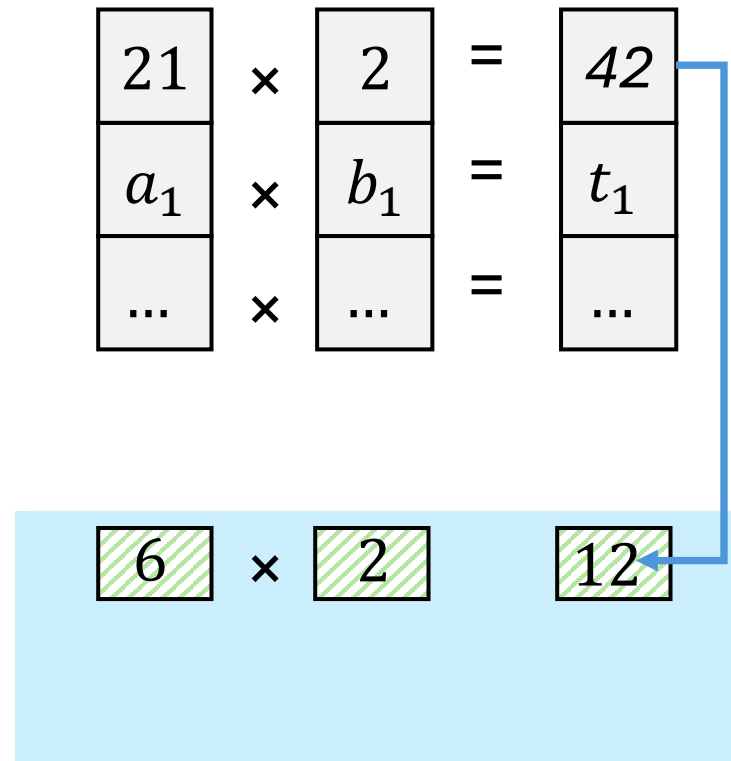
# Elementwise Operation (e.g. EvkMult)



$$2 = b0000\ 0000\ 0000\ 0010$$

$$\begin{aligned} 2_{\Sigma} &= b(0000 + 0000 + 0000 + 0010) \\ &= b0010 \\ &= 2 \end{aligned}$$

# Elementwise Operation (e.g. EvkMult)



$$42 = b0000\ 0000\ 0010\ 1010$$

$$\begin{aligned}
 42_{\Sigma} &= b(0000 + 0000 + 0010 + 1010) \\
 &= b1100 \\
 &= 12
 \end{aligned}$$

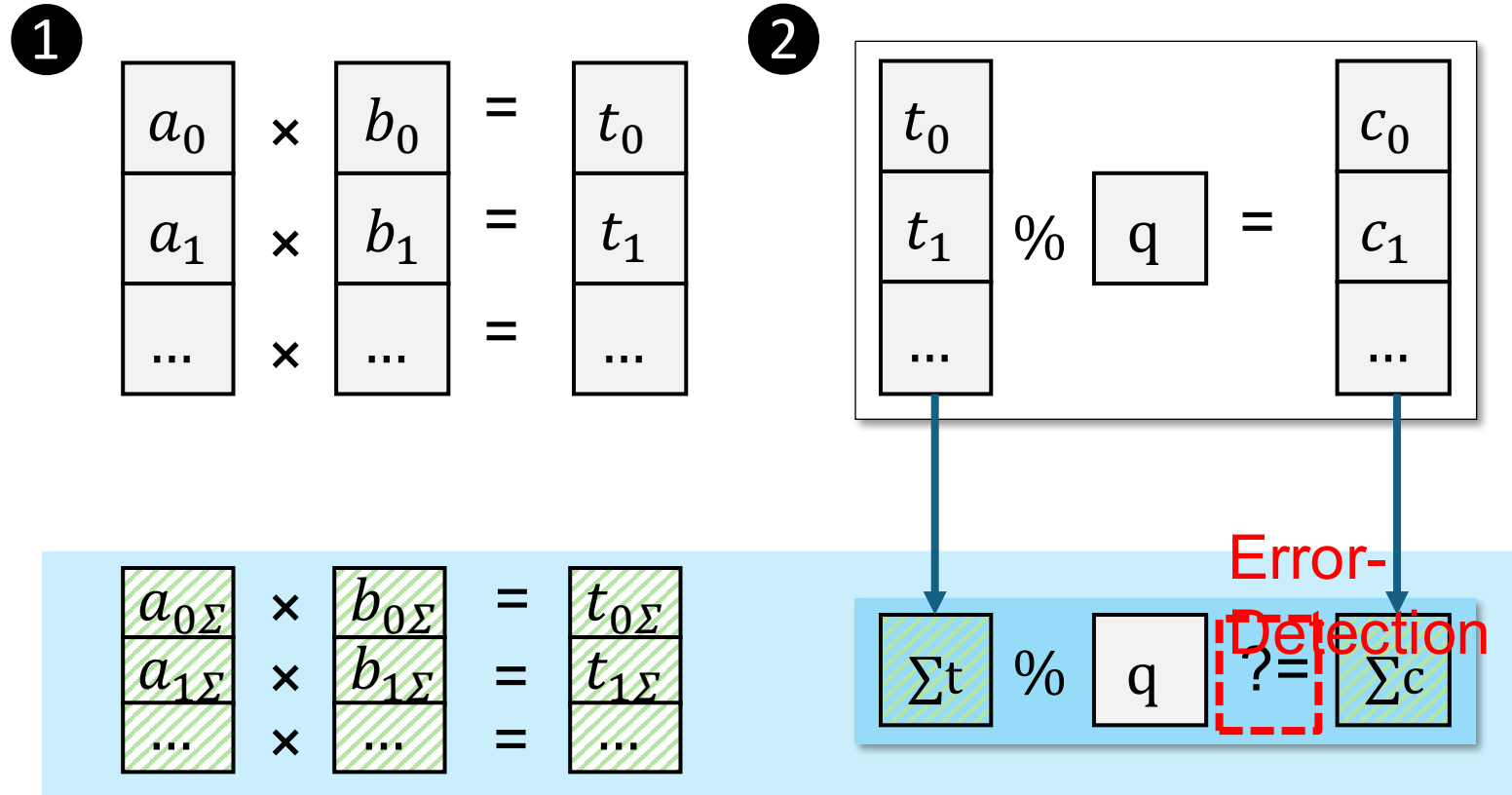
# Elementwise Operation (e.g. EvkMult)

$$\begin{array}{|c|} \hline 21 \\ \hline a_1 \\ \hline \dots \\ \hline \end{array} \times \begin{array}{|c|} \hline 2 \\ \hline b_1 \\ \hline \dots \\ \hline \end{array} = \begin{array}{|c|} \hline 42 \\ \hline t_1 \\ \hline \dots \\ \hline \end{array}$$

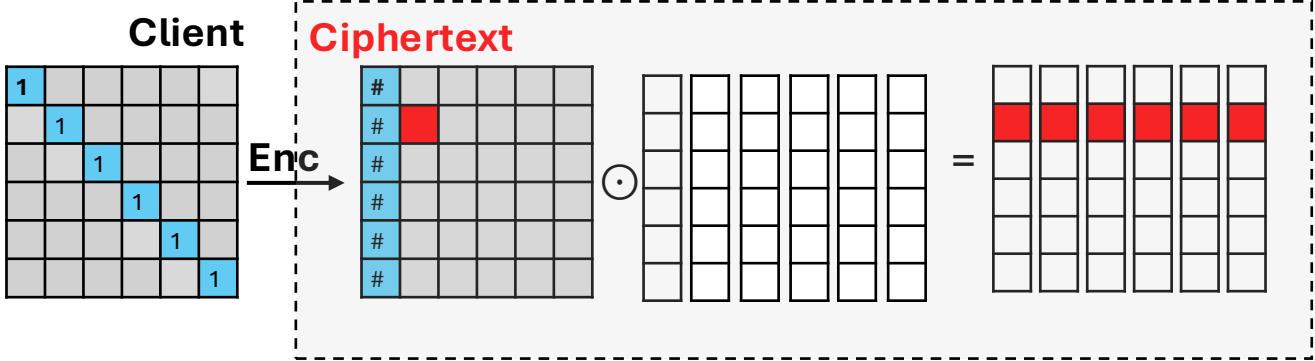
Error-Detection

$$\begin{array}{|c|} \hline 6 \\ \hline a_{1\Sigma} \\ \hline \dots \\ \hline \end{array} \times \begin{array}{|c|} \hline 2 \\ \hline b_{1\Sigma} \\ \hline \dots \\ \hline \end{array} = \begin{array}{|c|} \hline 12 \\ \hline t_{1\Sigma} \\ \hline \dots \\ \hline \end{array}$$

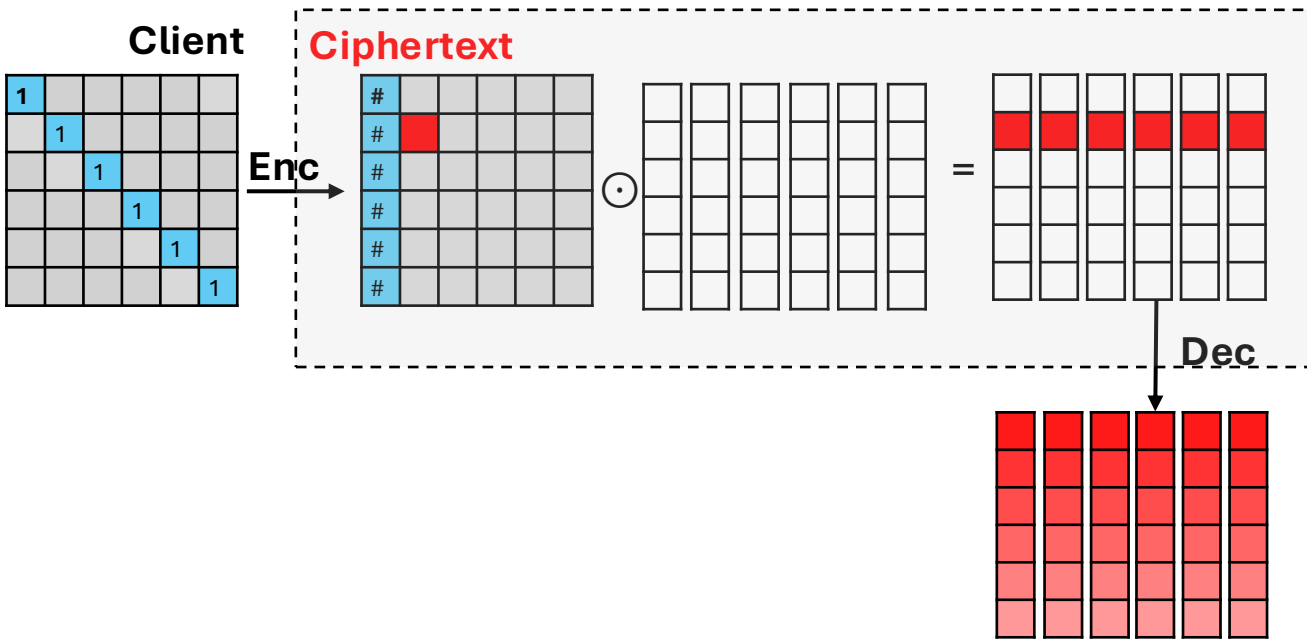
# Elementwise Operation (e.g. EvkMult)



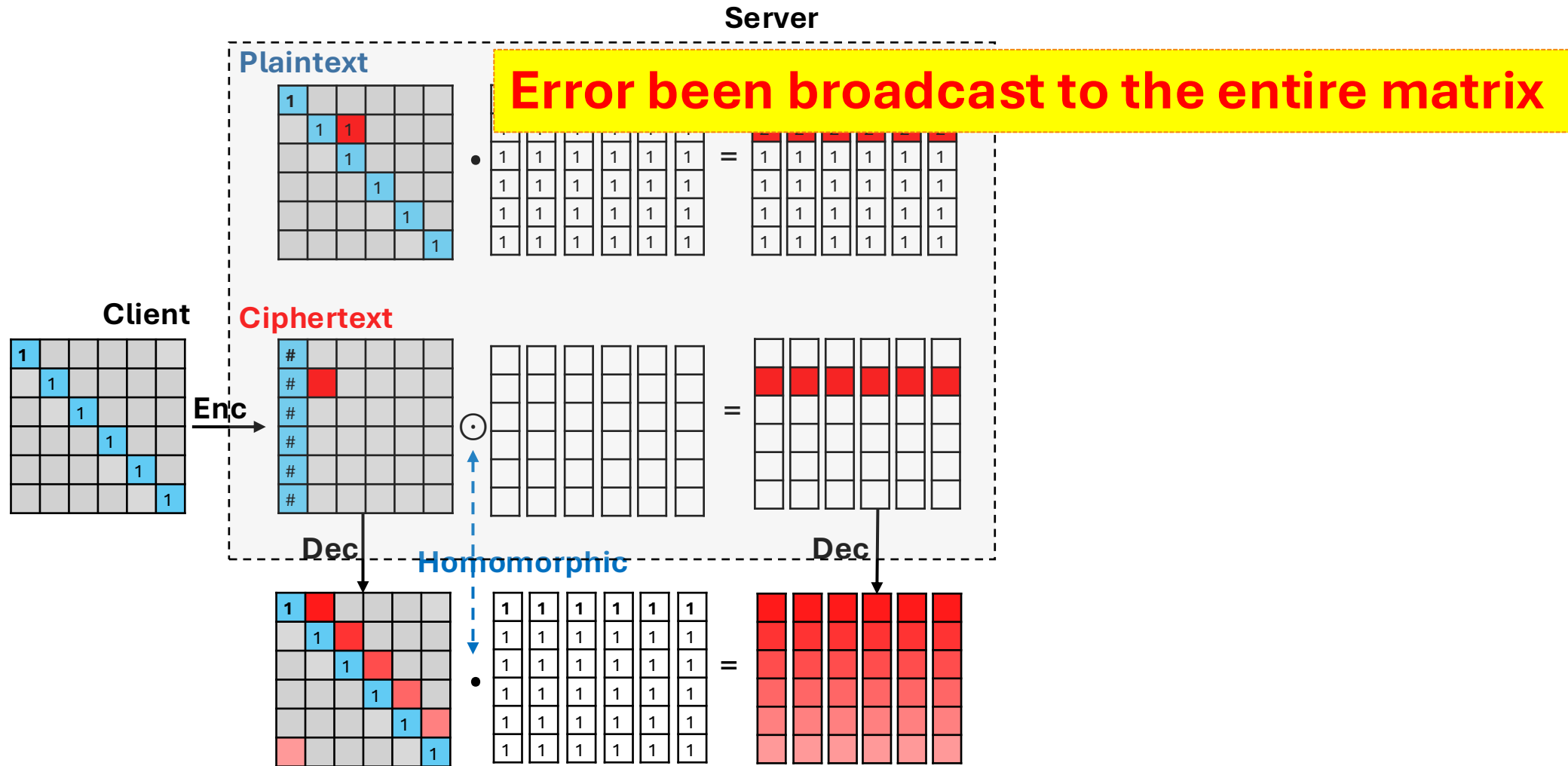
# Algorithm-SIMD Packing



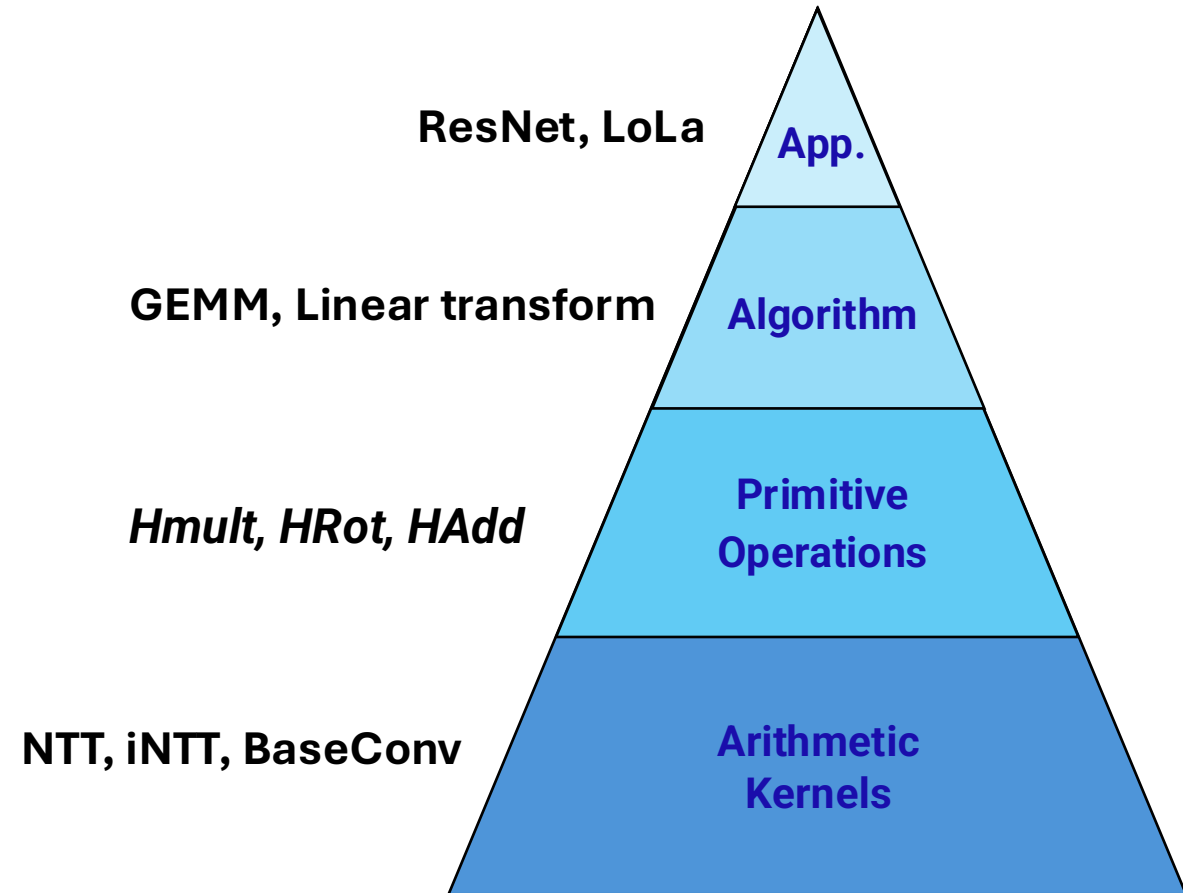
# Algorithm-SIMD Packing



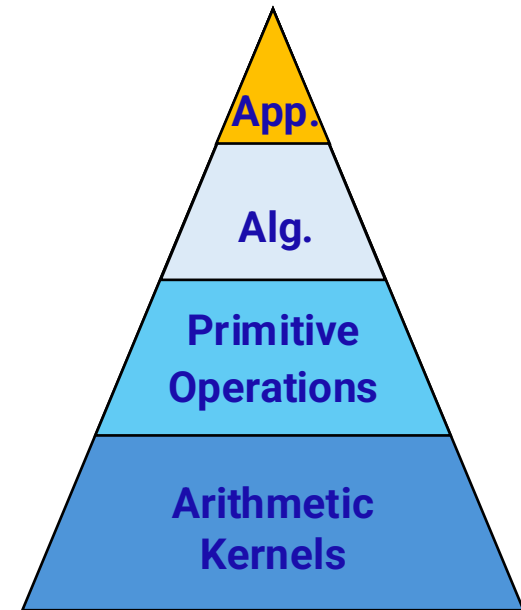
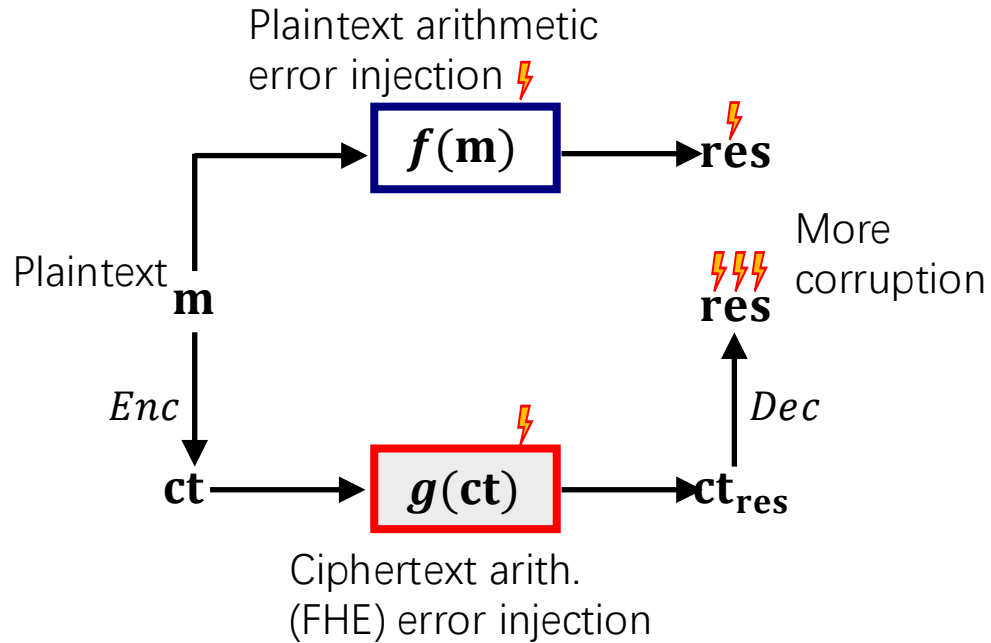
# Algorithm-SIMD Packing



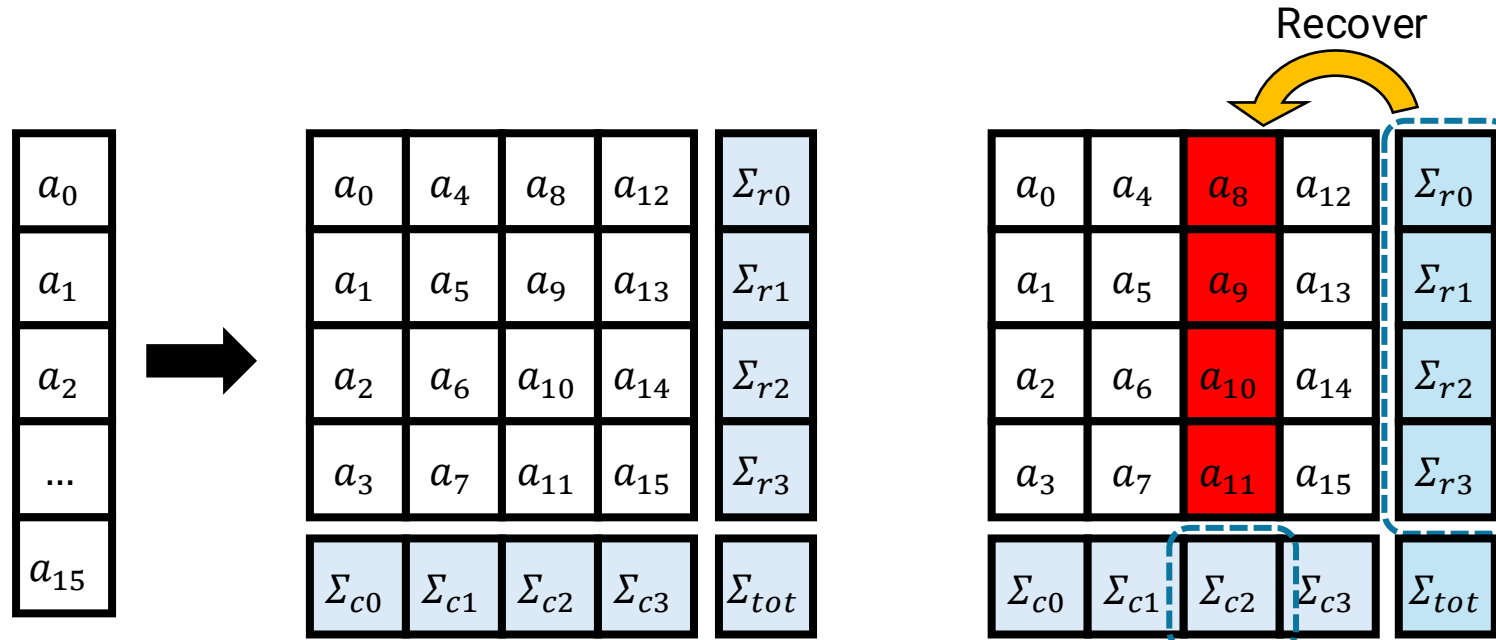
# Operations



# Application



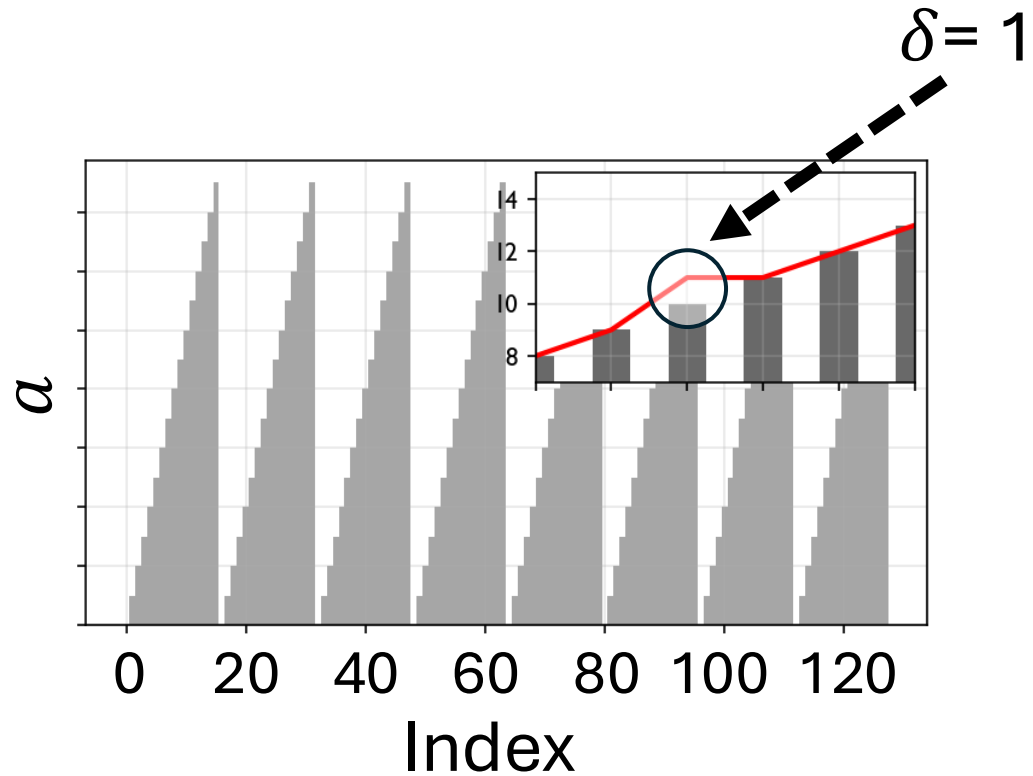
# Storage Protection(e.g. coefficient)



$\Sigma_{tot}$ : total Checksum  
 $\Sigma_r, \Sigma_c$ : row and col-wise Checksum

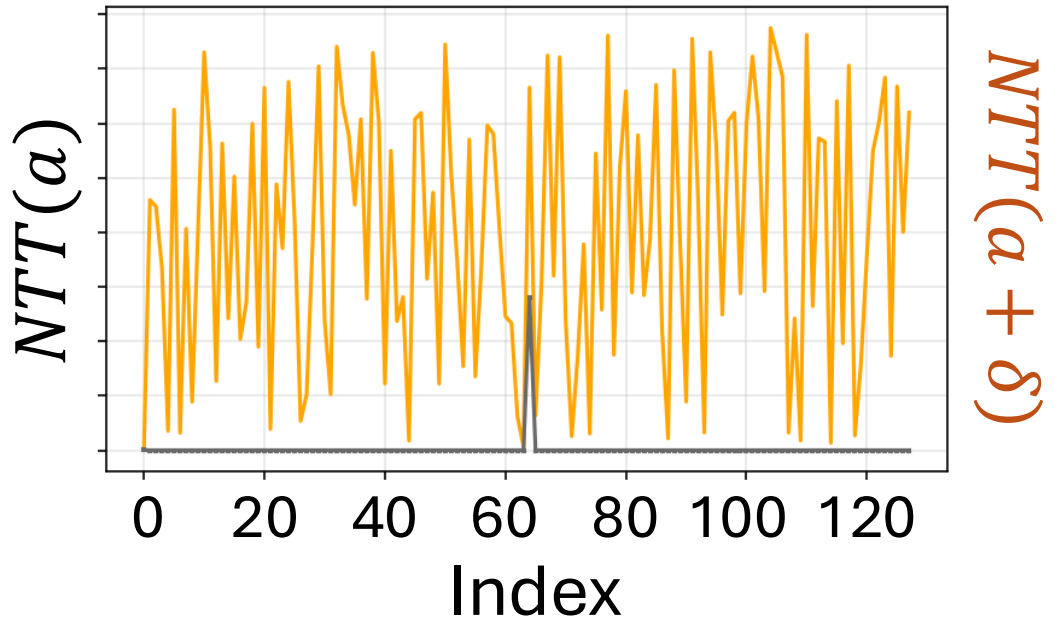
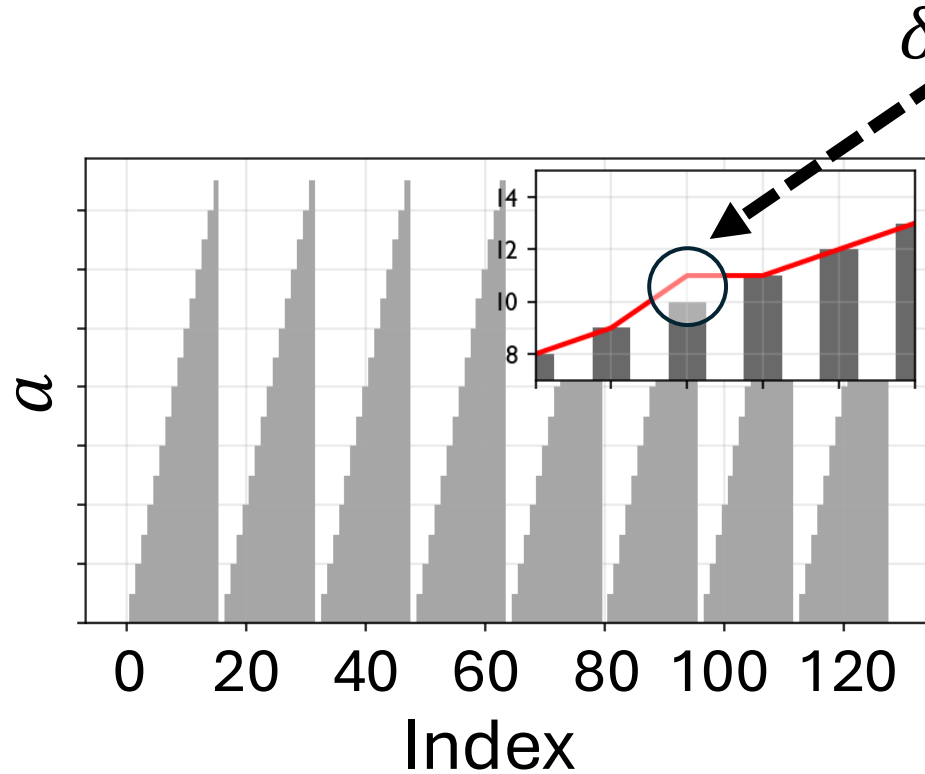
Detect  
 $\Sigma_{c2} \neq a_8 + a_9 + a_{10} + a_{11}$

# Arithmetic Kernel-(i)NTT



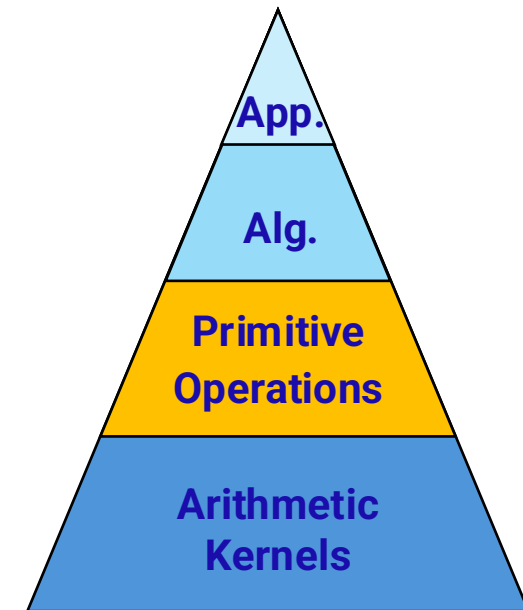
# Arithmetic Kernel-(i)NTT

**Error been broadcasted and amplified**

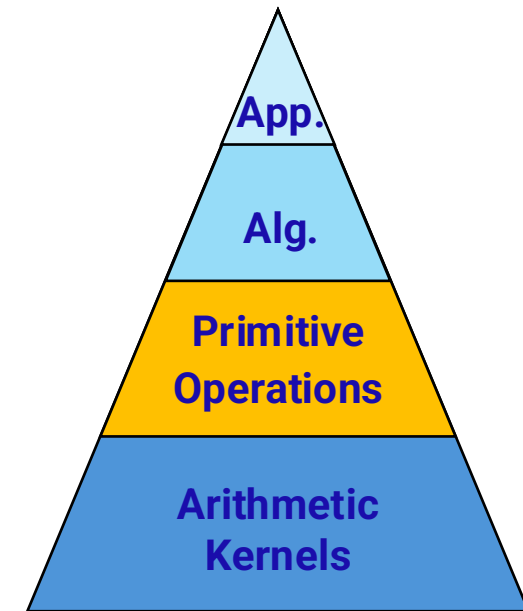
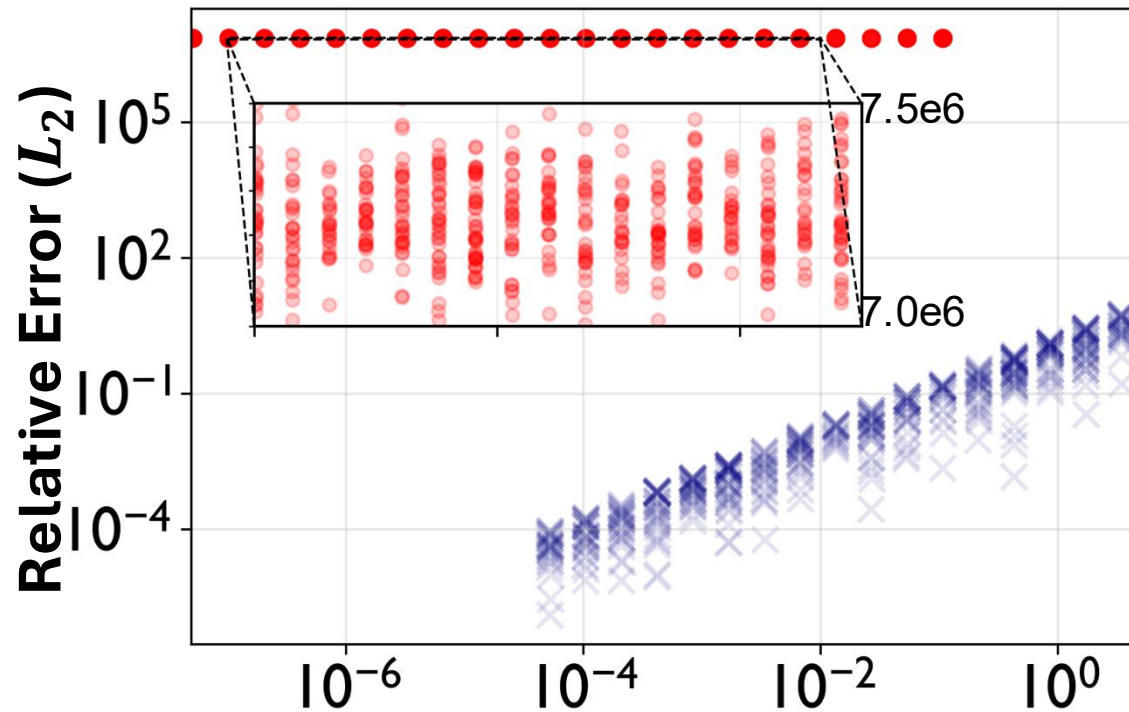


# Primitive Operations-Hadamard Product

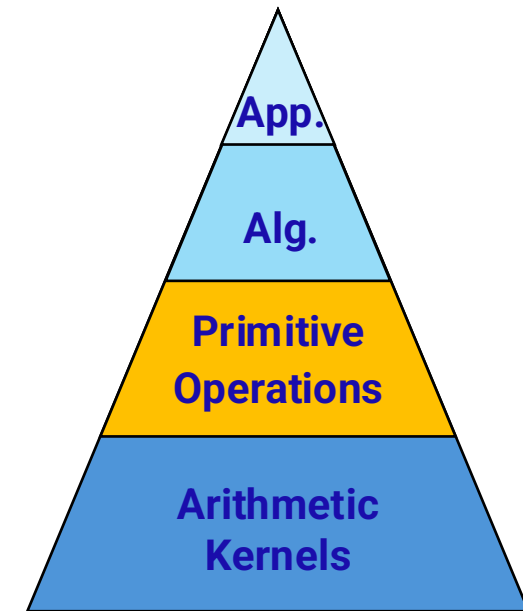
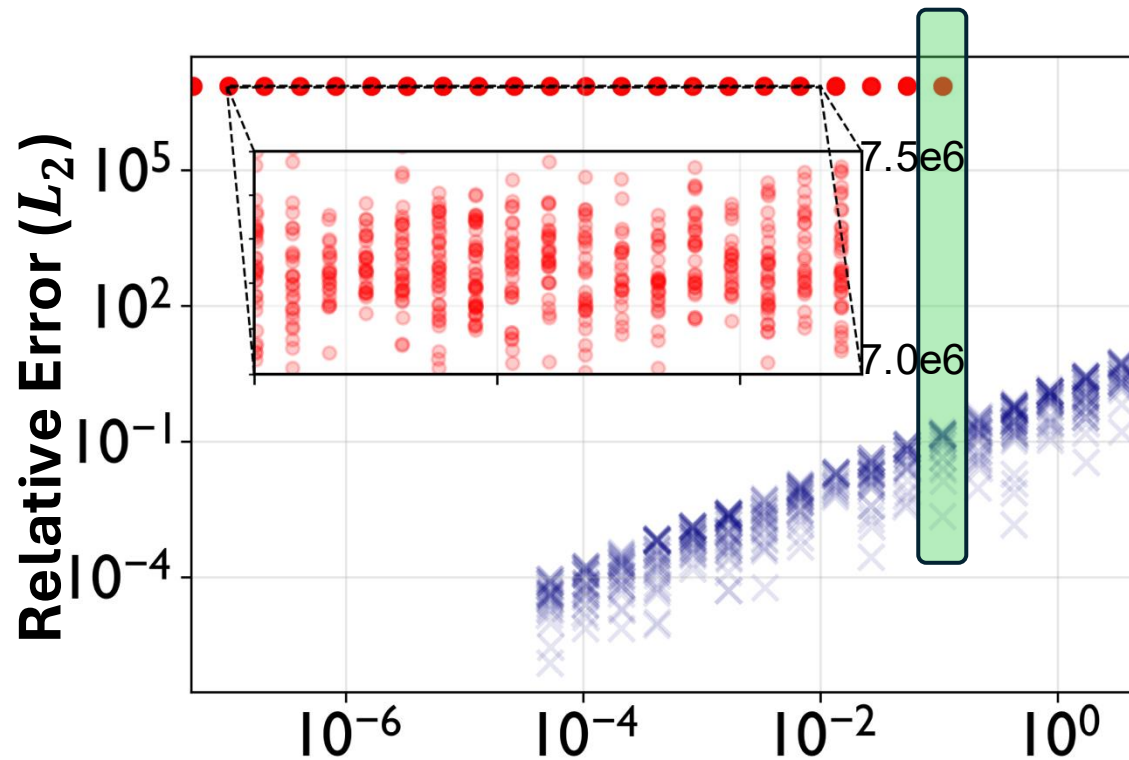
P1  $\odot$  P2 **Versus** C1  $\odot$  C2



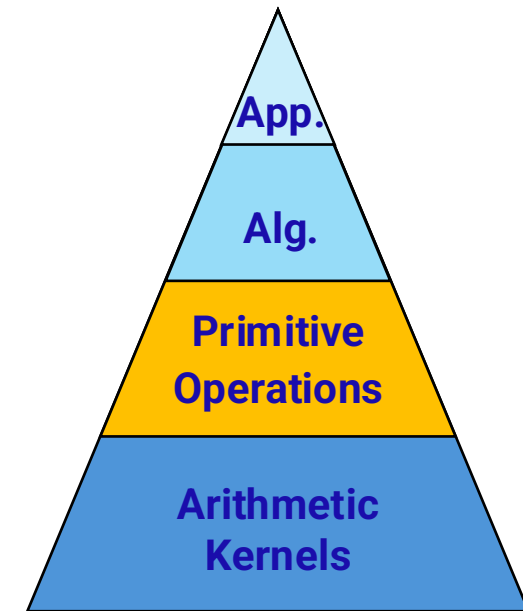
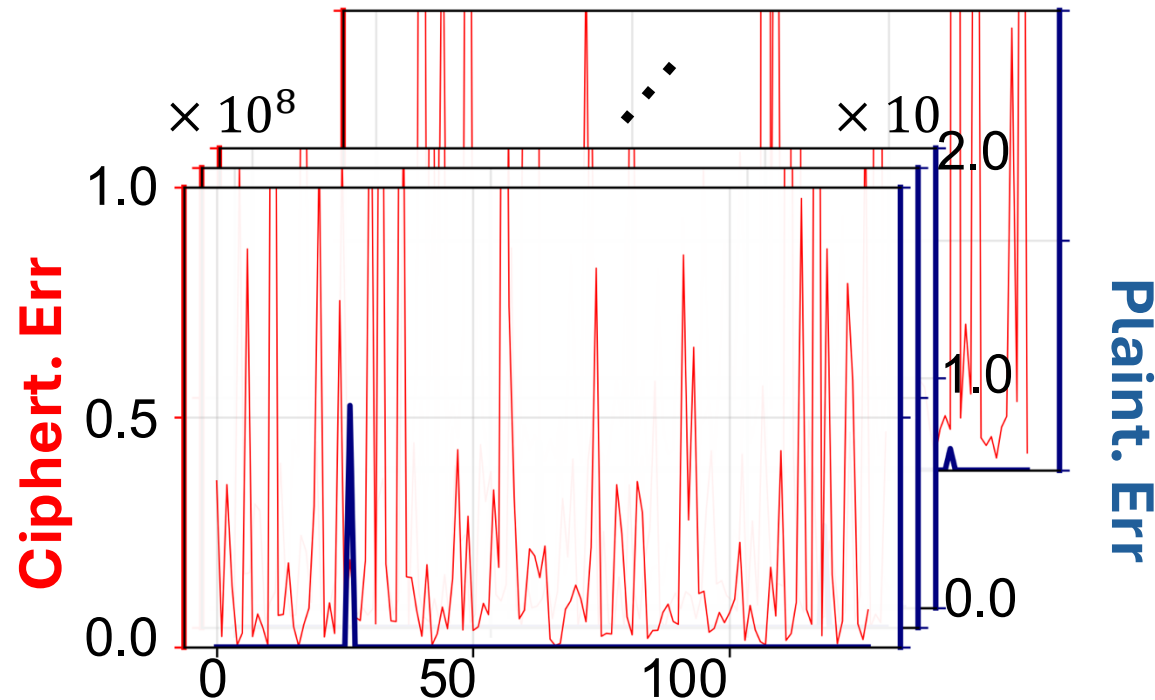
# Primitive Operations-Hadamard Product



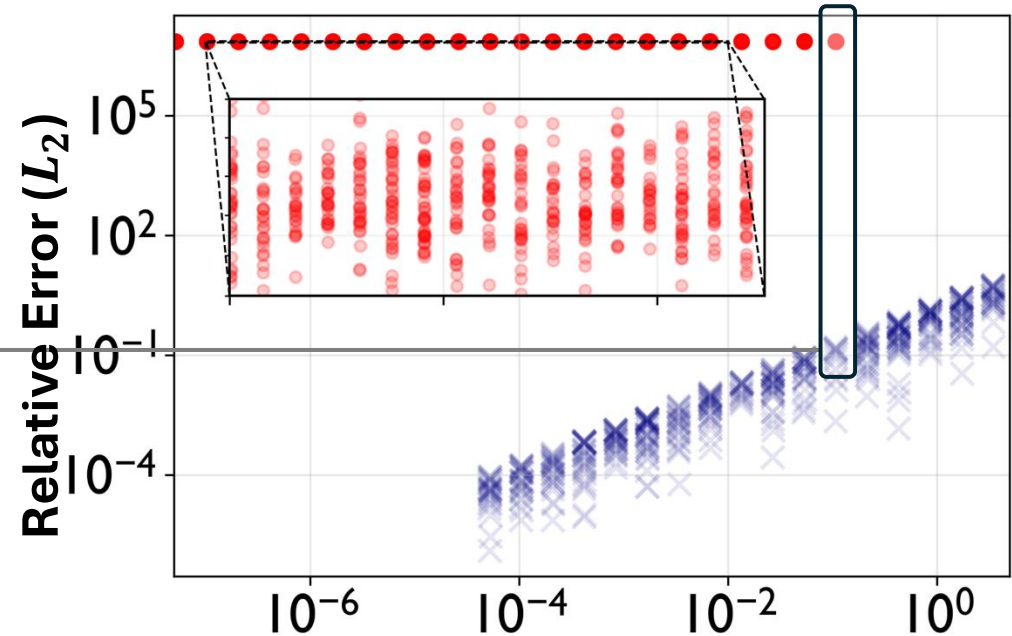
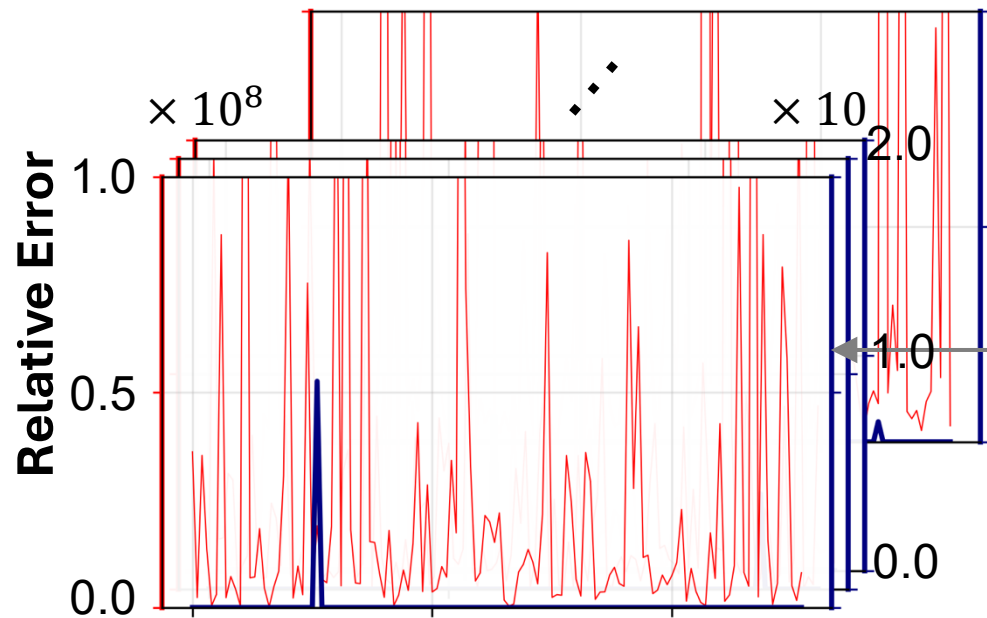
# Primitive Operations-Hadamard Product



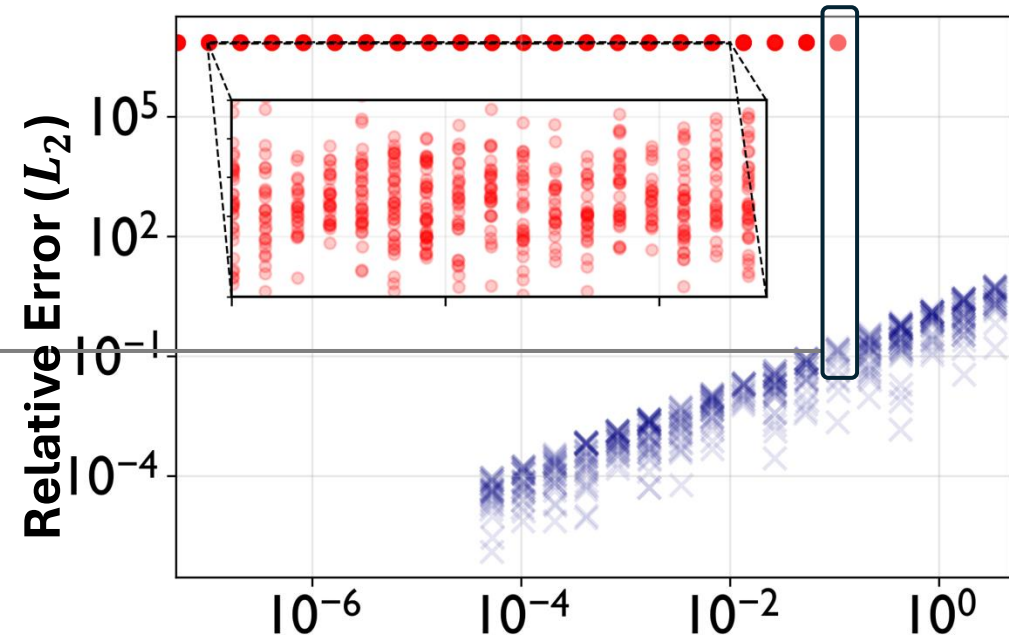
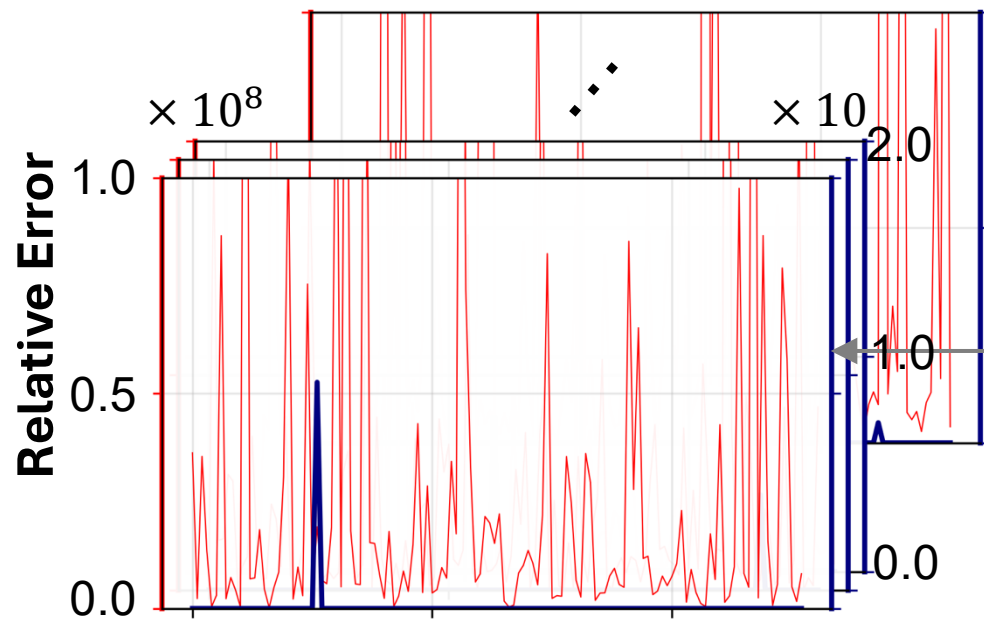
# Primitive Operations-Hadamard Product



# Primitive Operations-Hadamard Product



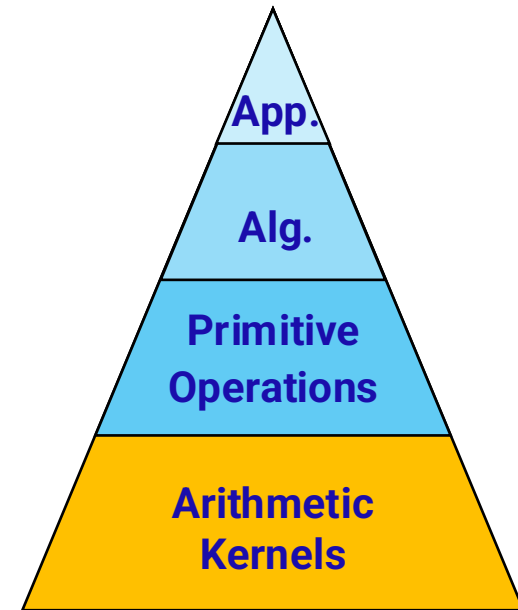
# Primitive Operations-Hadamard Product



**Error been broadcast to all elements**

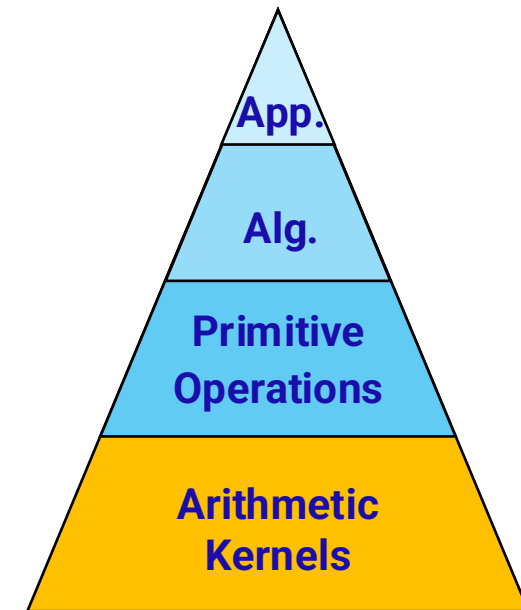
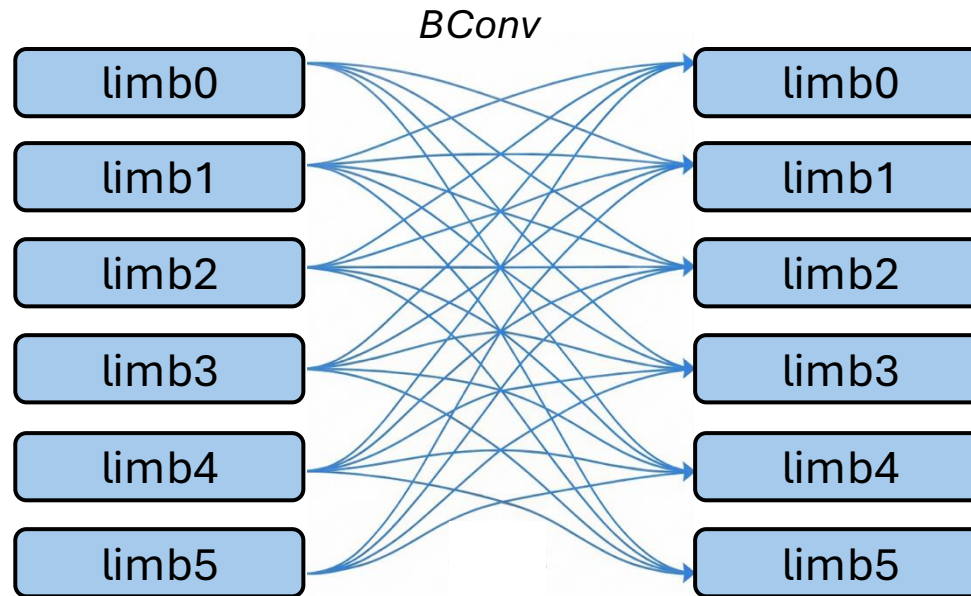
# Arithmetic Kernel-Base Conversion

$$[\mathbf{a}]_{\mathcal{B}} = \text{FastBaseConv}_{\{\mathcal{C} \rightarrow \mathcal{B}\}}([\mathbf{a}]_{\mathcal{C}}) = \sum_{i=0}^L \left( [\mathbf{a}]_{q_i} \cdot \frac{q_i}{Q} \bmod q_i \right) \cdot \frac{Q}{q_i} \bmod p_j$$



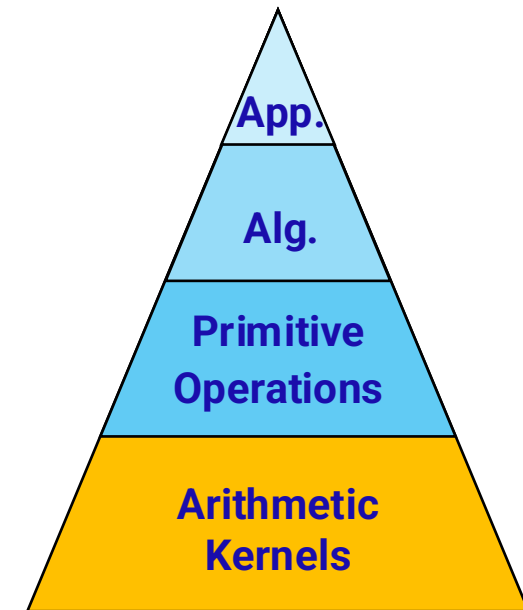
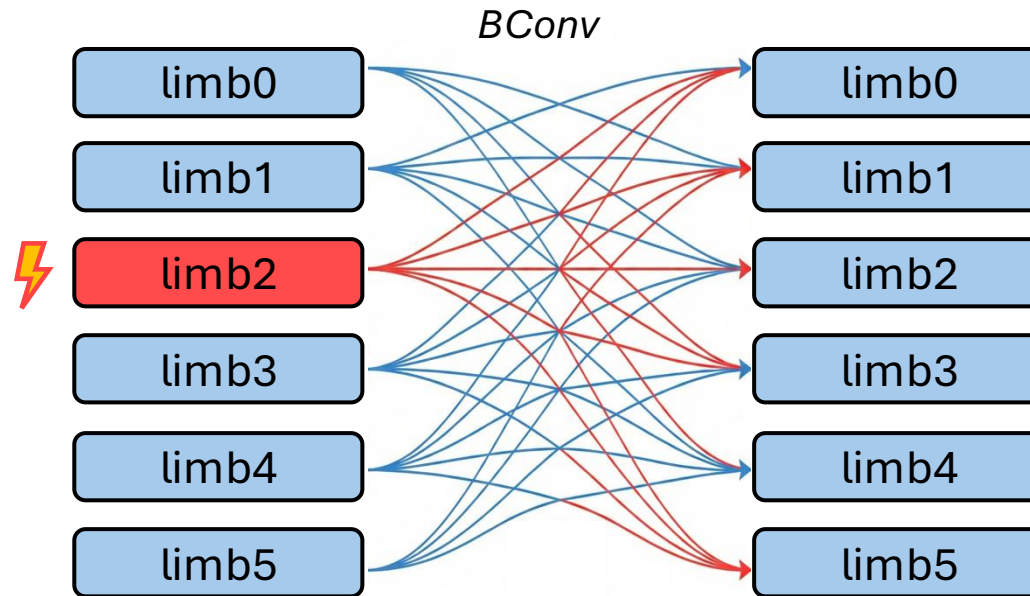
# Arithmetic Kernel-BaseConversion

$$[\mathbf{a}]_{\mathcal{B}} = \text{FastBaseConv}_{\{c \rightarrow \mathcal{B}\}}([\mathbf{a}]_c) = \sum_{i=0}^L \left( [\mathbf{a}]_{q_i} \cdot \frac{q_i}{Q} \bmod q_i \right) \cdot \frac{Q}{q_i} \bmod p_j$$



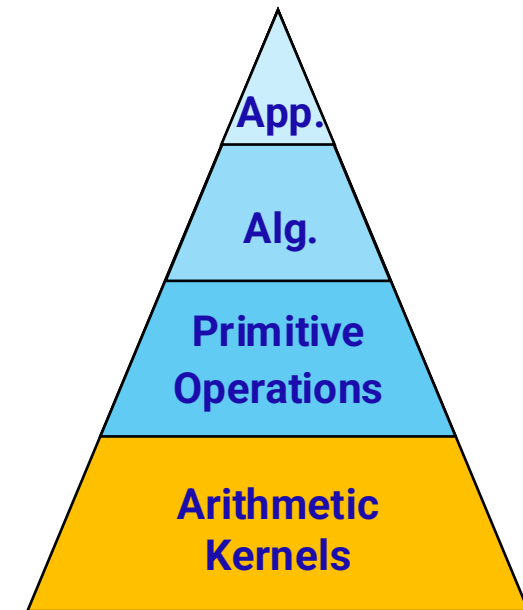
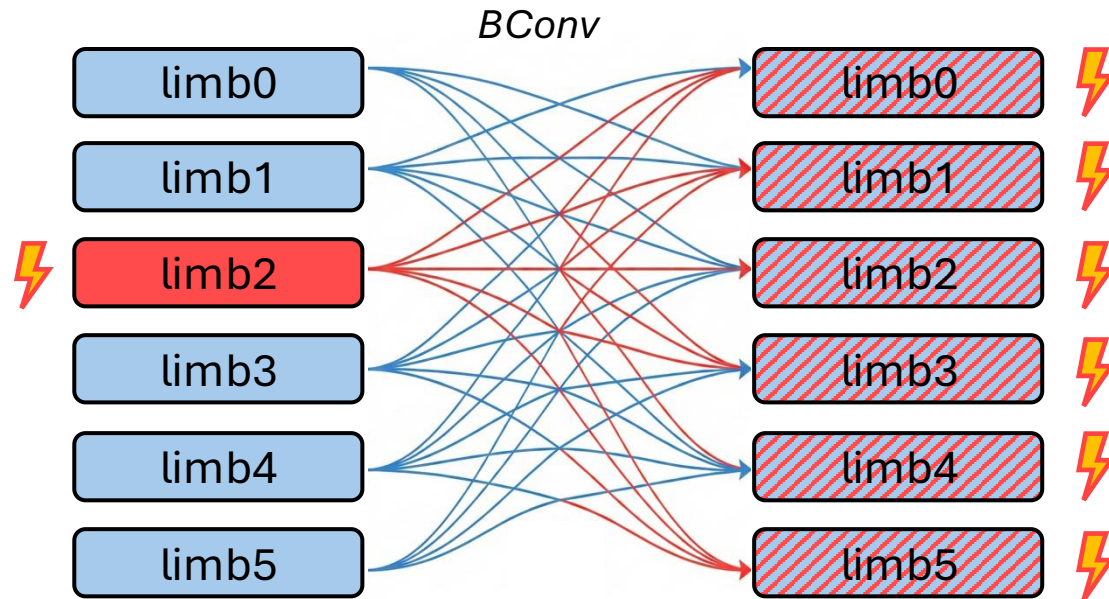
# Arithmetic Kernel-BaseConversion

$$[\mathbf{a}]_B = \text{FastBaseConv}_{\{C \rightarrow B\}}([\mathbf{a}]_C) = \sum_{i=0}^L \left( [\mathbf{a}]_{q_i} \cdot \frac{q_i}{Q} \bmod q_i \right) \cdot \frac{Q}{q_i} \bmod p_j$$



# Arithmetic Kernel-BaseConversion

$$[\mathbf{a}]_{\mathcal{B}} = \text{FastBaseConv}_{\{\mathcal{C} \rightarrow \mathcal{B}\}}([\mathbf{a}]_{\mathcal{C}}) = \sum_{i=0}^L \left( [\mathbf{a}]_{q_i} \cdot \frac{q_i}{Q} \text{mod } q_i \right) \cdot \frac{Q}{q_i} \text{mod } p_j$$



**Error been broadcast to all limbs**

# Elementwise Operation -- Modulus

---

**Algorithm 1** Resilient Barrett Reduction ( $c \equiv t \pmod q$ )

---

**Notation:** modulus  $q$ ;  $t \in \mathbb{Z}_{2^{2K}}^N$

**Precomputed:**  $K = \lceil \log_2 q \rceil$ ,  $\mu = \lfloor 2^{2K} / q \rfloor$

**Input:**  $t, t_{tot}$ ; **Output:**  $c, c_{tot}$ ,

```

1: procedure RELIA-BARRETT-REDUCE( $t, t_{tot}$ )
2:   for each  $t_i$  in  $t$  do
3:      $s_i \leftarrow (t_i \times \mu) \gg 2K$ 
4:      $c_i \leftarrow t_i - s_i \times q$ 
5:     if  $c_i \geq q$  then  $c_i \leftarrow c_i - q$  end if
6:     assert  $0 \leq c_i < q$  //range check
7:   end for
8:    $c_{tot} \leftarrow t_{tot} \pmod q$ 
9:   assert  $c_{tot} == \sum c$  //verify inter-elem sum
10:  return  $c, c_{tot}$ 
11: end procedure

```

---